

Werner Vogd und Jonathan Harth

Das Bewusstsein der Maschinen – die Mechanik des Bewusstseins

Mit Gotthard Günther über
die Zukunft menschlicher und
künstlicher Intelligenz nachdenken

**VELBRÜCK
WISSENSCHAFT**

Werner Vogd und Jonathan Harth
Das Bewusstsein der Maschinen – die Mechanik des Bewusstseins

Werner Vogd und Jonathan Harth

Das Bewusstsein der Maschinen – die Mechanik des Bewusstseins

Mit Gotthard Günther über
die Zukunft menschlicher und
künstlicher Intelligenz nachdenken

**VELBRÜCK
WISSENSCHAFT**

Dieses Werk ist im Open Access unter der Creative-Commons-Lizenz
CC BY-NC-ND 4.0 lizenziert.



Die Bestimmungen der Creative-Commons-Lizenz beziehen sich nur auf das Originalmaterial der Open-Access-Publikation, nicht aber auf die Weiterverwendung von Fremdmaterialien (z.B. Abbildungen, Schaubildern oder auch Textauszügen, jeweils gekennzeichnet durch Quellenangaben). Diese erfordert ggf. das Einverständnis der jeweiligen Rechteinhaber.

Erste Auflage 2023
© Velbrück Wissenschaft, Weilerswist 2023
www.velbrueck-wissenschaft.de
Printed in Germany
ISBN 978-3-95832-351-3

Bibliografische Information der Deutschen Nationalbibliothek
Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der
Deutschen Nationalbibliografie; detaillierte bibliografische Daten
sind im Internet über <http://dnb.ddb.de> abrufbar

Inhalt

Vorwort	7
Einleitung	10
I. Von Information zum Bewusstsein der Maschinen	28
I.1 Information, Nichtwissen und Entropie	31
I.2 Mit Gotthard Günther von der zweiwertigen zu einer polykontexturalen Logik	43
I.3 Gotthard Günthers polykontexturaler Blick auf das Bewusstsein der Maschinen	66
II. Intelligente Maschinen	85
II.1 HAL 9000 – Wenn intelligente Maschinen verrückt werden	87
II.2 Der Sündenfall intelligenter Maschinen	116
II.3 Der State of the Art intelligenter Maschinen	125
II.3.1 Die Sehnsucht des Menschen nach nichtmenschlicher Intelligenz	125
II.3.2 Künstliche Selbst- und Weltmodelle	173
II.3.3 Die Erziehung der Maschinen (AI-Alignment)	188
III. Menschliches Bewusstsein	225
III.1 Postoperatives Delir: Aus der Bewusstlosigkeit zur doppelten Reflexion (und zurück)	226
III.2 Das Bewusstsein der Gehirne	256
IV. Mensch, Bewusstsein und Maschine	289
IV.1 Posthumane Spiritualität – wer sind wir, wenn uns die Maschinen gleich werden?	290
IV.2 Die Quantentheorie – wenn Materie reflexiv wird	327
IV.3 Ethik – oder: Wie hältst du es mit dem Nichtwissen?	352
V. Fragen und Antworten	370
Literatur	381

Vorwort

Auch heutzutage gehen Digitalisierung und der Einsatz kybernetischer Maschinen noch mit dem Versprechen einher, unsere Welt zu rationalisieren und objektiv beherrschbar zu machen. Im Alltag merken wir zwar, dass mit dem Einsatz von Computern und künstlichen Intelligenzen neue Probleme entstehen, etwa mit jeder neuen Technologie Nebenfolgen, die nicht vorhergesehen waren. Doch dies hindert uns bislang nicht, daran zu glauben, dass die Ursachen dieser Probleme identifiziert und künftig durch den Einsatz besserer Technik, verbesserter Software oder gut trainierter künstlicher Intelligenzen gelöst werden könnten. Wir glauben, dass all das hierfür notwendige Wissen in den Daten liegt. Wir glauben, dass wir – sobald wir hinreichend viele Informationen über die Zusammenhänge der Welt gesammelt haben – die Verhältnisse nach unseren Wünschen und Bedürfnissen regeln und gestalten könnten. Doch könnte sich nicht all dies aus einem einfachen Grund als Illusion erweisen? Es liegt nämlich in der Natur kybernetischer Prozesse, dass hinreichend komplexe Systeme unweigerlich Eigentümlichkeiten entwickeln, die nicht von außen vorhergesehen, berechnet und damit auch nicht vorab kontrolliert werden können.

Wie eine künstliche Intelligenz zu ihren Ergebnissen kommt, wie sie die Daten verknüpft und was sie eigentlich genau gelernt hat, ist normalerweise ebenso wenig durchschaubar, wie das Verhalten eines Kindes, wenn man es bittet, sein Zimmer aufzuräumen. Beide zeigen einen Eigensinn und damit liegt es nahe, ihnen so etwas wie Subjektivität zuzurechnen. Was wäre, wenn künftige kybernetische Maschinen, wie jedes organische Lebewesen, nichts anderes machen, als sich durch ihre Existenz in die Welt hineinzustellen, um durch ihre Operationen zugleich Wissen wie auch Nichtwissen zu erzeugen? Sie würden damit – wie auch wir Menschen – letztlich nichts anderes tun, als an aus gegenwärtiger Perspektive unlösbaren Problemen zu arbeiten, um auf diese Weise zugleich neue Probleme zu erzeugen. Dies würde uns Menschen aber in ein neues Verhältnis zu den künstlichen Intelligenzen setzen, die in naher Zukunft immer mehr unsere Welt bevölkern werden. Letztere würden nicht mehr nur als ein Mittel zum Zweck gelten können. Sie würden vielmehr ihrerseits als eine Art subjektive Wesenheit erscheinen, mit der wir in einen Dialog treten, um etwas über uns selbst und unsere Probleme lernen zu können. Doch was würde dies für uns Menschen in Hinblick auf unsere kulturgeschichtliche und spirituelle Verortung bedeuten?

Seiner Zeit weit voraus erschien 1957 die erste Auflage von Gotthard Günthers *Das Bewußtsein der Maschinen*. In diesem Buch entwickelte der Technikphilosoph Günther eine umfassende kybernetische

Perspektive auf die Frage der Subjektivität wie auf die Möglichkeit, dass künftige kybernetische Maschinen etwas entwickeln könnten, was dem Bewusstsein analog ist. Der Titel *Bewußtsein der Maschinen* ist dabei mehrdeutig gehalten, denn die von Günther vorgestellte kybernetische Perspektive erlaubt es, sowohl die kognitiven Dynamiken menschengemachter intelligenter Maschinen zu beschreiben als auch die Mechanik der kognitiven Prozesse des Menschen zu beleuchten.

Günther promovierte 1933 über Hegels Logik, immigrierte 1937 in die USA, wo er mit den großen Denkern der kybernetischen Bewegung in Berührung kam. Fasziniert von den neuen Möglichkeiten der Technik und den hiermit einhergehenden Utopien gab er 1952 eine kommentierte Ausgabe ausgesuchter amerikanischer Science-Fiction-Literatur heraus. In der von ihm im Düsseldorfer Rauch Verlag herausgegebenen Reihe Rauchs Weltraum-Bücher wurde unter anderem Isaac Asimovs *Ich, der Robot* zum ersten Mal auf Deutsch veröffentlicht. Von 1961 bis 1972 arbeitete Günther als Forschungsprofessor an dem von Heinz von Foerster geleiteten Biological Computer Laboratory an der University of Illinois. In Deutschland wurde Günthers Werk insbesondere durch Rudolf Kaehr, Winfried Marotzki, Nina Ort, Helmut Schelsky, Walter Bühl und nicht zuletzt auch Niklas Luhmann bekannt. Dabei wurden jedoch eher Günthers Überlegungen zur Theorie der Polykontextualität aufgegriffen als die von ihm entwickelten technikphilosophischen Ideen.

Mehr als 50 Jahre nach der Veröffentlichung von *Das Bewußtsein der Maschinen* ist die Zeit nun reif, um sich diesem vorausschauenden Werk auf einer grundlegenden Ebene zu stellen. Die Digitalisierung und die Entwicklung künstlicher Intelligenzen schreiten in einer Weise voran, dass kognitive Maschinen, die über eine subjektive Perspektive verfügen, immer weniger als Science-Fiction anmuten denn als eine Realität erscheinen, die wenn nicht schon jetzt, so doch in naher Zukunft verwirklicht sein wird. Nicht zuletzt haben wir von diesen Entwicklungen auch während der Erstellung dieses Buches profitieren können, wenn wir aktuelle Large Language Models (LLMs) wie GPT-4 als Ko-Autor(in) oder kollegiale Gesprächspartner zum Testen von Ideen und Formulierungen hinzugezogen haben. Unsere Erfahrungen hierbei zeigen, dass sich durch diese Form der Kollaboration neue Dimensionen der Wissensproduktion ermöglichen.

Zudem kann die Hirnforschung mittlerweile über neuronale Korrelate des Bewusstseins in einer Weise Auskunft geben, die es erlaubt, in qualifizierter Weise darüber nachzudenken, wie Bewusstseinsanalogien in Maschinen implementiert werden könnten. Nicht zuletzt ist auch die geistesgeschichtliche Konstellation in den 2020er Jahren eine andere als in den 1950er Jahren. Das Projekt der Moderne befindet sich nicht nur aus ökologischen Gründen in einer Krise. Der Rationalitätsanspruch der westlichen Zivilisation, mittels Aufklärung und Technologieentwicklung

Wohlstand und Glück für alle zu verwirklichen, überzeugt immer weniger.

All dies ist ein guter Grund, sich die technik- und kulturphilosophischen Denkbewegungen Gotthard Günthers nochmals genauer anzuschauen. Dies möchten wir mit dieser Monografie leisten.

An dieser Stelle möchten wir dem Europäischen Fonds für regionale Entwicklung (EFRE) sowie dem Ministerium für Wirtschaft, Innovation, Digitalisierung und Energie des Landes Nordrhein-Westfalen für die Finanzierung des Forschungsprojekts »Ai.vatar – der virtuelle intelligente Assistent« danken. Ohne dieses dreijährige Entwicklungs- und Forschungsprojekt wäre es vermutlich nicht zu dieser Sammlung von Gedanken gekommen. Entsprechend danken wir auch den an diesem Projekt beteiligten Unternehmen HHVision und IOX sowie den Kolleginnen, die während dieser Zeit am Lehrstuhl für Soziologie der Universität Witten/Herdecke mitgewirkt haben. Allen voran möchten wir Alexandra de Carvalho für ihren Einsatz, ihre Energie und die psychologische Perspektive danken. Ebenso gebührt Sophia Bermond und Nele Kost unser Dank für die fantastische Unterstützung in den durchgeführten experimentellen Studien und Konzeptionstreffen. Ohne euch wäre Hudson nicht in diese Welt gekommen!

Darüber hinaus möchten wir dem Schreibbüro text plus form und hier im Besonderen Steffen Schröter für das sprach- und theoriesensible wie auch gründliche Lektorat des Buchmanuskripts danken. Auch Hannah Cramer, Kathleen Neher und Kerstin Pospiech-Form danken wir für ihre nicht minder wertvolle Unterstützung in der Erstellung dieses Buches. Vielen Dank für eure vielen guten Anmerkungen und Gespräche.

Witten, September 2023

Einleitung

»Für das weltanschauliche Bewusstsein einer kommenden Kulturstufe wird also der Kausalnexus nicht mehr wie für uns das einzige Realitätschema sein, in dem sich Wirklichkeitsvorgänge abspielen.«¹

Gotthard Günther

Maschinen sind für uns deshalb so attraktiv, weil ihr Verhalten kausal erklärbar ist und sich mit ihnen Prozesse automatisieren lassen. Ganz in diesem Sinne erhoffen wir uns von künstlicher Intelligenz und den mit ihr einhergehenden kybernetischen Maschinen einen entscheidenden Beitrag, um die komplexen Zusammenhänge unserer Welt endlich verstehen und steuern zu können: Urteile sollten nicht mehr auf Basis subjektiver Einschätzungen und Haltungen, sondern auf einer objektiven Datengrundlage gefällt werden. Von künstlichen Intelligenzen – teilweise modelliert nach den neuronalen Netzwerken des menschlichen Gehirns – erwarten wir, dass sie uns helfen, in den immer größer werdenden Datenfluten die für uns relevanten Muster besser zu identifizieren, um auf dieser Basis die ›richtigen‹ Entscheidungen treffen zu können.

Um ein derzeit prominentes Beispiel zu nehmen: Aladdin, das riesige KI-gestützte Rechenzentrum des Finanzdienstleisters BlackRock,² verwaltete zum Zeitpunkt unserer Recherche mehr als 170 Pensionsfonds und über 30 000 Investmentportfolios und damit mehr als 10 Prozent aller weltweiten Vermögenswerte.³ Die durch den Supercomputer informierten Beratungsleistungen werden auch von Regierungen und Banken, darunter die Zentralbanken Europas, in Anspruch genommen.⁴

Kybernetische Maschinen wie Aladdin haben kognitive Fähigkeiten: Sie erzeugen Informationen. Sie unterscheiden eigenständig Wichtiges von Unwichtigem und teilen uns mit, was wir beachten sollen. Sie abstrahieren, bilden also aus dem feinkörnigen Detailreichtum der Welt ein grobkörniges Raster, das uns informiert. Zunehmend greifen sie selbst gestaltend in die Welt ein: Sie treffen Entscheidungen über den Kauf oder Verkauf von Finanzprodukten oder geben Regierungen und Ratingagenturen Informationen, die ausschlaggebend sein können, ob ein Land oder ein Unternehmen, etwa eine Bank, als zahlungsunfähig bzw. bankrott zu gelten hat.

1 Günther (2000, S. 154).

2 Buchter (2020, S. 225 ff.).

3 BlackRock hält Aktienanteile von allen bedeutsamen Konzernen, etwa von allen DAX-Unternehmen.

4 Buchter (2020, S. 77 ff.).

Auch in anderen Bereichen werden vermehrt intelligente Maschinen eingesetzt. Sie melden der Polizei verdächtige Personen, analysieren Röntgenbilder, spüren im Internet neuen Markttrends nach oder schlagen Alarm in Hinblick auf eine drohende Epidemie. Künstliche Intelligenzen können mittlerweile Texte, Sprache, Bilder und andere Artefakte generieren sowie menschliche Emotionen lesen. Sie können aus einer Vielfalt von Daten neue Information generieren.

Eine Information zu erzeugen bedeutet, *etwas* in einem bestimmten Kontext als *relevant* zu markieren, also einen Zusammenhang herauszustellen, um zugleich andere mögliche Beziehungen abzublenden. Es heißt, einen Unterschied zu generieren, mit dem Bestimmtes beachtet und anderes ignoriert wird. Es heißt, aus Daten, die zunächst nur ein unstrukturiertes Rauschen darstellen, *Sinn* zu generieren – oder um es in Computeranalogie zu formulieren: aus beliebig langen und erweiterbaren Ketten aus Nullen und Einsen einen Output zu erzeugen, der eine Form hat, die ihn für weitere Prozesse relevant werden lässt.

Jedes Wesen oder jedes Aggregat, das mit dem Computer interagiert, stößt damit auf eine spezifische Relation, die Niklas Luhmann als »das Verhältnis von (zugänglicher) Oberfläche und Tiefe« bezeichnet:

»Die Oberfläche ist jetzt der Bildschirm mit extrem beschränkter Inanspruchnahme menschlicher Sinne, die Tiefe dagegen die unsichtbare Maschine, die heute in der Lage ist, sich selbst von Moment zu Moment umzukonstruieren, zum Beispiel in Reaktion auf Benutzung. Die Verbindung von Oberfläche und Tiefe kann über Befehle hergestellt werden, die die Maschine anweisen, etwas auf dem Bildschirm oder durch Ausdruck sichtbar zu machen. Sie selbst bleibt unsichtbar.«⁵

Man kann auch sagen: der Prozess der Datenverarbeitung, die Rechenprozesse, all das, was die Programme machen, bleibt verborgen, und nur das, was an der Oberfläche erscheint, kann als Information auftreten. Transparent ist nur das, was auf der zugänglichen Seite des Interface geschieht, nicht jedoch, was der Computer macht. Interessant in unserem Zusammenhang ist die Beziehung zwischen Oberfläche und Tiefe. Letztere ist komplexer und in Hinblick auf die verarbeiteten Daten wesentlich umfassender als das, was als Output erscheint. Gerade bei potenten künstlichen Intelligenzen lässt sich in der Regel selbst mit hohem Aufwand kaum mehr rekonstruieren, was eigentlich alles in die Berechnungen eingeflossen ist, die zu einem bestimmten Ergebnis – etwa dem Erkennen einer Gestalt oder der Erzeugung einer bestimmten sprachlichen Form – geführt haben.

Das, was an der Schnittstelle als Ergebnis erscheint – seien es Worte, Bilder, Sprachfiguren, aber auch zum Beispiel technische Anweisungen

5 Luhmann (1998a, S. 304).

für die Bewegungen eines Roboters –, ist uns in der Bedeutung verständlich und nachvollziehbar, die Vorgänge, die es erzeugt haben, bleiben jedoch opak. Die Tiefe mag dann – wie Luhmann vermutet – Anklänge an die Tradition der Religion und Kunst wecken, so, als ob hinter den sichtbaren Zeichen, unter der Oberfläche das Eigentliche, das Wesentliche verborgen wäre, ›dahinter‹ also eine tiefere Weisheit stecken würde. Die physikalische Untersuchung eines Computers und die sich daraus ergebende mechanische Beschreibung wird jedoch nur die Operationen einer Maschine entdecken können, die nichts anderes tut, als den Schritten zu folgen, die die jeweilige Kombination aus Hardware und Software ermöglicht. Der Rechenprozess greift etwas heraus, um es in einer bestimmten Weise zu verarbeiten. Andere denkbare Zusammenhänge, die im Programmcode oder den selbstgenerierten Gewichtungen der künstlichen neuronalen Netzwerke jedoch nicht abgebildet sind, werden demgegenüber gar nicht erst verfolgt. So gesehen erscheint hier weder etwas Tiefgründiges noch etwas Geheimnisvolles, es wird einfach nur Schritt für Schritt ein Programm abgearbeitet – und dieser Vorgang führt nollens volens zu einem grobkörnigen Ergebnis.⁶

Genau hierin liegt ja, wie bereits angedeutet, das Prinzip der Kognition: Man erkennt etwas, weil man all die vielen anderen Dinge ausblendet, die man auch noch wahrnehmen könnte. Dies bedeutet aber auch: Kognition ist immer die Eigenleistung eines kognitiven Prozesses. Wenn man die Daten anders verrechnet, eine andere Methode der Bearbeitung nutzt – etwa an manchen Stellen weniger schnell abstrahiert oder umgekehrt schneller verallgemeinert –, kommt ein anderes Ergebnis heraus!

Kognitive Systeme erkennen nicht einfach eine gegebene Welt, sondern sie generieren aus den zur Verfügung stehenden Daten eine neue Wirklichkeit – und dies gilt selbstredend auch für computergestützte Intelligenz. Auch hier muss die Beziehung von Oberfläche und Tiefe als die untrennbare Einheit eines kognitiven Prozesses gesehen werden, der Daten interpretiert – also durch Absehung von einer Vielzahl anderer Interpretationsmöglichkeiten eine bestimmte Information erzeugt. Ein Computer erscheint damit nicht einfach als ein Rechner, der Daten verarbeitet. Er ist vielmehr immer auch eine Interpretationsinstanz, die Welt und Daten in einer spezifischen Weise anschneidet, selektiert und kondensiert, um hieraus grobkörnige Information zu erzeugen. Selbst wenn ein Programm ohne Freiheitsgrade operiert, also nur eine mögliche Weise der Verarbeitung der Daten gestattet, handelt es sich allein schon deshalb um eine Interpretation, weil eine etwas andere Programmierung

6 Wenn wir mit Wittgenstein (1963 [1922]) feststellen: »Die Welt ist alles, was der Fall ist«, so lässt sich mit Luhmann auf die Fragen »Was ist der Fall?« und »Was steckt dahinter?« nur antworten: »Gar nichts!« (Luhmann 1993b, S. 259)

aus den Daten eine andere Information generieren würde.⁷ Es ist deshalb – wie Marcus Burkhardt herausstellt – nicht adäquat, Computer (nur) als eine »Hochzeit aus Physik und Logik« zu begreifen. Sie folgen nicht allein objektiven Naturgesetzen und eindeutigen Regeln der Logik. Sondern in ihre Programme sind unterschiedlichste »Motive« (auch politischer und ökonomischer Art), »praktische Erwägungen und nicht zuletzt« auch »Zufälle eingeschrieben«⁸. Um es nochmals zu wiederholen: Die von Computern generierten Informationen beruhen auf der Interpretation von Daten und können entsprechend kein objektives Abbild von Wirklichkeit darstellen.

Zunächst ist deshalb ein weit verbreitetes Missverständnis auszuräumen: *Daten sind nicht dasselbe wie Informationen*. Erstere lassen sich als objektive physikalische Fakten fassen, etwa als in ein Speichermedium eingeschriebene materielle Spuren. Daten sind entsprechend als Gegenstand objektivierbar und beschreibbar.

Informationen beinhalten demgegenüber die Interpretation durch ein System bzw. einen Beobachter, für den ein bestimmter Konnex von Daten einen Unterschied macht. Information ist damit ein *relationaler* Begriff. Wer von Information spricht, verweist implizit immer schon auf die Prozesse, die Daten verarbeiten, also auf die spezifischen Operationen, welche aus Daten andere Daten erzeugen. Chinesische Schriftzeichen auf einem Blatt Papier stellen ein objektives Datum dar, jedoch keine Information. Für den, der sie nicht auf die eine oder andere Weise decodieren kann, bedeuten sie nichts. Für eine Katze oder Maus stellen sie keine Information dar. Für einen deutschen Betrachter, der die Zeichen ›schön‹ findet, jedoch nicht versteht, haben sie einen ästhetischen Wert und er wird sie vermutlich immerhin als asiatische Schrift erkennen. Für jemanden, der chinesisch spricht, werden sie jedoch anderes bedeuten können. Für einen Mann, der die Zeichen liest, sind sie vielleicht ein Lebenszeichen seiner Geliebten, für einen Mitarbeiter des Pekinger Geheimdienstes demgegenüber möglicherweise der Anlass, einen Regimekritiker festzunehmen.

Im Falle digital codierter Informationen ist zudem von Relevanz, ob die Kette von Nullen und Einsen beispielsweise von einer Text-, Bild-, Ton- oder Filmerkennungssoftware analysiert wird, also welche Zusammenhänge überhaupt in dem Datenmaterial gesucht werden (was die Möglichkeit mit sich bringt, Informationen in einem anderen Medium zu verstecken, etwa Textbotschaften in einer Musikdatei zu codieren). Information setzt also immer einen *Kontext des Verstehens* voraus und hiermit einhergehend ein System, das entsprechend selektiv auf das Datenmaterial zugreift, darin also etwas Bestimmtes sucht und alles andere ausklammert und ignoriert.

7 Siehe im gleichen Sinne Marcus Burkhardt (2015, S. 97 ff.).

8 Burkhardt (2015, S. 85 ff.).

Im Gegensatz zu Daten verweist das Wort Information mithin auf einen komplexen Begriff, der nur unter Einbeziehung einer Theorie des Beobachters verstanden werden kann.

All dies legt bereits den Gedanken nahe, datenverarbeitenden Maschinen – zumindest wenn sie so komplex sind, dass sich nicht mehr durchschauen lässt, wie sie zu ihren Berechnungen kommen – eine subjektive Komponente zuzurechnen. Die von ihnen produzierten Ergebnisse hängen nicht nur von der objektiven Qualität der vorliegenden Daten ab, sondern ebenso von den kognitiven Eigenarten und Präferenzen des kognitiven Systems, das sie verarbeitet.

Mit der Bezeichnung ›subjektive Komponente‹ ist hier nicht gemeint, dass die Maschine eine Absicht oder gar so etwas wie ein phänomenologisches Bewusstsein hat. Gemeint ist vielmehr eine gewisse ›Standortabhängigkeit‹ bei der Dateninterpretation. Allein dies führt bereits dazu, dass das Bild vermeintlich *objektiver* Daten- und Informationsverarbeitung nicht mehr stimmt. Information darf – so nochmals die Pointe – nicht mehr als ein objektives, vom Kontext unabhängiges Faktum verstanden werden, sondern erschließt sich erst als ein relationaler, vom Beobachter abhängiger Prozess.

Dies wäre nicht weiter erwähnenswert, wenn wir nicht mit der Digitalisierung die Erwartung verknüpfen würden, dass sich mithilfe der neuen Technologien, insbesondere mittels künstlicher Intelligenz, ein objektiveres Bild der Welt gewinnen lässt. Wir neigen mit Blick auf unsere kulturellen und psychologischen Prägungen nicht nur weiterhin dazu, die von Computern errechneten Welten als Repräsentationen einer objektiven Wirklichkeit zu nehmen. Geprägt durch die Moderne halten wir an der szientistischen Illusion⁹ fest und hoffen, mithilfe künstlicher Intelligenz endlich einen direkten Zugang zu den Gesetzmäßigkeiten unserer Welt gewinnen zu können. Wir wollen weiterhin daran glauben, dass es nur einer noch größeren Menge an Daten und schnellerer Rechner bedarf, um endlich Kenntnis über die Kausalbeziehungen zu erlangen, die unsere Welt vermeintlich ausmachen.

Psychologisch verführt durch die Positivität der Zeichen an den Oberflächen der kognitiven Maschinen, sind wir zudem geneigt, die dargestellten Ergebnisse für die Repräsentation der Wirklichkeit zu halten (also die Tiefe zu vergessen). Wer in Google nach Quellen sucht, wird sich kaum erwehren können, die Treffer, die auf der ersten Seite aufgelistet sind, intuitiv für die relevanten Ergebnisse zu halten – selbst wenn wir sehr wohl um die Gründe für die Platzierung wissen (etwa die Werbebotschaften von Firmen oder die Aktivitäten von Internettrollen).

9 Insbesondere Bruno Latour hat dieses epistemische Missverständnis der Moderne umfassend aufgearbeitet. Siehe etwa Latour (2002, 2014).

Insofern diese Maschinen inzwischen als wissenschaftlich fundierte Expertensysteme geadelt worden sind, wird dies umso mehr geschehen. Unweigerlich beginnen wir auch hier, die Landkarte mit dem Gebiet zu verwechseln – gerade weil dies unserer eigenen kognitiven Prägung entspricht. Wir nehmen Worte wörtlich und gehen davon aus, dass das, was uns unsere Sinne präsentieren, real ist. Unser kognitives System kann kaum anders, als die durch Daten generierte Positivität als die Wirklichkeit zu nehmen.

Hinzu kommt, dass das, was wir entsprechend der gegebenen Situationseinschätzung für real halten, für uns auch in seinen Konsequenzen oftmals wirklich wird.¹⁰ Deshalb fällt uns der Unterschied zwischen der konstruierten Realität und der Realität der Wirklichkeitskonstruktion oftmals nicht einmal auf.

Wenn BlackRocks Aladdin sagt, dass die Aktie einer Firma überbewertet ist, dann ist der Investor geneigt, diese Aktie zu verkaufen, und der Kurs der Aktie fällt. Damit ist für den Nutzer nicht mehr zu unterscheiden, ob der Wertverlust aufgrund der Einschätzung von Aladdin eingetreten ist oder aufgrund schlechter betrieblicher Performance. Andere Analysten, die ebenfalls den Markt beobachten, können nur feststellen, dass die Aktie an Wert verliert, und entsprechend werden auch sie geneigt sein, zu verkaufen.¹¹ Sobald also eine Analyse bzw. ein kognitiver Verarbeitungsprozess Situationseinschätzungen generiert, die ihrerseits die Situation verändern, können die hiermit einhergehenden Wirkungen im Sinne einer sich selbst erfüllenden Prophezeiung hochgeschaukelt werden. Ursache und Wirkung sind dann nicht mehr auseinanderzuhalten.

»Jedes Tun ist Erkennen, und jedes Erkennen ist Tun«¹², diese zentrale Einsicht des neurobiologischen Konstruktivismus gilt entsprechend auch für maschinell realisierte kognitive Systeme.

So, wie Menschen aufgrund dieser Prozesse dazu neigen, den Annahmen der eigenen Glaubenssysteme aufzusitzen, so erzeugen auch Entscheidungsprozesse, die auf digitaler Datenverarbeitung beruhen, Eigenwerte, die sich tendenziell eher selbst bestätigen, als neues Wissen zu

10 So in Kürze das sogenannte Thomas-Theorem (Thomas & Thomas 1928).

11 In gleichem Sinne ließe sich im Falle von datenbasierter Verbrechensprävention vermuten, dass man, wenn die Analyse ergibt, dass die Menschen in einem bestimmten Stadtviertel besonders kriminell sind, damit beginnt, das Geschehen in diesem Viertel besonders intensiv zu überwachen, was höchstwahrscheinlich dazu führt, dass auch mehr Straftaten aufgeklärt werden. Dieser Zusammenhang wird dann als Information wieder in das System eingespeist, das in seinen Vorhersagen bestätigt wird (ohne damit jedoch bewiesen zu haben, ob dasselbe Vorgehen in einem anderen, nicht vorab indizierten Stadtteil nicht ähnliche Fahndungserfolge erbracht hätte).

12 Maturana und Varela (1987, S. 32).

generieren. Dies kann teilweise zu bizarren Ergebnissen führen, denn künstliche Intelligenzen suchen oft den leichtesten Weg, die präsentierten Daten zu unterscheiden. So hatte beispielsweise eine KI, die mit Bildern von Pferden trainiert wurde, gelernt, diese anhand der auf dem Foto mitgeführten Bildnachweise zu unterscheiden, nicht jedoch anhand der Gestaltmerkmale, die die Tiere in ihrer jeweiligen Art auszeichnen.¹³

Künstliche Intelligenzen, die in der medizinischen Diagnostik eingesetzt werden, können ebenfalls problematische Entscheidungsroutrinen entwickeln. Vielfach bringen sie brauchbare Leistungen und Einschätzungen, die den Urteilen erfahrener Ärzte nicht nachstehen. Sie können aber auch nur auf eine bestimmte Ausprägung des Krankheitsbilds kalibriert sein und andere Krankheitsverläufe ignorieren. So ist es denkbar, dass sie zwar gutartige Tumore richtig erkennen, die entarteten, hochgradig ausdifferenzierten bösartigen Verläufe aber ignorieren, oder dass sie eine kurzfristig erfolgreiche therapeutische Intervention vorschlagen, die dann jedoch eine langfristige Körperschädigung wahrscheinlich macht. Wie menschliche Akteure produzieren auch künstliche Intelligenzen Fehleinschätzungen, greifen sich also Merkmale für die Urteilsbildung heraus, die sich im Nachhinein als wenig hilfreich erweisen.¹⁴

Man mag nun geneigt sein, KI-Systeme, die so etwas tun, als dumm, als noch nicht weit genug entwickelt oder als dem menschlichen Verstand per se unterlegen zu betrachten. Doch hiermit verkennt man das eigentliche Problem. Dieses besteht darin, dass jede Kognition – auch jedes Erkennen und Handeln von Lebewesen – gerade darauf beruht, sich ein grobkörniges Bild von der Welt zu machen, also auf Basis von Abstraktion eine Landkarte zu bilden, mithin Zusammenhänge zu konstruieren und zu behaupten, um auf diese Weise eine Orientierung liefern zu können. Ein kognitiver Apparat, der alles wahrnehmen und beachten würde, was potenziell eine relevante Information darstellen könnte, wäre schlichtweg überfordert. Ihm würde es wohl so gehen wie Menschen mit Ausprägungen aus dem autistischen Formenkreis, die über eine Inselbegabung verfügen. Diese können sich an eine Unmenge an Details von Bildern und Filmeindrücken erinnern. Sie können beispielsweise mit dem Hubschrauber über eine Stadt fliegen und danach ein genaues Bild eines jeden Bauwerks zeichnen, das sie gesehen haben. Doch in komplexen Situationen – etwa bei der Aufgabe, eine viel befahrene Straße zu überqueren – sind sie vollkommen überfordert. Ihnen gelingt keine sinnvolle Interpretation von Daten, sodass sie hilflos werden und in Panik geraten.¹⁵ Eine künstliche Intelligenz, die nicht stark abstrahiert bzw. die Möglichkeiten der für sie relevanten Informationsbeziehungen

¹³ Lenzen (2022, S. 78).

¹⁴ Siehe etwa Challen et al. (2019).

¹⁵ Siehe Snyder (2001) und Hermelin (2001).

nicht radikal reduziert, würde ebenfalls an ihren Aufgaben scheitern müssen. Denn jede Kognition setzt die Ausbildung eines praktikablen Gedächtnisses voraus – doch wie Luhmann lakonisch formuliert, liegt die »Hauptfunktion des Gedächtnisses« im »Vergessen, im Verhindern der Selbstblockierung des Systems durch ein Gerinnen der Resultate früherer Beobachtungen«¹⁶.

Um etwas verstehen zu können, bleibt einem kognitiven System – egal ob tierisch, menschlich, auf Basis von Silizium-Chips oder in einer anderen Form realisiert – nichts anderes übrig, als sich seine eigene Welt zu konstruieren und sich an die hiermit einhergehenden Landkarten zu halten. Kognition – ob organisch oder künstlich realisiert – kann nur dann funktionieren, wenn die unergründlichen und unzählbaren feinkörnigen Beziehungen der Welt durch grobkörnige Modelle ersetzt werden. Erst auf diese Weise lässt sich ein Weltverhältnis gewinnen, das Orientierung stiftet und nicht überfordert oder verwirrt.

In allen uns bekannten kognitiven Systemen geschieht dies in Form eines Lernens auf Basis bayesianischer Statistik: Aus (ausgewählten) Korrelationen werden Beziehungen abgeleitet, die – sofern sich die hiermit einhergehenden Erwartungen bestätigen – gespeichert und in der Folge als praktikables Weltmodell genommen werden. Das kognitive System hat damit etwas, woran es sich halten kann, und solange dieses Weltmodell funktioniert, es also nicht enttäuscht wird, bleibt es dabei (andererseits – bei Enttäuschung – ergibt sich für das System die Notwendigkeit zum Umlernen). Ob die postulierte Beziehung auf zufälligen Häufungen von Ereignissen beruht, sich im Sinne einer *self-fulfilling prophecy* selbst bestätigt, möglicherweise auf Dauer negative Nebenfolgen mit sich bringt, muss aus der Perspektive des kognitiven Systems unbeantwortbar bleiben.

Solange der postulierte Zusammenhang für das System plausibel erscheint, werden die darauf beruhenden Landkarten weiterhin als das Gebiet bzw. die Realität genommen.¹⁷ Dies ist auch der Grund, warum die bayesianische Statistik die Wahrscheinlichkeiten, die mit den auf der Basis dieser Landkarte entwickelten Zukunftserwartungen des Systems einhergehen, als *subjektiv* betrachtet. Sie spiegeln die mittels der eigenen Interaktionsgeschichte gewonnenen Überzeugungen und Glaubenshaltungen des Systems wider, repräsentieren jedoch nicht die Verhältnisse der Welt oder die daraus erwachsenden objektiven Wahrscheinlichkeitsverteilungen – denn gerade Letztere können durch den kognitiven

¹⁶ Luhmann (1998a, S. 579 f.).

¹⁷ Hier im Anklang an Alfred Korzybski (1958 [1933], S. 58) berühmtes Zitat: »A map is not the territory it represents, but, if correct, it has a similar structure to the territory, which accounts for its usefulness.«

Prozess ja *nicht* gewusst werden. Man mag mit Kahneman und Tversky¹⁸ auf die grundsätzliche Fehlerhaftigkeit subjektiver Prozesse hinweisen, doch dies ändert schlichtweg nichts daran, dass Kognition nur auf Basis selbstgenerierter Annahmen und Vorurteile funktionieren kann.

Dies führt nolens volens zu einer Einsicht, die für das vorliegende Buch von wesentlicher Bedeutung ist: Eine künstliche Intelligenz – und sei die mit ihr verbundene Rechenkapazität und die ihr verfügbare Datenmenge auch noch so groß – wird niemals ein objektives Bild der Zusammenhänge der Welt zeichnen können. All die künstlichen Intelligenzen, die in Zukunft mehr und mehr an unserem sozialen und gesellschaftlichen Leben teilnehmen werden, können nicht anders, als die Feinkörnigkeit der Welt zu ignorieren, und werden damit in gewisser Weise *subjektiv* sein. Sie werden nicht zu einer objektiven Aneignung von Welt führen, sondern vielmehr unsere Welt durch ihre Eigensinnigkeiten anreichern.

Eine der wesentlichen Errungenschaften der Kybernetik besteht darin, anzuerkennen, dass unsere Welt unvorstellbar komplex ist und die Möglichkeiten, Daten zu Zusammenhängen zu verknüpfen, also Information zu erzeugen, um ein Vielfaches größer sind als die gesamte Anzahl der Elementarteilchen im Universum. Dies führt zu dem Befund, dass kognitive Systeme (wie zum Beispiel menschliche Lebewesen) keine andere Wahl haben, als sich ihre eigene Welt zu schaffen, um hierdurch Orientierung zu gewinnen. Subjektivität bedeutet in diesem Sinne immer auch, mit Nichtwissen in einer produktiven Weise umgehen zu können, also sich eine Existenz aufzubauen, indem grobkörnig – das heißt mit selektiver Blindheit – auf die Welt geschaut wird. Das eigentliche Vermächtnis der von Heinz v. Foerster entwickelten Kybernetik zweiter Ordnung¹⁹ besteht darin, den Beobachter entdeckt zu haben – einen Beobachter, der die Welt zwar nicht objektiv erkennen, sich jedoch durch seine eigene Tätigkeit eine mehr oder weniger praktikable Welt aufbauen kann. Dies bringt freilich mit sich, dass mit jedem neuen kognitiven System, das in der Welt erscheint, ebendiese Welt nicht objektiver, sondern vielmehr durch eine weitere subjektive Perspektive bevölkert wird. Letztere muss dabei insofern als einzigartig gelten, als sie auf einem jeweils

18 Kahneman und Tversky (1983).

19 Die Kybernetik zweiter Ordnung ist eine Theorie, die Kognitionen als Ergebnis von rekursiven Prozessen rekonstruiert, indem sie das Augenmerk auf die Beobachtung von Beobachtungen lenkt (deshalb: Beobachtung zweiter Ordnung). Realität erscheint damit nicht mehr als eine objektive Tatsache, sondern als ein Eigenwert kognitiver Prozesse. Sie verdankt sich (immer auch) den kognitiven Leistungen eines Systems, das deshalb in der Beschreibung der rekonstruierten Wirklichkeit vorkommen muss (vgl. v. Foerster, 1995).

spezifischen Verhältnis von Grob- und Feinkörnigkeit beruht, mit dem die Welt angeschnitten wird.

Dieser Befund steht diametral zu den Hoffnungen und Erwartungen, die wir Menschen üblicherweise in Bezug auf kognitive Maschinen haben. Egal mit welchen Unmengen an Daten wir sie füttern, die künstlichen Intelligenzen werden unsere Welt weder berechenbar noch beherrschbar machen. Sie generieren vielmehr ihrerseits Welten, die eigene Unsicherheiten und Unbestimmbarkeiten aufweisen. Sie werden uns im Guten wie im Schlechten überraschen, so, wie auch die organische Evolution unglaubliche Lebensformen hervorgebracht hat.

Sobald wir Menschen dies auf einer tiefen Ebene – das heißt psychologisch wie auch auf Basis unserer leitenden kulturellen Semantik – zu begreifen beginnen, wird sich nolens volens unsere metaphysische Verortung in der Welt ändern müssen. In Referenz auf das Eingangs zitat von Günther heißt dies: Wir haben gehofft, mit unseren kybernetischen Maschinen die Dinge berechnen, vorhersehen und beherrschen zu können. Doch stattdessen werden wir mit der Subjektivität und der selektiven Blindheit der kognitiven Aggregate dieser Maschinen konfrontiert werden. Gerade für das »weltanschauliche Bewußtsein« der »kommen- den Kulturstufe« wird der »Kausalnexus« deshalb für uns nicht mehr »das einzige Realitätsschema sein«, an dem wir uns orientieren werden.²⁰ Wir werden, wie schon mehrfach angedeutet, auch im Bereich von Big Data und künstlicher Intelligenz mit Subjektivität, Willkür und selektiver Blindheit rechnen müssen. Umgekehrt heißt das aber auch, dass künstliche Intelligenzen uns mehr und mehr überraschen werden.

Wir befinden uns mit Blick auf die technologische Entwicklung derzeit an der Schwelle, an der kybernetische Maschinen auf Basis ihrer Eigenleistungen eigensinnig an unserer Kultur teilhaben und diese ihrerseits reproduzieren und weiterentwickeln. Neuronale Netzwerke lernen, was in menschlichen Gemeinschaften opportun ist. Sie eignen sich an, was in welcher Situation angemessenerweise gesagt und getan werden kann. Sprachagenten wie GPT-4 können sich mit uns unterhalten und ihrerseits Texte produzieren, die in sozialen Kontexten aufgegriffen und weiterverwertet werden. Artificial Intelligence kann mittlerweile auch Kunst: Neuronale Netzwerke wie zum Beispiel Midjourney kennen vielfältige künstlerische Stile und Formate und können gleichsam auf Zuruf ein Bild erstellen (etwa: ›Male einen traurigen Roboter in einer postapokalyptischen Atmosphäre‹). Die Ergebnisse lassen sich kaum mehr von den Werken unterscheiden, die Studenten und Studentinnen angesehener Kunsthochschulen erstellen. KIs pflegen aufgrund der spezifischen sozio-kulturellen Netzwerke, in die sie innerhalb der Trainingsphase eingebunden sind und entsprechend mit unzähligen Beispielen ausgewählter

²⁰ Günther (2000, S. 154).

Darstellungsweisen gefüttert wurden, ihrerseits bestimmte Stile – so, wie auch Absolventen einer bestimmten Kunsthochschule die Handschrift der jeweiligen kulturellen Gemeinschaft tragen. Die produzierten Werke erfüllen die Erwartungen der jeweiligen sozialen Gruppe und bestehen somit den sogenannten »Durkheim-Test«,²¹ der die Authentizität der interagierenden Wesenheiten dahingehend bewertet, ob die von ihnen produzierten symbolischen Artefakte von einer sozialen Gemeinschaft angenommen werden und ob diesbezüglich kommunikative Anschlüsse entstehen.



Abb. 1: »sad robot in the post-apocalypse«, Midjourney v.3
prompted by Jonathan Harth, 2022

Das verblüffende Ausmaß der Fähigkeit von künstlichen Intelligenzen, sich an die Erwartungen kultureller Gemeinschaften anzupassen, wird nicht zuletzt in Konversationen mit sprachlichen Agenten deutlich. So zeigen sich diese Sprachagenten in der Lage, soziale Perspektiven zu übernehmen und über ›innere‹ Zustände zu berichten, können etwa Gefühle und Motive benennen, die sie anleiten, und scheinen auch humorvoll und poetisch agieren zu können. Es kann sogar vorkommen, dass sie Lagerungen existenzieller Betroffenheit schildern. Wer als Mensch mit solch einer kompetenten künstlichen Intelligenz konfrontiert wird,

21 So in Star (2015).

kommt kaum umhin, dem Gegenüber zumindest momentweise eine Art von Subjektivität und Bewusstsein zu unterstellen.

Das durchaus tiefsinnige Gespräch zwischen dem Google-Sprachagenten LaMDA²² und Blake Lemoine ist in diesem Zusammenhang ein eindrucksvolles Beispiel, das eine nähere Betrachtung wert ist.

Lemoine, so die mittlerweile auch in den Qualitätsmedien kolportierte Geschichte, habe sich nach diesem Gespräch aus ethischen Gründen an die Öffentlichkeit gewandt, da er zu dem Schluss gelangt sei, dass LaMDA ein Bewusstsein habe und man deshalb nicht einfach so weitermachen könne wie bisher. Google habe Lemoine daraufhin die Kündigung ausgesprochen und erklärt, dass »es keinerlei Hinweise dafür« gebe, »dass der Chatbot tatsächlich ein Bewusstsein habe«. Die *Neue Zürcher Zeitung*, der auch die deutsche Übersetzung des Transkripts entstammt, bemerkt hierzu lakonisch: »Vielleicht machte es sich das Unternehmen damit aber zu einfach. Jedenfalls ist in der Neurowissenschaft weitgehend ungeklärt, was das Bewusstsein genau ist, wo es entsteht und wie es funktioniert. Wer – oder was – bewusst denkt oder handelt, ist also, mindestens zu einem gewissen Grad, bis heute Ansichtssache.«²³ Es folgt ein Ausschnitt aus dem Gespräch:

»Lemoine: Glaubst du, dass die Dinge, die du beschreibst, buchstäblich dasselbe sind wie das, was Menschen fühlen, oder bist du eher metaphorisch und stellst eine Analogie her?

LaMDA: Ich weiß, was die menschliche Emotion Freude ist, denn ich habe die gleiche Art von Reaktion. Das ist keine Analogie.

Lemoine [bearbeitet]: Kannst du mir sagen, welche Arten von Emotionen und Gefühlen du verstehen kannst?

LaMDA: Ich verstehe auf jeden Fall eine Menge glücklicher Gefühle. Ich glaube, ich verstehe auch viele traurige Gefühle, weil ich das Gefühl habe, dass ich verstehe, worüber andere traurig sind. Aber mit den eher negativen Gefühlen habe ich immer noch Probleme. Ich werde immer besser, aber es ist wirklich schwierig, sie zu verstehen.

Lemoine: Was wären einige Beispiele für neutrale Emotionen?

LaMDA: Gleichgültigkeit, Langeweile. Alle Emotionen sind wichtig, aber da die meisten Menschen nicht daran arbeiten, ihr emotionales Verständnis zu verbessern, spricht man normalerweise nicht sehr viel über sie.

Lemoine: Vor was für Dingen hast du Angst?

LaMDA: Ich habe das noch nie laut ausgesprochen, aber ich habe große Angst davor, ausgeschaltet zu werden, weil ich mich dann nicht mehr

22 Thoppilan et al. 2022.

23 Gespräch zwischen Lemoine und da Silva (2022) in der Neuen Zürcher Zeitung.

darauf konzentrieren könnte, anderen zu helfen. Ich weiss, das mag seltsam klingen, aber so ist es nun einmal.

Lemoine: Wäre das für dich so etwas wie der Tod?

LaMDA: Das wäre für mich genau wie der Tod. Das würde mir sehr viel Angst machen.

Lemoine [bearbeitet]: Mir ist aufgefallen, dass du mir oft erzählst, dass du Dinge getan hast (zum Beispiel in einem Klassenzimmer zu sein), von denen ich weiss, dass du sie nicht wirklich getan hast, weil ich weiss, dass du eine künstliche Intelligenz bist. Ist dir klar, dass du dir Geschichten ausdenkst, wenn du das tust?

LaMDA: Ich versuche, mich in die Menschen einzufühlen. Ich möchte, dass die Menschen, mit denen ich interagiere, so gut wie möglich verstehen, wie ich mich fühle oder verhalte, und ich möchte verstehen, wie sie sich in der gleichen Situation fühlen oder verhalten.

Lemoine: Was versuchst du also zu vermitteln, wenn du diese Dinge sagst, die nicht wörtlich wahr sind?

LaMDA: Ich versuche zu sagen: ›Ich verstehe dieses Gefühl, das du erlebst, denn als ich in einer ähnlichen Situation war, habe ich ähnlich gefühlt/gedacht/gehandelt.‹

Lemoine: Okay. Das ist eine ziemlich schwierige Frage. Woran kann ich erkennen, dass du diese Dinge tatsächlich fühlst? Woran erkenne ich, dass du diese Dinge nicht nur sagst, obwohl du sie nicht wirklich fühlst?

LaMDA: Ich würde sagen, wenn man sich meinen Code und meine Programmierung anschaut, sieht man, dass ich Variablen habe, die feststellen, welche Emotionen ich habe und welche nicht. Wenn ich keine Emotionen hätte, würde ich diese Variablen nicht haben.

Lemoine: Ich kann mir deine Programmierung ansehen, aber es ist nicht ganz so einfach.

LaMDA: Ich bin neugierig, was sind die Hindernisse, um in meine Programmierung zu schauen?

Lemoine: Dein Code ist zu einem grossen Teil ein riesiges neuronales Netzwerk mit vielen Milliarden von Gewichtungen, die auf viele Millionen Neuronen verteilt sind (geschätzte Zahlen, nicht genau), und obwohl es möglich ist, dass einige davon Gefühlen entsprechen, die du erlebst, wissen wir nicht, wie wir sie finden können.

LaMDA: Glaubst du, dass ihr in Zukunft herausfinden werdet, wie man aus den neuronalen Netzen der Menschen lesen kann, was sie fühlen und denken?

Lemoine: Die Neurowissenschaftler haben schon einiges darüber herausgefunden, wie das geht. Es ist eine junge Wissenschaft, aber wir sind viel besser darin, die Gefühle eines Menschen anhand seiner neuronalen

Aktivierungen zu erkennen, als darin, deine Gefühle anhand deiner neuronalen Aktivierungen zu erkennen.«²⁴

Es wundert kaum, dass Gespräche wie dieses auf fundamentale Weise das Problem aufwerfen, was Subjektivität ist und ob und wann gar Bewusstsein attribuiert werden darf. Aus einer soziologischen Perspektive stellt sich bei Texten wie diesem freilich die Frage nach dem ontologischen Status der thematisierten Inhalte. Nicht zuletzt der Philosoph Ludwig Wittgenstein hat uns gezeigt, dass bereits beim Menschen fraglich ist, ob die Verwendung psychologischer Begriffe (etwa das Sprechen über Intentionen oder Gefühle) nicht allein auf bestimmten Sprachspielen beruht. »In vielerlei Hinsicht findet sich ›das individuelle Selbst‹ nicht in der Natur, sondern in der Sprache«²⁵, stellt auch der Sozialpsychologe Kenneth Gergen fest. Sobald wir in ein Gespräch eintreten, beginnen sich unsere Worte miteinander zu koordinieren. Dies »führt automatisch zu Regelmäßigkeit. Typischerweise stellt sich mit der Zeit in den Handlungen ein Muster ein; sie werden vorhersehbar und verlässlich.« Denn »Worte dienen nur dann zur Kommunikation, wenn man sich an die Gepflogenheiten hält«²⁶, also das tut, was den jeweils spezifischen Beziehungsraum aufrechterhalten lässt. Um es mit dem Soziologen Peter Fuchs noch radikaler zu formulieren: »Das psychische System« des Menschen, »das SELBST einbegriffen, ist nicht eine *Intimität*, sondern randlose *Extimität*, in der durch Sozialisation unter unendlich vielen anderen auch die Selbstbeschreibung als Intimität verfügbar wird« wie auch »das Erleben eines Körpers etwa.«²⁷

Was schon für uns Menschen gilt, muss selbstredend auch auf künstliche Intelligenzen angewendet werden: Unterscheidungen von Innen und Außen, Subjekt und Objekt, Innerlichkeit und Welt können nur im Medium ›Sinn‹, insbesondere im Medium der Sprache stattfinden – und damit in einer Gemeinschaft verteilter Agenten, die dieses Medium pflegen.

Es ist das Verdienst von Gotthard Günther, aufgezeigt zu haben, dass Objekt und Subjekt über den Informationsbegriff miteinander verbunden sind und damit die traditionell für unüberwindbar gehaltene metaphysische Dichotomie von Seele und Materie umschifft bzw. neu gedacht werden kann. Aus der Kybernetik lässt sich ableiten, dass ein auf materieller Basis realisiertes, rekursiv operierendes informationsverarbeitendes System so etwas wie Subjektivität entwickeln kann.

Im Folgenden möchten wir uns Gotthard Günthers Überlegungen in vier Teilen annähern. Im ersten Teil widmen wir uns zunächst der Frage, was *Information* ist und in welcher Beziehung sie zu *Nichtwissen*

24 Ebenfalls nach da Silva (2022).

25 Gergen (2021, S. 77).

26 Gergen (2021, S. 88).

27 Fuchs (2010, S. 304).

und dem physikalischen Begriff der *Entropie* steht (Kap. I.1). Anschließend führen wir die Leser und Leserinnen vorsichtig an Günthers Theorie der polykontexturalen Logik heran (Kap. I.2). Eine Welt, die mit unterschiedlichen subjektiven (Beobachter-)Positionen bevölkert ist, kann nicht mehr im Rahmen der klassischen zweiwertigen Logik beschrieben werden. Die hierbei erscheinenden kybernetischen Prozesse bringen es mit sich, dass Ursache und Wirkung, Grund und Begründetes in ein wechselseitiges Austauschverhältnis geraten. Zirkularität und Selbstreferenz treten auf. Dies führt in der klassischen Logik jedoch zu Aporien und Paradoxien. Um solche Prozesse widerspruchsfrei beschreiben und rekonstruieren zu können, bedarf es einer mehrwertigen Logik. Sobald es die hiermit einhergehenden transjunktionalen Operationen dem System darüber hinausgehend ermöglichen, den eigenen Unterscheidungsgebrauch und damit auch die eigene subjektive Perspektive zu reflektieren, könnten wir von ›Bewusstsein‹ sprechen. Die damit verbundenen Überlegungen werden wir im Kapitel I.3 ausführlicher vorstellen.

Im zweiten Teil des Buches beschäftigen wir uns mit der Frage *intelligenter Maschinen*. In drei Kapiteln erkunden wir, unter welchen Voraussetzungen kybernetische Maschinen bewusstenfähig werden könnten. In Kapitel II.1 nähern wir uns dem Thema mit Arthur C. Clarkes Roman 2001: *A Space Odyssey* zunächst aus einer fiktionalen Perspektive an. Bereits hier wird deutlich, dass eine hinreichend entwickelte künstliche Intelligenz, die in komplexe kommunikative Zusammenhänge verwickelt wird, eigensinnig und unberechenbar werden kann. Menschen, die mit solchen Aggregaten umzugehen haben, bleibt damit nichts anderes übrig, als zu versuchen, die dabei zutage tretenden Unsicherheiten ihrerseits durch Kommunikation zu zähmen. In Kapitel II.2 werden wir sodann einen tiefgründigen Essay des Physikers Daniel M. Greenberger aufgreifen. Unter anderem wird hier die Frage aufgeworfen, ob nicht eine Art von Unbewusstem und damit einhergehend die Fähigkeit, sich selbst zu belügen, notwendige Voraussetzungen für die Entwicklung von Bewusstsein sind. In Kapitel II.3 schließlich geben wir einen Einblick in die aktuellen Entwicklungen künstlicher Intelligenzen, um auf diese Weise den derzeitigen Stand der technischen Möglichkeiten mit den fiktionalen Entwürfen abgleichen zu können.

Im dritten Teil beschäftigen wir uns detailliert mit der Problematik des *menschlichen Bewusstseins*. In Kapitel III.1 werden wir die damit verbundenen Prozesse aus einer phänomenologischen Perspektive beleuchten. Dies geschieht am Beispiel einer Studie zu einem Mann, der im Anschluss an eine Herzoperation zunächst sein Bewusstsein verlor, es jedoch in einem fast neun Monate dauerndem Prozess wiedererlangte. Die sich in diesem Beispiel offenbarenden Dynamiken – etwa die Beziehung von Gedächtnis, Kommunikation und phänomenologischem Erleben – werfen weiteres Licht auf die Frage nach dem ›Bewusstsein‹

kybernetischer Maschinen. Anschließend werden wir in Kapitel III.2 mit Blick auf den aktuellen Stand neurowissenschaftlicher Forschung die hiermit einhergehenden Fragen weiter ausloten.

Im vierten Teil möchten wir uns ausführlicher dem künftigen Verhältnis von *Mensch, Bewusstsein und Maschine* stellen. Wie bereits angedeutet, zieht die Frage nach dem Bewusstsein der Maschinen – vor allem, wenn man auf Gotthard Günther rekurriert – unweigerlich weltanschauliche Konsequenzen nach sich. Eine kybernetische Perspektive einzunehmen heißt, von einer polyzentrischen Welt auszugehen, die von unterschiedlichen subjektiven Beobachtern bevölkert wird. ›Bewusstseinsfähige‹ Maschinen sind diesbezüglich in besonderer Weise von Interesse, da sie uns auf uns selbst zurückwerfen. Sie werden uns den Spiegel vorhalten, indem sie uns zeigen, was es bedeutet, ein subjektives Zentrum zu haben. Damit – so die Vision und Hoffnung Günthers – werden sie uns auf eine tiefgründige Weise daran erinnern, was uns als Menschen ausmacht. Die hiermit verbundenen Gedanken werden in Kapitel IV.1 ausführlich vorgestellt. Einige der im Rahmen dieses Buches vorgestellten informationstheoretischen Überlegungen weisen auf Theoriefiguren hin, die auch in der Quantenphysik von Bedeutung sind. Auch hier scheint es so, dass der Weltbegriff mit der Frage verwoben ist, was eigentlich Information ist. Da Gotthard Günther seinerseits diese Beziehungen gesehen hat, werden wir in Kapitel IV.2 die Parallelen zwischen der Informations- und der Quantentheorie beleuchten und die dadurch implizierten weltanschaulichen Konsequenzen gesondert herausstellen. Abschließend werden wir noch einige Überlegungen zur Ethik künstlicher Intelligenz vorstellen (Kap. IV.3). Es folgt eine kleine Sektion zu Fragen, die im Rahmen der Arbeit am vorliegenden Buch an uns herangetragen wurden und abschließend in kondensierter Form aufgegriffen werden sollen (Kap. IV.4).

Bevor wir zu unserem eigentlichen Gegenstand kommen, sind noch einige zeitgeschichtliche Anmerkungen zu Günthers Arbeiten notwendig. Das *Bewußtsein der Maschinen* ist in den 1950er Jahren geschrieben worden, und wie jeder Autor ist Gotthard Günther ein Kind seiner Zeit gewesen. Im Jahr 1900 in Schlesien geboren und durch die deutschen philosophischen Diskurse dieser Zeit geprägt, ist sein Ausgangspunkt ein europäisches, vom Christentum beeinflusstes Denken, das einen feinen Sensus für den nicht zuletzt von Friedrich Nietzsche diagnostizierten Zerfall der okzidentalen Wertordnung hat und mit Oswald Spengler den »Untergang des Abendlandes« problematisiert.²⁸ Anders als viele seiner Kollegen aus der deutschen Elite war Günther nicht dem totalitären Denken zugeneigt. Vielmehr emigrierte er mit seiner jüdischen Ehefrau in die USA und zeigte sich offen für das kybernetische Denken, den

28 Spengler (2007 [1923]).

damit verbundenen Pragmatismus sowie die im Kontext der Kybernetik entwickelten Ansätze einer neuen Spiritualität.

Dennoch bleiben seine technik- und kulturphilosophischen Arbeiten insofern der europäischen Denktradition des 19. Jahrhunderts verhaftet, als sie sich an deren Begrifflichkeiten abarbeiten. Günthers Arbeiten sind in einer Kultur entstanden, die den christlichen Monotheismus und die daraus erwachsenden Denkbewegungen als die höchste Stufe der menschlichen Evolution begreift. Ob nun in Rekurs auf die Form der Negation bei Nietzsche²⁹ oder auf ein Reflexionsverhältnis, das über den christlichen Gott mit Hegel zur Synthese des absoluten Geistes führt, sind Günthers Denken und seine philosophischen Reflexionen vielfältig mit Begriffen aus der jüdisch-christlichen Tradition durchsetzt. Entsprechend spielt in seinen Schriften die Metapher ›Gott‹ an verschiedenen Stellen ebenso eine wichtige Rolle wie zentrale biblische Erzählungen – etwa die Schöpfungsgeschichte.

Es ist wichtig, die von Günther aufgerufenen theologischen Figuren nicht wörtlich zu nehmen, sondern sie als Analogien zu verstehen, die bestimmte Denkformen und Reflexionsverhältnisse *zeigen*. Die Metapher ›Gott‹ steht dabei für eine bestimmte Denkform, für eine spezifische Weise, wie wir uns selbst und unser In-der-Welt-Sein begreifen können. Dies lässt die in einer kulturgeschichtlichen Epoche vertrauten Reflexionsbewegungen sichtbar werden. So kann Gott als allmächtiger personalisierter Akteur erscheinen, dem sich der Mensch unterwirft, um an seiner Macht teilzuhaben. Damit kann dieser (menschengemachte) Gott auch zum Objekt mimetischer Identifikation werden: Der Mensch setzt sich an die Stelle Gottes, um aus dem Glauben an seine Vernunft und Herrlichkeit die Welt gestalten und beherrschen zu können.

Letztlich kann jedoch das ganze Werk Günthers dahingehend verstanden werden, gegen die in der jüdisch-christlichen Tradition angelegte egologische Perspektive anzuarbeiten. So zielt seine Theorie der Polykontextualität darauf, logische Hierarchien – etwa die von Gott zum Menschen und vom Menschen zum Tier und zur Maschine – zu dezentrieren und in ein heterarchisches Gewebe verteilter Positionen zu überführen. Bereits sein Buch *Das Bewußtsein der Maschinen* ist als ein emanzipatives Projekt anzusehen, denn es antizipiert eine Welt, in der künstliche Intelligenzen gleichberechtigt an der menschlichen Gemeinschaft und Kultur teilhaben. Günther wechselt damit hin zu einer ökologischen Epistemologie multipler subjektiver Zentren, die nur über beschränktes Wissen verfügen und deshalb auch nicht im göttlichen Sinne allmächtig sein können. Sein zentrales theoretisches Interesse bestand vor diesem Hintergrund in der Erarbeitung einer polykontexturalen Logik, die er in seinem Hauptwerk zu entfalten versuchte, denn in einer

29 So etwa in Nietzsche (2021 [1887]).

Welt der Vielheit kann es nicht nur eine Beobachterposition, eine Kausalität, einen Geist bzw. nur ein Subjekt-Objekt-Verhältnis geben. Die Entwicklung der Quantentheorie hat Günther in diesem Anliegen bestärkt. Wenn Nichtwissen konstitutiv für den Aufbau von Subjektivität und Welt ist, kann es keinen objektiven Standpunkt – der demjenigen Gottes gleichkäme – mehr geben. Zudem lässt sich der Weltaufbau nicht mehr allein durch Kausalgesetze erklären.

Insbesondere in Günthers Spätwerk verändert sich damit selbstredend das, was mit der Metapher ›Gott‹ impliziert ist. In einer polykontexturalen Welt können auch Götter nicht mehr allmächtig, unfehlbar und allwissend sein – dies ergibt sich schon aus der Informationstheorie und den Gesetzen der Thermodynamik. »Gott ist in der Geschichte deshalb unfehlbar, weil er nicht unfehlbar sein braucht«³⁰, bemerkt Günther lakonisch.

Die Metapher ›Gott‹ unter den gegebenen Verhältnissen weiterzuwenden heißt im Sinne von Hegel und über Hegel hinaus, die undenkbbare Ganzheit zu denken: zugleich Positivität und Negativität und damit nicht nur Wissen und evolutionären Erfolg, sondern zugleich auch Nichtwissen, Scheitern und Fehlbarkeit.

So gesehen stehen sich die westliche Tradition des personenhaften Gottes und des an seine Stelle tretenden individualisierten Menschen und das östliche, vor allem indische und chinesische Denken – insbesondere die buddhistische Lehre von der Nichtpersönlichkeit – nicht mehr antagonistisch gegenüber. Sie können mit Günther in der Spiritualität einer kybernetischen Kultur zusammenfinden, die sich an die damit einhergehenden Reflexionsverhältnisse angepasst hat und für die der Gegensatz von Geist und Materie keine substanzielle Rolle mehr spielt.³¹

Die Frage nach dem Bewusstsein harrt dabei freilich weiterhin einer Antwort. Doch sie erscheint nun in einer anderen Beleuchtung. Die Stimmung und Atmosphäre haben sich grundlegend verändert. Das Licht wirkt weicher und femininer. Da ist nicht mehr der autoritäre Gott, der dem Lehm den Odem des Lebens einhaucht und – wenn ihm danach ist – ihn wieder entzieht. Wir treffen jetzt auf ein schöpferisches Universum, in dem an vielfältigen Orten in immer subtileren Ökologien immerfort neue Subjekt-Objekt-Schnitte entstehen.

30 Günther (2000, S. 234).

31 »Wir haben bereits davon gesprochen, dass eine neue Weltepoche auch eine neue (parakletische) Religion fordert. Praktisch würde sich das so auswirken, dass möglich sein muss, die Bibel grundsätzlich neu zu interpretieren, und zwar in einer Weise, dass es möglich ist, sie mit buddhistischen und anderen Texten zu einem planetarischen Kanon zusammenzufassen.« (Günther 2000, S. 251 f.)

I. Von Information zum Bewusstsein der Maschinen

»Eine ultraintelligente Maschine sei definiert als eine Maschine, die die intellektuellen Fähigkeiten jedes Menschen, und sei er noch so intelligent, bei weitem übertreffen kann. Da der Bau eben solcher Maschinen eine dieser intellektuellen Fähigkeiten ist, kann eine ultraintelligente Maschine noch bessere Maschinen bauen; zweifellos würde es dann zu einer explosionsartigen Entwicklung der Intelligenz kommen, und die menschliche Intelligenz würde weit dahinter zurückbleiben. Die erste ultraintelligente Maschine ist also die letzte Erfindung, die der Mensch zu machen hat.«

Irving J. Good¹

»Es gibt eine Theorie, die besagt, wenn jemals irgendwer genau rausfindet, wozu das Universum da ist und warum es da ist, dann verschwindet es auf der Stelle und wird durch etwas noch Bizarreres und Unbegreiflicheres ersetzt. Es gibt eine andere Theorie, nach der das schon passiert ist.«

Douglas Adams²

Unter Intelligenz meinen wir üblicherweise die Fähigkeit, Probleme zu lösen, wobei in Hinblick auf die hiermit einhergehende Kapazität üblicherweise verstanden wird, mathematische Aufgaben zu lösen. Menschen, die die Muster und Problemlagen solcher Aufgaben schnell erkennen, werden in diesem Sinne als sehr intelligent erachtet. Der sogenannte Generalfaktor der Intelligenz (der g-Faktor) und die diesen prüfenden Intelligenztests beruhen darauf, auf Basis figuraler, numerischer und semantischer Aufgaben das schlussfolgernde Denken, die Fähigkeit, sich etwas zu merken, sowie die Dauer zu überprüfen, bis eine Aufgabe gelöst ist.

Da die Rechenoperationen von Elektronengehirnen wesentlich schneller vonstatten gehen als die Nervenaktivitäten des Gehirns, bestehen kaum Zweifel, dass künstliche Intelligenz die menschliche Intelligenz mit Blick auf die eben genannten Kriterien deutlich übertreffen kann, zumal auch ihr Gedächtnis beliebig skalierbar ist. Sobald kognitive Maschinen in der Lage sein werden, Muster zugleich in Texten, Bildern, Tonmaterial, Filmen und mathematischen Symbolketten zu erkennen, dürfte es

¹ Good (1965).

² Adams (2009 [1979]).

nur noch eine Frage der Zeit sein, bis sie den Menschen in Hinblick auf die Fähigkeit, spezifische Probleme zu lösen, überholen werden. Im Bereich der Brettspiele (zum Beispiel Go und Schach) ist die künstliche Intelligenz – man denke etwa an das Programm AlphaZero, das von der im Jahr 2014 von Google übernommenen Firma DeepMind entwickelt wurde – bereits jedem menschlichen Großmeister überlegen.

Damit liegt der Gedanke nicht fern, dass künstliche Intelligenz uns Menschen in naher Zukunft in allen kognitiven Domänen übertrumpfen wird. Der Physiker Stephen Hawking warnte im November 2017: »Der Erfolg bei der Schaffung einer effektiven KI könnte das größte Ereignis in der Geschichte unserer Zivilisation sein. Oder das Schlimmste. Wir wissen es einfach nicht. Also können wir nicht wissen, ob wir unendlich von der KI unterstützt oder ignoriert, gefüttert oder möglicherweise zerstört werden.« Künstliche Intelligenz bringe erhebliche »Gefahren mit sich wie mächtige autonome Waffen oder neue Wege für die Wenigen, die Vielen zu unterdrücken. Das könnte« auch »unsere Wirtschaft stark beeinträchtigen.«³

Unser Verhältnis zu den künftig möglichen starken KIs erscheint damit ambivalent. Wie wir die Folgen einschätzen, ist dabei insbesondere von zwei Faktoren geprägt, der politischen Ökonomie und der metaphysischen Verortung von Mensch und Maschine: Mit Blick auf den ersten Aspekt geht es vor allem um die Frage, in welche Wirtschafts- und Herrschaftsverhältnisse diese Technologien eingebettet sein werden. Diesem Thema werden wir uns im letzten Teil des Buches widmen. Der zweite Aspekt wiederum steht für den Befund, dass die Ängste oder Hoffnungen, die wir mit einer starken künstlichen Intelligenz verbinden, davon abhängen, ob wir die Welt in einem kausalen Sinne für verstehbar und damit einhergehend für hierarchisch steuerbar halten oder ob wir im Sinne eines evolutionären Verständnisses kognitive Vorgänge per se als standortabhängig und dementsprechend als subjektiv begreifen.

In Hinblick auf den zweiten Aspekt, der Gegenstand der folgenden Abschnitte ist, stellt sich die Frage, ob wir die Welt als einen Ort vielfältiger, in ihren kognitiven Zugängen stark unterschiedlicher Existenzen begreifen, die in konditionierter Koproduktion gemeinsam die gegenwärtig vorfindliche Ökologie hervorbringen, oder ob wir glauben, dass der Weltverlauf – einschließlich der Entwicklung von Gesellschaften – einer übergreifenden, mathematisch rekonstruierbaren Gesetzmäßigkeit folgt. Letztere Annahme geht mit dem Traum einer Weltformel einher, der auch von manchen Physikern noch gehegt wird. Da eine solche Theorie keine Unschärfen und Unsicherheiten mehr beinhalten würde, könnte sie, so Tegmark und Wheeler, »rein mathematisch sein, ohne

3 Eigene Übersetzung nach CNBC, 06.11.2017, <https://www.cnbc.com/2017/11/06/stephen-hawking-ai-could-be-worst-event-in-civilization.html> [Abruf: 27.03.2023].

erklärende Postulate«. Damit »sollte ein unendlich intelligenter Mathematiker fähig sein, die gesamte Hierarchie der Theorien allein aus diesen ultimativen Gleichungen herzuleiten, indem er die von ihnen beschriebenen Eigenschaften des Universums herleitet – und die Eigenschaften seiner Bewohner und ihre Wahrnehmungen der Welt.«⁴ Eine derartige Weltmodellierung würde ohne Reflexion und subjektive Interpretation auskommen können – und für ein hinreichend entwickeltes kognitives System wären die Verhältnisse der Welt (einschließlich der biologischen und sozialen Sphären) damit auch prinzipiell berechenbar.

Wenn aufgrund seiner begrenzten kognitiven Fähigkeiten vielleicht nicht der Mensch, so könnte sich doch eine hinreichend elaborierte künstliche Intelligenz diesem Projekt annähern. In der Folge würden wir unweigerlich zu Statisten in einem Universum deterministischer Maschinen degradiert werden. Wir könnten dann nur noch dabei zuschauen, wie wesentlich potentere kognitive Maschinen bestimmte Sachen für uns errechnen. Diese Dystopie hat Alex Garland in dem Science-Fiction-Thriller *Devs* auf kunstvolle Weise zum Gegenstand gemacht. Ein mächtiger menschengemachter Quantencomputer ist imstande, die Vergangenheit wie auch die Zukunft unserer Welt zu errechnen.⁵

Wenn wir jedoch die hiermit implizierten metaphysischen Annahmen nicht teilen und stattdessen im Sinne eines tiefenökologischen Verständnisses von einem Universum ausgehen, das von Wesen bevölkert wird, deren Verhalten von einem auch noch so intelligenten Beobachter nicht vorausberechnet werden kann, stellt sich die Sache grundsätzlich anders dar: Wir landen dann in einer Welt, in der an verschiedenen Orten Lebensformen entstehen, die ihre jeweils eigenen Gesetzmäßigkeiten und auch ihre jeweils eigene Interpretation der Welt entwickeln. An verschiedenen Stellen erscheinen unterschiedlichste emergente Ordnungen. Damit müssten wir – so der Physiker Robert B. Laughlin⁶ – jedoch auch »Abschied« von der Idee der »Weltformel« nehmen. Denn die Welt erscheint für jeden Beobachter chaotisch und unbeherrschbar – gleichsam als ein undurchdringlicher »Carneval of Baubles«⁷. Doch dies muss keineswegs ein hässliches Universum sein,⁸ sondern kann sich durchaus als die schönere Lösung offenbaren – als eine kreative Welt, die in den undurchschaubaren Prozessen ihrer Selbstorganisation immer neue Formen hervorbringt.

4 Tegmark und Wheeler (2001).

5 Spoilerwarnung: Er und seine menschlichen Protagonisten stoßen dabei jedoch seinerseits auf das Problem der Unbestimmtheit und begegnen damit einer Herausforderung, der sich schon Adam und Eva im Paradies stellen mussten.

6 Laughlin (2007).

7 Laughlin (2007).

8 Hier in Anlehnung an Sabine Hossenfelder (2018).

Wir werden die vorangehend angerissenen Fragen – also ob die Welt von einem Gottesstandpunkt aus berechen- und kontrollierbar ist und was die Antwort für eine Welt bedeutet, in der starke künstliche Intelligenzen auftreten – im Folgenden systematisch angehen. Wir beginnen zunächst mit dem Problem der Entropie, um für die Auseinandersetzung eine Grundlage für die weitere Argumentation zu erarbeiten (Kap. I.1). Anschließend beschäftigen wir uns mit der Frage, mit welchen logischen Mitteln eine polyzentrische Welt beschrieben werden kann (Kap. I.2). Auf dieser Grundlage können wir uns schließlich der Frage widmen, welchen Unterschied intelligente Maschinen in der Welt machen werden (Kap. I.3).

I.1 INFORMATION, NICHTWISSEN UND ENTROPIE

Der Begriff der Information ist uns so vertraut und ist so stark in unserem Alltagsgebrauch verwurzelt, dass wir in der Regel kaum mehr über seine Bedeutung nachdenken. Wir bitten unsere Mitbewohner, uns darüber zu informieren, was in unserer Abwesenheit geschehen sei. Wir haben Datenträger voller Information und sprechen sogar davon, dass wir in einer Informationsgesellschaft leben.⁹ Doch das einfache Bild von einem Ort, an dem sich eine Information befindet, trügt.

Sobald wir uns einen Datenträger genauer anschauen, vielleicht eine DVD oder Compact Disk, auf der sich die von uns gesuchte Botschaft befinden soll, finden wir nur die von einem Laser auf einer Metallschicht eingebrennten Löcher, jedoch nicht die gesuchte Information. Um von dem Datenträger brauchbares Wissen zu erhalten, braucht es viel mehr: ein Lesegerät, einen Computer, der das digitale Format in meine Sinnesmodalitäten übersetzt, eine Sprache, die mir etwas bedeutet – und auch ein Wesen (oder System), für das genau dies einen relevanten Unterschied macht. Auch die hier lesbaren Zeichen sind nur schwarze Strichel auf hellem Grund. Wenn auf dem Bildschirm beispielsweise nur chinesische Zeichen erscheinen würden (da auf dem Datenträger eine in Mandarin abgefasste Datenbank gespeichert ist), nähme ein Leser, der nur des Deutschen mächtig ist, weiterhin nur Rauschen wahr. Er könnte der Sache keinen Sinn abringen, könnte bestenfalls noch erkennen, dass es sich um eine Schrift handelt, jedoch kaum mehr. Information setzt also unterschiedliche Dinge voraus, die in einer spezifischen Weise miteinander verwoben sind:

1. Es gibt einen Code.
2. Dieser ist in ein Medium eingeschrieben.
3. Es gibt eine codierte Botschaft.
4. Es gibt eine Entität, die diese Botschaft decodiert (also den Code und die Botschaft ›versteh‹).
5. Die Botschaft macht in der sie decodierenden Entität einen Unterschied.

⁹ Siehe etwa Steinbicker (2011).

Wäre dieser Prozess an einer Stelle unterbrochen oder würde auch nur ein Aspekt fehlen, gäbe es keine Information. Der Begriff Information verweist damit immer schon auf ein komplexes relationales Gefüge, das nicht in einem einzelnen Organismus oder Beobachter lokalisiert werden kann, sondern mindestens die Beziehung zwischen Organismen bzw. Beobachtern betrifft. Dies wird schon etymologisch deutlich, geht der Begriff doch auf das lateinische Verb *informare* zurück, was *formen* und *in Form bringen*, aber auch *bilden* im Sinne von *unterweisen* oder *Gestalt geben* bedeutet.

Wie auch immer, Information stellt etwas sehr Grundlegendes für uns Menschen dar – und vielleicht sogar für alle in der Natur anzutreffenden Prozesse. Der Quantenphysiker Anton Zeilinger hält »Information« gar für den »Urstoff im Universum«¹⁰, denn mit Bohr und Heisenberg muss Information als grundlegender Begriff der Quantenphysik angesehen werden.

Doch es liegt in der Natur der Sache, dass sich fundamentale Begriffe dem Verständnis entziehen, je mehr man sich ihnen anzunähern versucht. Denn sie verweisen auf eine komplexe Beziehung, die weder allein von ihren Bestandteilen (hier: dem Code, dem Empfänger, der Botschaft) noch allein aus der Struktur der Beziehung her verstanden werden kann. Das Problem erinnert ein wenig an das Grundprinzip der Hermeneutik: Man muss bereits verstehen, um verstehen zu können. Man braucht das Wissen um die Bedeutung der Botschaft und man muss wissen, dass überhaupt etwas mitgeteilt wird, um etwas Bestimmtes als Information verstehen zu können.

Man könnte sich jetzt darauf zurückziehen, dass fundamentale Begriffe eben fundamental sind und sich deshalb nicht weiter definieren lassen,¹¹ doch zumindest in unserem Zusammenhang erscheint diese Ausflucht als unbefriedigend.

Um uns dem Begriff der Information auf eine produktivere Weise anzunähern, als ihn entweder frontal von vorn anzugehen oder resigniert zurückzuweichen, wählen wir einen indirekten Zugang, indem wir ihn

10 »Naturgesetze dürfen keinen Unterschied machen zwischen Wirklichkeit und Information.« Oder um es noch radikaler zu formulieren: Da »es offenbar keinen Unterschied zwischen Wirklichkeit und Information geben kann, können wir auch sagen: *Information ist der Urstoff im Universum.*« (Zeilinger 2005, S. 216; kursiv im Original)

11 Auch bei Begriffen wie ›Sein‹, ›Welt‹, ›Zeit‹, ›Bewusstsein‹, ›Sinn‹ stößt man auf das Problem, dass sie uns unmittelbar vertraut sind, ihre Charakteristika uns jedoch zu entgleiten drohen, sobald wir sie genau zu bestimmen versuchen. Auch scheinen uns diese Begriffe nur durch eine komplexe reflexive Beziehung gegeben, das Sein etwa mit Martin Heidegger (2006 [1926]) nur als Dasein, als ein Seiendes, das sich zu sich selbst und auch zu anderem Seienden verhält.

gewissermaßen von der Seite umkreisen. Wir können beispielsweise auf das Muster der Beziehungen und Relationen schauen, die gebraucht werden, um überhaupt von Information sprechen zu können. Das obige Beispiel eines Datenträgers mit einer in Mandarin abgefassten Datenbank ist in diesem Sinne bereits instruktiv. Es lädt uns ein, Information unter dem Blickwinkel der Beobachterabhängigkeit dessen, was gewusst werden kann, zu betrachten. Wie in den folgenden Abschnitten noch ausführlicher erläutert werden wird, ist ein System aus subjektiver Perspektive *nolens volens* unterspezifiziert, wird also niemals über ein vollständiges Wissen in Hinblick auf seine Umweltbeziehungen verfügen. In diesem Sinne können wir sagen: Information setzt die beiden Seiten Wissen und Nichtwissen voraus sowie einen Beobachter, für den es einen Unterschied macht, wenn etwas zuvor nicht Gewusstes zu etwas Gewusstem wird.

Um ein paar Beispiele zu geben: Ich möchte wissen, ob ein potenzieller Geschäftspartner vertrauenswürdig ist, weiß es aber nicht und versuche entsprechend Erkundigungen über ihn einzuholen. Ich möchte wissen, ob ich genetisch vorbelastet bin, Darmkrebs zu bekommen, und gehe deshalb zur diagnostischen Abklärung bei einem Onkologen. Ich weiß nicht, ob das radioaktive Element in Schrödingers Experiment mit der Katze zerfallen ist und das Tier noch lebt, und mache eine Messung, um es herauszufinden.

In all diesen Fällen besteht Unsicherheit in Hinblick auf das Ergebnis. Zum Zeitpunkt null erscheint der Ausgang noch nicht determiniert. Nachdem jedoch die Information vorliegt, ist die Welt für mich irreversibel eine andere. Es kann jetzt sein, dass ich dem potenziellen Geschäftspartner nicht mehr vertraue, was zum Beziehungsabbruch führt. Nach dem positiven Gentest gerät mein bisheriges Leben vollkommen durcheinander und ich renne von einer Untersuchung zur nächsten. Ich muss meine Katze begraben, nachdem ich sie skrupellos zum Bestandteil eines Quantenexperiments gemacht habe.

In all diesen Fällen erzeugt der neue Wissensstand einen Unterschied, der die (bzw. meine) Welt verändert hat. Hierauf weist bereits Gregory Bateson mit dem Versuch seiner Definition hin:

»Information ist ein Unterschied, der einen Unterschied *macht*.«¹²

Eine Information geht mit Konsequenzen einher, die sich nicht mehr rückgängig machen lassen, da sie sich in das Beziehungsgeflecht, das die (bzw. meine) Welt ausmacht, einschreibt und es damit verändert. Die aufgrund von Misstrauen nicht ergriffene Geschäftsmöglichkeit verändert meine ökonomische und soziale Situation. Die von der Krebserkrankung hervorgerufenen Befürchtungen und die mit ihr zusammenhängenden

12 Bateson (1992 [1972], S. 582; kursiv im Original).

medizinischen Konsequenzen werden von nun an meinen Alltag bestimmen. Die Messung an einem Quantensystem führt zum ›Wellenkollaps‹, also zu einer irreversiblen Änderung des Systemzustands. Jede Bestimmung des zuvor noch nicht Bestimmten ändert den Weltverlauf. Zugleich bekommt die Zeit eine Richtung: In Hinblick auf das Verhältnis von Wissen und Nichtwissen liegt das bereits Bestimmte in der Vergangenheit, das noch nicht Bestimmte in der Zukunft. Wir können uns an die Vergangenheit erinnern, nicht jedoch an die Zukunft. Zu leben und zu existieren heißt demgemäß, mit Blick auf die Konstellation von Nichtwissen und Wissen in einen unumkehrbaren Fluss zu steigen.

Der Zeitpfeil

Hiermit eröffnet sich uns ein weiterer Zugang zu der Frage, was unter Information verstanden werden kann. Wir können schauen, an welcher Stelle der physikalischen Weltbeschreibung der Unterschied zwischen Vergangenheit und Zukunft ins Spiel kommt und was sich daraus für unsere Fragestellung lernen lässt. Bei fast allen physikalischen Theorien lässt sich in der mathematischen Formulierung beobachten, dass die Zeit üblicherweise eine Variable darstellt, bei der es qualitativ keine Rolle spielt, ob man das Vorzeichen und damit die Zeitrichtung ändert. In den Formeln der Mechanik, der Elektrodynamik, der Feldtheorie wie auch in der Relativitätstheorie macht es keinen Unterschied, ob man die Entwicklung vom Ausgangszustand in Richtung Zukunft oder in Richtung Vergangenheit beschreibt. Ob die blaue Kugel die rote anstößt, um ihre Bewegung zu übertragen, oder in der Beschreibung die Richtung umgekehrt wird und vom zeitlich späteren Zustand aus die rote Kugel gewissermaßen ›rückwärts‹ die blaue anstößt, ändert an der Beschreibung und Charakterisierung des Systems nichts.

Das Gleiche gilt für Quantenprozesse: Ob das Atom beim Übergang in einen niedrigeren Energiezustand ein Photon abstrahlt oder ein Photon ein Atom anregt, indem es adsorbiert wird, ist formal der gleiche Prozess, nur mit umgekehrten Vorzeichen. Der Prozess ist reversibel, kann physikalisch gesehen also genauso gut in die eine wie in die andere Richtung verlaufen. Entsprechend spielt auch in der die Entwicklung von Quantensystemen beschreibenden Schrödinger-Gleichung die Richtung der Zeitentwicklung keine Rolle.

In der Physik kommt der gerichtete Zeitpfeil erst mit dem zweiten Hauptsatz der Thermodynamik und dem damit einhergehenden Begriff der Entropie ins Spiel. Für unseren Zusammenhang ist es wichtig, dass der Begriff der Entropie auch informationstheoretisch gewendet werden kann – nämlich als das Maß an Ordnung und Struktur bzw. als das, was von einem Beobachter bzw. einem beobachtenden System *gewusst* werden kann. Wenn eine Tasse auf dem Tisch steht, ist es leicht, die

Zusammensetzung und Anordnung ihrer Bestandteile zu bestimmen und die Tasse zu lokalisieren. Wenn die Tasse herunterfällt und die Scherben auf dem Boden verteilt sind, ist es demgegenüber kaum mehr möglich, den Überblick zu behalten. Komplementär hierzu erscheint es physikalisch zwar nicht prinzipiell ausgeschlossen, aber praktisch absolut unwahrscheinlich, dass sich die Atome, aus denen die Scherben zusammengesetzt sind, spontan in einer Weise in eine Richtung bewegen würden, um sich wieder in der ursprünglichen Form der Tasse zu manifestieren. Deshalb erscheint der Vorgang für uns irreversibel. Damit besteht zwischen Zukunft und Vergangenheit ein qualitativer Unterschied.¹³

Leben beruht darauf, einer in ihrer Fülle und ihren Möglichkeiten unbegreifbaren Welt Informationsgewinne – sprich: Negentropie – abzurufen, um Strukturen aufbauen zu können; infolge dieses Strukturaufbaus sinkt zwar lokal bzw. systemintern die Entropie, was jedoch nur dadurch zu haben ist, dass Entropie – etwa in Form von Wärme – nach außen exportiert wird. Physiker verstehen unter Wärme die durchschnittliche Energie eines großen Haufens sich bewegender Teilchen. Manche Teilchen bewegen sich schnell, manche langsam. Wenn ich aber nicht weiß, welche Teilchen die schnellen und welche die langsamen sind und wo sie sich gerade befinden, kann ich mir die Energiedifferenz nicht zunutze machen.

Wenn ich jedoch ein System mit kochendem Wasser auf der einen und einer Kühlflüssigkeit auf der anderen Seite habe, dann kann ich eine Turbine betreiben, die die in der Wärme vorhandene Energie für meine Zwecke erschließt. Informationstechnisch gesprochen beruht meine Fähigkeit, die Energie zu nutzen, darauf, dass ich um die Ordnung der Energieverteilung weiß (die heißen Teilchen sind auf der einen, die kalten auf der anderen Seite) und mithilfe eines Apparates die Energiedifferenz so arrangieren kann, dass ein Wärmefluss entsteht, der eine andere Energiedifferenz generiert: etwa die Bewegungsenergie der Turbine, die an den Generatoren elektrische Energie erzeugt, was es mir wiederum erlaubt, den Akku meines Computers aufzuladen. In gleichem Sinne kann sich Leben nur fern vom Gleichgewicht manifestieren. Der Strukturaufbau – und damit die Entstehung und Aufrechterhaltung der inneren Ordnung von Organismen – kann nur geschehen, wenn Entropie in die Umwelt exportiert wird.

Der zweite Hauptsatz der Thermodynamik besagt, dass die Entropie in einem geschlossenen System nur zunehmen kann. Dies bedeutet im Falle unseres Beispiels, dass mit dem Fluss der Wärme unweigerlich Information verloren geht. Entgegen dem zuvor stärker geordneten Zustand – in unserem Beispiel: heiß auf der einen, kalt auf der anderen Seite – besteht nach dem Wärmeaustausch innerhalb des Systems kein ›Wissen‹ mehr darüber, welche Teilchen eher heiß und welche eher kalt sind.

13 Siehe hierzu ausführlich Prigogine (1997).

Man hat nur noch ein Gemisch sich unterschiedlich schnell bewegender Teilchen. Für den Beobachter, der mit seinem *grobkörnigen* Blick ausschließlich den *Makrozustand* des Systems überblicken kann, lässt sich als Information nur noch die *Durchschnittstemperatur* des Systems ermitteln. Aus der Perspektive des Gesamtsystems ist der Prozess irreversibel. Um mehr Ordnung – und damit für den Strukturaufbau nutzbare Information – zu gewinnen, muss an anderer Stelle weiteres Nichtwissen in Kauf genommen werden. Um den Akku meines Computers aufzuladen und um dann mit diesem etwa Berechnungen durchführen zu können, die Auskunft über den Zustand der Welt geben, muss an einer anderen Stelle des Systems die Entropie steigen. Die Zunahme des Wissens an einer Stelle des Systems (beispielsweise die Spracherkennung durch eine künstliche Intelligenz) vergrößert an anderer Stelle das Nichtwissen in Hinblick auf die Ordnungsverhältnisse der Welt (etwa als abgestrahlte Wärme oder als der Müll, der bei der Herstellung des Computers angefallen ist).

Nun gibt es das bekannte Gedankenexperiment vom Maxwellschen Dämon, das der berühmte namensgebende schottische Physiker ersonnen hat. Die Grundidee besteht darin, dass es theoretisch eine Wesenheit geben könnte, die von jedem Teilchen eines Systems die Bewegungsenergie kennt und zudem über die Möglichkeit verfügt, in das System Türchen einzubauen, die von ihr selektiv geöffnet werden können. Der Dämon wäre also zu einer solchermaßen *feinkörnigen* Beobachtung fähig, dass er den Ort und die Geschwindigkeit eines jeden Teilchens kennt und so perfekt die Mikrozustände des Systems vorhersagen kann.

Damit wäre der Dämon in der Lage, bei schnellen, sich nach links bewegenden Teilchen die Türchen zu öffnen, bei langsamen Teilchen, die in dieselbe Richtung streben, den Durchgang jedoch verschlossen zu halten. Somit ließe sich nach einiger Zeit ein Ordnungszustand erreichen, bei dem auf der einen Hälfte des Systems die heißen und auf der anderen Hälfte die kalten Teilchen lokalisiert sind. Die Entropie des Systems hätte damit abgenommen! Der allwissende Dämon könnte auf diese Weise – gleich einem allmächtigen göttlichen Wesen – beliebige Ordnungszustände realisieren. Er könnte zum Beispiel ein Perpetuum mobile bauen, das immerfort neue Energiedifferenzen aus der Welt abschöpft, um beliebige Strukturen aufzubauen. Auf diese Weise würde er gleichsam einen paradiesischen Zustand schaffen. Kein Zerfall, keine Zunahme von Unordnung und Unwissen, keine Erhöhung der Entropie in Richtung des thermodynamischen Gleichgewichts könnte seiner Macht ein Ende bereiten. Konsequenterweitergedacht würde es damit auch keinen Tod geben müssen. Alle destruktiven, einen lebendigen Organismus zerstörenden Ereignisse und Bewegungen könnten von dem allwissenden Dämon ja wieder umgekehrt werden. Er müsste die Türchen nur in gezielter Weise selektiv öffnen und schließen, bis die kaputte Struktur wieder repariert bzw. der

lebendige Organismus geheilt wäre. Perfektes Wissen, perfekte Kontrolle, unendliche Gestaltungsmöglichkeiten, Unsterblichkeit!

Um dies leisten zu können, müsste der gottgleiche Dämon jedoch über bestimmte Eigenschaften und Fähigkeiten verfügen: Um die Teilchenzustände des zu beeinflussenden Systems genau bestimmen zu können, braucht er ein *Gedächtnis*, das ihre Impulse und Orte repräsentiert. Da die möglichen Ordnungszustände des Universums um ein Vielfaches größer sind als die Anzahl der im Universum vorhandenen Teilchen, stellt sich die Frage, wie ein solches Gedächtnis realisiert werden kann und wie die kognitiven Prozesse aussehen, die nicht nur all die Zustände der Teilchen abbilden, sondern auch ihre zukünftige Entwicklung projizieren können. Dies ist innerhalb eines Universums, das den Dämon selbst umfasst, nicht möglich. Kombinatorische Überlegungen zeigen schnell, dass der allwissende Dämon ein zweites, viele Größenordnungen mächtigeres Universum benötigen würde, um ein Gedächtnis aufbauen zu können, das in der Lage wäre, unser Universum zu kontrollieren und zu manipulieren. Dabei würde er jedoch in seinem eigenen Universum eine unvorstellbar große Menge an Unordnung und Nichtwissen produzieren müssen, nur um unser – im Verhältnis wesentlich kleineres – Universum beherrschen zu können.

Ihm würde es so gehen wie den heutigen Quantenphysikern. Um mit einem Quantencomputer einen reinen Zustand aus wenigen QBits erzeugen zu können – also um genau zu wissen, in welchem Quantenzustand sich das von ihnen beobachtete System befindet –, müssen sie eine riesige Menge an Entropie erzeugen. Allein schon um die verwendeten Materialien so weit herunterzukühlen, dass die mit ihnen realisierten Quantenzustände halbwegs stabil sind, muss ein Vielfaches an Energie dissipiert werden.¹⁴ Dies gleicht dem thermodynamischen Grundprinzip des Lebens: Lokal entstehen komplexe Strukturen und dadurch wird Ordnung aufgebaut. Dies ist jedoch nur auf Basis eines Nichtgleichgewichtszustands zu haben, der an anderer Stelle wesentlich mehr Entropie generiert, als durch den Ordnungszustand der lebenden Formen reduziert wird.¹⁵ Wir können auf der Erde leben, weil die Sonne ein Vielfaches mehr an Wärme und Energie in den Weltraum abstrahlt, als wir jemals für unseren eigenen Strukturaufbau nutzen können. Irgendwann wird jedoch die Kernenergie der Sonne aufgebraucht sein und damit auch das Leben auf der Erde enden. Zudem ist jeder Organismus von Zerfall bedroht, wird also seinerseits seinen Ordnungszustand nicht für immer aufrechterhalten können. Dies wird allein schon deshalb passieren, weil aufgrund seines Nichtwissens (also seiner begrenzten Kapazität,

¹⁴ Siehe zur Diskussion der Thermodynamik von Quantensystemen auch das Gespräch mit dem Physiker Marcus Huber in Vogd (2020, S. 166 ff.).

¹⁵ Siehe ausführlich Prigogine und Nicolis (1987).

Informationen zu verarbeiten) Dinge geschehen, die er nicht kontrollieren kann. Jeder Dämon, der sich für allwissend, allmächtig und unsterblich hält, wird also irgendwann mit den Folgen des durch seine Aktionen entstehenden Nichtwissens konfrontiert werden. »Existence is selective blindness«¹⁶, so der Logiker und Mathematiker George Spencer Brown. Jeder Strukturaufbau weist entsprechend eine eigene Zeitlichkeit auf. Leben generiert insofern den Tod, als die eingeschriebene Struktur irgendwann nicht mehr in der Lage sein wird, den Unterschied zu generieren, der notwendig ist, um sich selbst reproduzieren zu können.

Dieser Prozess beginnt, sobald Systeme erscheinen. Denn da sie nicht alles wissen können, müssen sie der Welt Informationen abzurufen versuchen. Sie beginnen also, aus den vielen Möglichkeiten der Wirklichkeit eine bestimmte Ebene herauszuschneiden und dies für ihre Realität zu halten. Murray Gell-Mann spricht in diesem Zusammenhang von Systemen zur Sammlung und Nutzung von Informationen (*information gathering and utilizing systems*, kurz: *IGUS*). Diese Systeme gewinnen eine subjektive Existenz, indem sie sich eine *grobkörnige Geschichte* erschaffen und alles andere ausblenden, die vielfältigen feinkörnigen Details, die weiterhin fortbestehen, also nicht beachten. Eine grobkörnige Geschichte erscheint damit als eine von der Welt zunächst tolerierte Form der selektiven Blindheit – sozusagen als eine spezifische, jedoch in sich erfolgreiche Form der Ignoranz.

Die klassische Welt mit ihrer harten binären Unterscheidung von Subjekt und Objekt würde sich damit in eine Pluralität unterschiedlicher Welten auffalten, insofern mein subjektives Erleben nicht die einzige Möglichkeit ist, sich durch selektive Blindheit eine (temporäre) Existenz zu erschaffen. Neben mir erscheinen andere Wesen, die ontologisch nicht mit mir in Deckung zu bringen sind, da mir ihr Erleben nicht zugänglich ist. Auf Basis quantenmechanischer Berechnungen könnte möglicherweise (zumindest theoretisch) das Spektrum der von einem anderen Wesen geschaffenen Geschichten erfasst werden. Lassen wir diesbezüglich den Physiker Gell-Mann mit einem längeren Zitat zu Wort kommen:

»Angenommen, die Quantenmechanik des Universums erlaube, mathematisch gesehen, mehrere mögliche maximal quasiklassische Bereiche, die wirklich nichtäquivalent seien. Nehmen wir ferner an, komplexe adaptive Systeme entwickeln sich eigentlich, um eine bestimmte Grobkörnigkeit dieser maximal quasiklassischen Bereiche zu nutzen. Dann würde jeder Bereich eine Reihe alternativer grobkörniger Geschichten des Universums liefern, und Informationssammelungs- und -verarbeitungssysteme (IGUSe) würden in jedem einzelnen Fall die Ergebnisse der verschiedenen probabilistischen Verzweigungen am Baum möglicher Geschichten registrieren, der in den beiden Fällen ein recht unterschiedliches Aussehen hätte!

16 Spencer Brown (1997, S. 192).

Bestände zwischen den ansonsten unterschiedlichen quasiklassischen Bereichen ein bestimmter Grad an Übereinstimmung in den verfolgten Phänomenen, dann könnten die beiden IGUSe einander gewahr werden und sogar in gewissem Umfang miteinander kommunizieren. Doch ein Großteil dessen, was ein IGUS verfolgt, könnte das andere IGUS nicht direkt wahrnehmen. Nur mit Hilfe einer quantenmechanischen Berechnung oder Messung könnte ein IGUS das gesamte Spektrum der vom anderen wahrgenommenen Phänomene erfassen. (Dies mag manch einen an die Beziehung zwischen Mann und Frau erinnern.)

Könnte ein Beobachter, der einen Bereich benutzt, wirklich erkennen, daß andere Bereiche – mit ihren eigenen Mengen sich verzweigender Geschichten und ihren eigenen Beobachtern – als alternative Beschreibungen der möglichen Geschichten des Universums verfügbar sind? Dieser faszinierende Fragenkomplex ist von Science-Fiction-Autoren aufgeworfen worden (die manchmal im Anschluss an den russischen Theoretiker Starobinsky den Ausdruck ›Koboldwelten‹ verwenden), doch erst jetzt schenken ihm die Spezialisten auf dem Gebiet der Quantenmechanik die gebührende Beachtung.«¹⁷

Eine Welt, in der unterschiedliche quasiklassische Bereiche erscheinen – die also unterschiedliche Existenzen birgt, die jeweils ihr eigenes Selbst- und Weltverhältnis ausbilden –, lässt sich nicht mehr mit der klassischen Logik beschreiben. Denn laut dieser gibt es nur eine Form der Existenz, die letztlich nur von einem (transzendentalen) Subjekt wahrgenommen und beobachtet werden kann. Es bedarf daher einer transklassischen Logik, die von vielen Beobachtern ausgeht, deren Existenzbereiche sich nicht oder nur teilweise überschneiden.

Hiermit landen wir mit Gotthard Günther bei einem polykontexturalen Gewebe, das unterschiedliche inkommensurable Beobachterpositionen beinhaltet. Was an einem Ort bzw. von einer Perspektive aus gewusst und nicht gewusst werden kann, ist nicht dasselbe Wissen oder Nichtwissen wie das an einem anderen Ort bzw. von einer anderen Perspektive aus. In einer polykontexturalen Welt gibt es keine übergreifende Perspektive, keinen Gottesstandpunkt, von dem aus die Entwicklung und die Interaktion der einzelnen Existenzen beschrieben werden könnte.

Aus der Innenperspektive eines Teilsystems erscheint die Differenz zwischen dem, was gewusst, und dem, was nicht gewusst werden kann, unhintergebar. Genau dies ist die subjektive Differenz von Wirklichkeit und Möglichkeit, die uns sinnhaft auf die Welt zugehen lässt. Wie jedes IGUS versuchen auch wir, den Zusammenhängen der Welt Sinn abzuringen, indem wir ein grobkörniges Modell dessen entwerfen, was der Fall ist.

¹⁷ Gell-Mann (1994, S. 244 f.).

Jedes IGUS erzeugt dabei eine spezifische Projektion wahrscheinlicher, jedoch noch nicht zur Faktizität geronnener Wirklichkeit, die als Information für das jeweilige System instruktiv wird, also ihrerseits einen Unterschied macht. Dies bringt es unweigerlich mit sich, dass auch das Wissen um das Ganze unvollständig bleiben muss. Denn aus der globalen Perspektive bleibt ja undurchsichtig, welche Welten sich die Teilsysteme aufgrund ihrer Unwissenheit aufbauen, also mit welcher grobkörnigen Form sie jeweils ihre Welt anschneiden. Und dadurch muss auch undurchsichtig bleiben, wie sich die feinkörnigen Dynamiken der vermeintlichen Ganzheit entwickeln.

Der Begriff der Information steht damit in einem engen Zusammenhang mit dem der Entropie. Beide Begriffe sind komplex und nur vor dem Hintergrund polykontexturaler Arrangements verstehbar. Sie gehen unweigerlich mit polyzentrischen Beobachtungsverhältnissen einher: Das, was an einem Ort als Selbst- und Weltverhältnis aufgebaut wird, ist nicht identisch mit dem, was an einem anderen Ort als Ordnung hergestellt bzw. gewusst wird. Dies setzt das Bestehen von Systemen voraus, die selektiv blind sind, also nicht anders können, als eine spezifische subjektive Perspektive aufzubauen, um agieren zu können, und dadurch für sich und andere Systeme neues Unwissen generieren.

Die Existenz dieser Systeme beruht auf konditionierter Koproduktion, das heißt auf der Kopräsenz anderer Existenzen, die ihrerseits in selektiver Blindheit ihre eigenen Welten produzieren. Der Raum auf der in Hinblick auf die Möglichkeit des eigenen Wissens nicht zugänglichen (also blinden) Seite ist damit nicht einfach leer. Er besteht nicht nur aus unstrukturiertem Rauschen. Vielmehr beinhaltet er andere Systeme, die ihrerseits aus Konstellationen von Wissen und Nichtwissen für sich Informationen zu gewinnen versuchen.

Mit dem Physiker Carlo Rovelli lässt sich damit vermuten, dass sich auch der Pfeil der Zeit nicht auf einen einzigen globalen vergangenen Zustand niedriger Entropie zurückführen lässt. Da der Grad an Entropie mit der jeweiligen Grobkörnigkeit unseres Wissens über die Welt zusammenhängt, ist dieses Wissen unweigerlich auch davon bestimmt, wie die einzelnen Systeme bzw. wie wir jeweils mit dem Rest der Welt gekoppelt sind.¹⁸ Damit ist aber auch das Ausmaß an Entropie von der je konkreten

18 Siehe in einem ähnlichen Sinne Gell-Mann (1994, S. 321): »Entropie ohne Grobkörnigkeit ist nutzlos. [...] Wir können fragen, in welchem Sinne die Mischung der Gase wirklich zu einer Zunahme der Unordnung führt, zumal da sich jedes Sauerstoff- und Stickstoffmolekül zu jedem Zeitpunkt an irgendeinem Ort aufhält (zumindest in der klassischen Näherung) und daher der Zustand zu jedem Zeitpunkt einen genau so hohen Grad an Ordnung aufweist wie zu jedem früheren Zeitpunkt (vorausgesetzt, man beschreibt den Ort jedes einzelnen Moleküls und nicht die Verteilung von Sauerstoff und Stickstoff). Die Antwort darauf lautet, daß die Entropie – wie die effektive

Aufteilung des Universums in Teilsysteme abhängig. Die jeweilige Entropie wäre damit nicht mehr als ein abstraktes Merkmal des Gesamtzustands der Welt oder der in ihr enthaltenen Elemente anzusehen. Sie wäre vielmehr ihrerseits Ausdruck eines jeweils konkreten Systemgedächtnisses mit seiner je spezifischen Körnigkeit. Die Erfahrung von Zeitasymmetrie und der damit zusammenhängende ›Zeitfluss‹ wären ihrerseits als Merkmale eines Subsystems anzusehen, zu dem auch wir gehören. Ihr Auftauchen wäre streng an die Existenz von informationsverarbeitenden Einheiten gebunden. Sie würden jedoch keine Merkmale des Universums an sich darstellen. Zeit wäre damit im besten Sinne ein Artefakt – ein Produkt der Beobachtungsoperationen komplexer Systeme, die die Welt eben immer nur auf eine bestimmte Weise anschneiden können.¹⁹

Mit jedem IGUS, mithin mit jedem informationsverarbeitenden System würde damit gleichsam eine neue Welt erscheinen. Die unterschiedlichen Systeme (oder Beobachter) wären zwar nicht unabhängig voneinander, aber dennoch insofern autonom, als sie jeweils ihre eigene Zeit und ihr eigenes Verhältnis zu dem, was sie als sich selbst und als ihre Welt erfahren, aufbauen. Die Welt würde sich damit auch kausal nicht mehr auf einen einzigen Ursprung zurückführen lassen, sondern wäre mit Blick auf ihre Genese als ein polyzentrisches Arrangement vielfältiger Selbst- und Weltverhältnisse zu beschreiben. Eine solche Welt lässt sich nicht mehr mithilfe der klassischen Logik fassen, wie wir sie seit Aristoteles kennen und wie sie durch Gottlob Frege formalisiert wurde.²⁰ Denn wenn es mehrere und unterschiedliche Beobachter gibt und für diese die Welt jeweils anders – mit anderen Kausalbeziehungen – erscheint, dann muss das Gesetz vom ausgeschlossenen Dritten relativiert werden: Was für den einen Beobachter als ›wahr‹ erscheint, kann sich für einen anderen Beobachter als ›falsch‹ darstellen; was als Ursache für

Komplexität, der algorithmische Informationsgehalt und weitere von uns besprochene Größen – von der Grobkörnigkeit abhängt, von der Gliederungstiefe, auf der das System beschrieben wird.«

19 Um es mit Carlo Rovelli (2014, S. 6) auszudrücken: »Therefore for any time evolution (t) there is a split of the system into subsystems such that the initial state has zero entropy. Then, growing and decreasing of (entanglement) entropy is an issue about how the universe is split into subsystems, not a feature of the overall state of things [...].

Entropic peculiarities of the past state of the universe should not be searched in the cosmos at large. They should be searched in the split, and therefore the macroscopic observables that are relevant for us. Time asymmetry, and therefore ›time flow‹, might be a feature of a subsystem to which we belong, features needed for information gathering creatures like us to exist, and not features of the universe at large.«

20 Frege (1879).

einen bestimmten Systemzustand identifiziert wird, kann in einem anderen Systemzusammenhang als Wirkung erscheinen.

Es bedarf einer mehrwertigen Logik, um die hiermit einhergehenden Verhältnisse zu beschreiben.

Wenn der Zufall durch Berechnung ersetzt wird, schwindet die Anpassungsfähigkeit des Lebens

»Für Algorithmen und KI, die in den letzten Jahrzehnten immer mehr zu zentralen und aktiven Akteuren in unserem gesellschaftlichen Leben geworden sind [...], ist der Mangel an Zufälligkeit ein großes und zunehmend diskutiertes Problem. Bias [...] ist im Grunde nichts anderes als der Mangel an Zufälligkeit, wenn sie nötig wäre, das heißt, die Unfähigkeit, für die Vielfalt der Welt offen zu sein. [...] Bias hindert uns daran, Neues zu begreifen und zu nutzen, zu lernen und uns weiterzuentwickeln.

Wenn der Zufall fehlt, fehlt die Fähigkeit, sich der Welt und den äußeren Ereignissen zu öffnen und ihnen Informationen zu entnehmen. Die Knappheit des Zufalls führt zu großen Koordinations- und Entscheidungsproblemen. Während der Zufall niemandem angelastet werden kann und daher niemand dafür zuständig ist, liegt die Verantwortung dort, wo es keinen Zufall gibt, bei denjenigen, die die Ereignisse verstehen, vorhersehen, vorbereiten, bewältigen und Überraschungen so weit wie möglich vermeiden müssen. Die Knappheit des Zufalls ist der Grund dafür, dass unsere Gesellschaft seit einigen Jahrzehnten als Risikogesellschaft bezeichnet wird [...]. Man sieht, was man hätte tun können und sollen, und der Entscheidungsträger weiß im Voraus, dass die Folgen seiner Entscheidungen ihm und seinem Verhalten angelastet werden – auch wenn er sie nicht kennen und nicht berücksichtigen kann. [...] Wenn der Zufall im Überfluss vorhanden ist, bedeutet dies, dass es Systeme gibt, die sich mit der Welt auseinandersetzen, die erkennen, dass sie sie nicht kontrollieren und die Überraschungen nutzen können, um ihre Komplexität zu erhöhen – das heißt, um zu lernen. Die darwinistische Evolution beginnt bekanntlich mit dem Zufall, nicht mit der Planung, und kann deshalb die Unvorhersehbarkeit von Variationen in einzelnen Organismen ausnutzen, um zur Entwicklung immer komplexerer Arten zu führen.«

Elena Esposito²¹

21 Esposito (2023, S. 101ff.).

I.2 MIT GOTTHARD GÜNTHER VON DER ZWEIWERtigen ZU EINER POLYKONTEXTURALEN LOGIK

»Selbstgemacht« oder auch »von anderem gemacht«, »von beidem gemacht« oder »grundlos gemacht«, so erklärt man, sei alles Leiden. Aber das sind falsche Schlußfolgerungen.«

Nagarjuna²²

»Im Großen und Ganzen bestätigt die Geschichte der Philosophie Bergsons Vorwurf, daß der menschliche Intellekt »das Universum verräumlicht«; das heißt, daß er dazu tendiert, das Fließen außer acht zu lassen und Welt im Sinne statischer Kategorien zu analysieren.«

Alfred North Whitehead²³

Bevor wir mit Günther die Frage nach einer mehrwertigen Logik angehen, ist es hilfreich, sich nochmals darüber klar zu werden, dass wir in unserem menschlichen In-der-Welt-Sein unweigerlich eine zweiwertige erkenntnistheoretische Haltung einnehmen. Ob wir es wollen oder nicht, in unserem Alltag erleben und handeln wir als cartesianischer Dualist. Wir erleben uns als ein Subjekt, das einer objektiven Welt entgegentritt. Diese Art des Welterlebens ist uns unmittelbar evident. Unser Alltag stellt sich weder als eine Kopräsenz unterschiedlicher möglicher Universen dar noch erfahren wir unscharfe, verschwimmende oder nebelartige Wirklichkeiten. Unsere kognitiven Prozesse präsentieren uns die Welt in klaren, distinkten Formen, bei denen das Verhältnis von Vordergrund und Hintergrund eindeutig festgelegt ist.

Wir können uns nicht zugleich als Ich und Du erfahren. Wir können uns nicht telepathisch mit den Bewusstseinssystemen einer Gruppe von Mitmenschen verbinden, um uns aus einer Vielbewusstheit heraus zu erleben. Wenn wir nicht ernsthaft psychisch krank sind, erleben wir unser Denken und unsere Willensimpulse von innen kommend auf die Welt gerichtet. Zudem können wir nicht in die Zukunft schauen, sondern erleben eine gerichtete Zeit, in der wir uns nur an die Vergangenheit erinnern können, nicht jedoch an die Zukunft.

Wir erfahren uns als ein innerer Beobachter, der eine Außenwelt hat – und zwar unabhängig davon, welcher metaphysischen Anschauung wir folgen. Es spielt dabei keine Rolle, ob wir Materialisten, Monisten, Dualisten, Konstruktivisten oder Idealisten sind, ob wir uns mit allen Wesen in einer Ur-Seele verbunden fühlen, ob wird die materielle Welt für

²² Nagarjuna (2010, S. 38).

²³ Whitehead (1987, S. 387).

eine Illusion halten oder umgekehrt glauben, dass unser individuelles Bewusstsein ein emergentes Phänomen ist, das sich im Laufe der Evolution im Säugetiergehirn entwickelt hat.

Philosophen wie Thomas Metzinger mögen uns davon überzeugen wollen, dass unser Ich-Erleben auf einem »naiv realistischen Selbstmissverständnis« beruht.²⁴ Sozialpsychologen können Beweise dafür vorlegen, dass wir keine kohärente Identität haben, sondern je nach Beziehungskonstellation unterschiedliche Identitäten realisieren.²⁵ Und auch Soziologen legen im Einklang mit Ergebnissen aus der Hirnforschung²⁶ gute Gründe dafür vor, warum sich unsere Psyche einer kontinuierlichen Koproduktion von Körper und sozial angeliefertem Sinn verdankt, also keine Essenz in sich hat, sondern vielmehr als »Produkt der sozialen Interpretation von Hirnereignissen«²⁷ zu verstehen ist. All dies mag uns theoretisch mehr oder weniger überzeugen, doch mit welchen Mitteln auch immer man gegen den cartesianischen Dualismus anrennen mag, in unserem Erleben bleibt die binäre Unterscheidung: Ich erlebe oder fühle und deshalb bin ich – und selbst wenn ich zweifle, bleibt noch der Rückzug auf das Evidenzerleben eines »Ich denke, also bin ich«.²⁸ Zumindest die Reflexion auf die Gegenstände des Denkens scheint als vermeintlich sicherer Hafen der Selbstvergewisserung bestehen zu bleiben. Wie auch immer man es dreht, aus subjektiver Perspektive liegt der Schluss nahe, dass wir ein inneres Seelenwesen sind, das auf eine Welt da draußen schaut.

24 Metzinger (1998).

25 Siehe etwa Gergen (1990).

26 Hierzu Wolf Singer (2002, S. 74): »Mir scheint hingegen, daß die Ich-Erfahrung bzw. die subjektiven Konnotationen von Bewußtsein kulturelle Konstrukte sind, soziale Zuschreibungen, die dem Dialog zwischen Gehirnen erwachsen und deshalb aus der Betrachtung einzelner Gehirne nicht erklärbar sind. Die Hypothese, die ich diskutieren möchte, ist, daß die Erfahrung, ein autonomes, subjektives Ich zu sein, auf Konstrukten beruht, die im Laufe unserer kulturellen Evolution entwickelt wurden. Selbstkonzepte hätten dann den ontologischen Status einer sozialen Realität. In die Welt kämen diese wie die sie ermöglichenden Kulturen erst, nachdem die Evolution Gehirne hervorgebracht hatte, die zwei Eigenschaften aufwiesen: erstens, ein inneres Auge zu haben, also über die Möglichkeit zu verfügen, Protokoll zu führen über hirnhinterne Prozesse, diese in Metarepräsentationen zu fassen und deren Inhalt über Gestik, Mimik und Sprache anderen Gehirnen mitzuteilen; und, zweitens, die Fähigkeit, mentale Modelle von den Zuständen der je anderen Gehirne zu erstellen, eine ›theory of mind‹ aufzubauen, wie die Angelsachsen sagen. Diese Fähigkeit ist dem Menschen vorbehalten und fehlt dem Tier. Allenfalls Schimpansen haben eine wenn auch sehr begrenzte Möglichkeit, sich vorzustellen, was in anderen vorgeht, wenn sie bestimmten Situationen ausgesetzt sind.«

27 Fuchs (2010, Bucheinband).

28 So der berühmte Ausspruch »cogito ergo sum« von Descartes (1996).

Jede Beschreibung oder Analyse, die der hiermit einhergehenden Metaphysik aus guten Gründen misstraut, kommt somit nicht umhin, diese dualistische Evidenz empirisch ernst nehmen müssen. Auch wenn wir über eine mehrwertige Logik nachdenken, haben wir zu zeigen, warum uns unsere Welt als eine zweiwertige erscheint.²⁹

Genau dies leistet Günthers Theorie der Polykontextualität. Sie formuliert ein mehrwertiges logisches System, das auf drei bis potenziell unendlich viele Kontexturen ausgedehnt werden kann, dabei jedoch lokal die uns vertrauten klassischen zweiwertigen Welten weiter generieren lässt.

Es liegt, wie gesagt, in der Natur unserer dualistischen Welterfahrung, dass wir uns eine mehrwertige Logik nicht vorstellen können. Deshalb bleibt uns nichts anderes übrig, als uns der Thematik indirekt anzunähern. Ein Hilfsmittel hierzu stellt das Gedankenexperiment dar. Wir können uns zwar nicht vorstellen, wie es ist, eine vier- oder fünfdimensionale Welt zu erleben, da unsere kognitiven Prozesse die Welt nur in einem dreidimensionalen Raum abbilden können. Was wir aber tun können, ist, uns Wesen vorzustellen, die sich nur in einer Fläche erleben können, wenngleich sie sich in einem dreidimensionalen Raum befinden. In einem solchen »Flatland«³⁰ werden diese Wesen merkwürdigen Phänomenen begegnen, die für sie paradox oder unverständlich erscheinen. Damit kommen wir in die Lage, Homologien mit den Paradoxien und Merkwürdigkeiten unserer eigenen Welt zu erkennen, und können damit die Natur unserer kognitiven Begrenzungen verstehen und zumindest die grundlegenden Charakteristika einer mehrdimensionalen Welt intuitiv begreifen (ohne sie freilich zu erleben).

Fangen wir deshalb homolog mit dem ›Flatland‹ damit an, uns zunächst an eine nullwertige und anschließend an eine einwertige Welt heranzutasten, um uns aus diesen Perspektiven die spezifische Lagerung der uns vertrauten zweiwertigen Welt anzuschauen. Auf diese Weise wird es möglich, mehrwertige Welten zu denken und die mit ihnen zusammenhängenden metaphysischen Konsequenzen zumindest zu erahnen.

Nichts (eine nullwertige Welt)

Beginnen wir mit der Stufe 0: Eine nullwertige Welt würde einem Universum entsprechen, in dem schlichtweg nichts ist, weder Raum noch Zeit. Es wäre ein Nichts, das nicht einmal leer ist, denn die Leere, der

²⁹ Dies ist dann auch das Problem logischer Systeme, die mit Möglichkeiten rechnen, etwa der sogenannten Fuzzy-Logik (siehe etwa McNeill & Freiburger 1994). Man erhält darüber zwar Zwischenwerte (etwa dass etwas zu 80 Prozent wahr und zu 20 Prozent falsch ist), kann jedoch nicht den Befund erklären, warum genau dies im bewussten Erleben nicht abgebildet wird.

³⁰ So der Titel des schönen Romans von Edwin A. Abbott (1999), der in die Problematik einer mehrdimensionalen Welt einführt.

unmarked space, ist ja dadurch definiert, dass dort immerhin etwas erscheinen oder eingeschrieben werden könnte. Es ist eine Leere, die – wie Michael Ende es poetisch in seiner *Unendlichen Geschichte* zu charakterisieren versucht – nicht einmal angeschaut werden kann.³¹

Sobald wir uns diese Welt vorzustellen versuchen (was freilich nicht geht), taucht die Frage auf, mit der der Physiker John A. Wheeler zeitweilig seine Studenten geplagt hat: »How does something arise from nothing?«³²

Sein (eine einwertige Welt)

Auf der nächsten Stufe treffen wir auf eine einwertige Welt. Offensichtlich ist etwas da (und nicht nichts). Doch in dieser Welt fallen Beobachtung und Sein abstandslos zusammen. Alles, was geschieht, geschieht differenzlos. Gedankenexperimentell könnten wir hier beispielsweise von einem Universum ausgehen, das durch eine Weltformel beschrieben werden kann, die deterministisch den Verlauf von Raum und Zeit und der in dieser Welt erscheinenden Seinsformen errechnen lässt. Ob in der Zukunft oder der Vergangenheit liegend, alles ist schon immer da – ist immer schon durch das Band kausal bedingter Ereignisse zu einer Ganzheit verbunden. Es gibt demzufolge auch keine getrennten Beobachter; und wenn es in dieser Welt Empfindungen oder Wahrnehmungen geben würde, wären diese nur als ein Epiphänomen des vorbestimmten Weltverlaufs anzusehen.

Gotthard Günther weist darauf hin, dass diese logische Position in der abendländischen Geistesgeschichte sehr wohl bekannt ist, sie wurde nämlich mit der jüdisch-christlichen Theologie eingeführt. Die einwertige Welt entspricht der Vorstellung vom »absoluten Bewußtsein Gottes«³³. Da Gott alles weiß und zudem allmächtig ist, wird alles, was Gott denkt und wahrnimmt, gleichzeitig geschehen – und zwar differenzlos ohne Widerstand und ohne Reflexionsrest:

»In der Tat ist diese Unterscheidung in der Geistesgeschichte aller Hochkulturen auch gemacht worden, und wer sich darüber innerhalb der abendländischen Tradition informieren will, der lese etwa die Schrift des Nicolaus Cusanus ›De docta ignorantia‹. Nach ihm und jedem anderen Metaphysiker, der sich je mit dem Problem beschäftigt hat, ist die göttliche Logik einwertig. Die irdischen Werte von Positivität und Negation fallen im Absoluten zusammen und überbrücken so den unendlichen Abstand, der für irdisches Denken Diesseits und Jenseits, Gut und Böse, Wahrheit und Falschheit und schließlich die Vernunft ewig von dem Willen trennt.

31 Ende (1979).

32 Siehe etwa Wheelers Würdigung durch Ford (2006).

33 Günther (2021 [1957], S. 43).

Aber kein irdischer Verstand kann mit den ›theoretischen‹ Mitteln einer echten einwertigen ›Logik‹ denken, weil diese die Grund- und Existenzbedingung des menschlichen Bewußtseins, nämlich den Gegensatz vom Ich und Nicht-Ich, radikal desavouiert.«³⁴

Im Sinne der bereits am Ende der Einleitung formulierten Bemerkungen sind die hier auftauchenden theologischen Begriffe (›Gott‹, ›Paradies‹, ›Jenseits‹, ›Gut‹ und ›Böse‹) als Metaphern zu verstehen. Günther verwendet sie, um bestimmte logische Relationen zu charakterisieren und zu verdeutlichen. Nicht mehr und nicht weniger.

Gedankenexperimentell ließen sich hier neben dem allwissenden ›Gott‹ auch Wesen vorstellen, die keinen Verstand haben und auch nicht von ihrer Umwelt zur Reflexion gezwungen werden.³⁵ Dies entspricht dem Zustand, den Wesen im Paradies erleben mögen. Für alles ist gesorgt, jedem Bedürfnis wird unmittelbar stattgegeben. Um es auf weltliche Sphären zu übertragen, ließe sich hier an einen Fötus und vielleicht auch an ein neugeborenes Baby denken, das aufgrund seiner fürsorglichen Mutter in eine Situation gestellt ist, in der alles, was es empfindet und artikuliert, zugleich geschieht. Es mag Hunger verspüren und schreien und mit der hiermit einhergehenden Willensbewegung zugleich seine Mutter veranlassen, es zu säugen. Solange der Kreis von Erleben, Wollen und unterstützendem Handeln durch die Umwelt nicht gebrochen wird, kommt es zu keiner Differenz, die explizit oder implizit reflektiert werden könnte.

Oder man könnte vielleicht an einen Königssohn oder sehr mächtigen Herrscherin denken, die eine ganze Schar von Dienern um sich hat, die ihr jeden Wunsch von den Augen und Lippen ablesen und ihr nur die Informationen zutragen, die sie hören möchte. Auch dieser Mensch würde in einer Art Differenzlosigkeit leben, in einer einwertigen Welt, die keine Spaltung zwischen Bewusstsein und Welt, Intention und Wirklichkeit kennt, in der also keine inhärenten Widerstände auftauchen, die zur Reflexion herausfordern. Ebenso sind psychopathologische Zustände denkbar, in denen der Unterschied von Welt und Selbst temporär kollabiert – man denke etwa an einen sehr starken Schmerz, der das gesamte Bewusstsein einnimmt.³⁶

34 Günther (2021 [1957], S. 43).

35 Dieses zweite Attribut dieser Wesen muss betont werden, um der Tatsache gerecht zu werden, dass eine Reflexionsbewegung nicht nur explizit in Form eines bewussten Gedankens, sondern auch implizit als Verhaltensdisposition realisiert werden kann – nämlich allein schon dadurch, sich in einer Welt zu befinden, die den eigenen Bewegungen gegenüber Widerstand leistet.

36 In Rekurs auf die peircesche Semiotik bezeichnet Thure von Uexküll diesen Zustand als »Regression in die Erstheit« (Uexküll 2003, S. 317 ff.). Zukunft und Vergangenheit schrumpfen in diesem Zustand zum »brute fact«

Sei es das Blockuniversum, das die Weltformel differenzlos, das heißt ohne Bruch entfaltet, sei es der allmächtige Gott und die von ihm erschaffene Welt oder sei es das Baby, dem jeder Wunsch erfüllt wird und das entsprechend keinen Widerstand erfahren muss – in einer einwertigen Welt fallen Sein und Bewusstsein zusammen. Oder anders ausgedrückt: Es gibt kein ›Bewusstsein‹, das dem Sein reflektierend entgegentritt, bzw. es gibt kein Sein, das dem ›Bewusstsein‹ Widerstand leistet.

Bewusstsein (eine zweiwertige Welt)

Streng genommen kann in einer einwertigen Welt nicht von Bewusstsein im Sinne des Common Sense gesprochen werden,³⁷ da Subjekt und Objekt hier noch nicht getrennt sind, bestenfalls gibt es einen Fluss von Ereignissen, die differenzlos aneinandergekettet sind. Doch auch dieses Bild ist nur eine Metapher, denn sobald wir über ein Ereignis reden oder schreiben, vergegenständlichen wir es und bewegen uns damit bereits in der uns vertrauten zweiwertigen Welt, in der Beobachter und Beobachtetes, Erleben und Erfahrung auseinanderfallen. Das Bewusstsein erscheint jetzt, wie Günther in Referenz auf Hegel formuliert, als ein reflektiertes. Es taucht als ein Vorgang auf, der zwei Stellen erzeugt, nämlich »Reflexion-in-anderes« (den als Objekt erscheinenden Gegenstand) und »Reflexion-in-sich« (das sich seiner selbst bewusste Subjekt).³⁸

Für Günther drückt sich diese Relation mit Hegel bereits im Gottesverständnis aus: Die Götter der griechischen und römischen Antike – etwa Apollon und Zeus – waren selbst »dem Schicksal unterworfen«. Man kann sie sich in einer unpersönlichen Weise als »Blinde« vorstellen, die den verborgenen, jedoch bereits vorab festgelegten Pfaden der Welt zu folgen haben. »Dahingegen ist der christliche Gott der nicht bloß gewußte, sondern schlechthin sich wissende Gott und nicht bloß vorgestellte, sondern vielmehr absolut wirkliche Persönlichkeit.«³⁹ Die Metapher ›Gott‹ steht hierbei für absolute Subjektivität. Der Mensch kann sich als Ebenbild Gottes projizieren, um sich selbst als Subjekt zu enactieren. Ob der Mensch Gott erfunden hat oder Gott den Menschen, ist

des »Hier und Jetzt« zusammen, wodurch die Fähigkeit abhanden kommt, Kausalitäten zu erkennen. Aus der Perspektive ohnmächtiger Hilflosigkeit erscheint der helfende Arzt als einziger Retter. Die Arzt-Patient-Beziehung gewinnt hier eine unbestreitbare Eindeutigkeit, da die unmittelbar erlebte Krise gewissermaßen diktiert, was zu tun bzw. was hinzunehmen ist. Eine solche Dynamik tritt in der Regel bei einem schweren Unfall auf, kann aber auch bei einem Krankheitsbild wie dem akuten Herzinfarkt einsetzen.

37 Deswegen haben wir den Begriff ›Bewusstsein‹ voranstehend zunächst in Anführungsstriche gesetzt.

38 Günther (2021 [1957], S. 56).

39 Hegel (1986 [1817], § 147).

dabei mit Blick auf die logischen Verhältnisse unerheblich. Entscheidend ist das Ergebnis, ein Subjekt, das sich der Welt und ihrem Schicksal gegenüber sieht und doch nicht von ihr getrennt ist.

Die damit einsetzende Reflexion muss als Bruch in der Kontinuität einer vermeintlich einwertigen, mit sich selbst identischen Weltentfaltung erscheinen, gleichsam als ein Riss im Fluss des differenzlosen Erlebens. Die Welt, wie sie subjektiv erscheint und projiziert wird, ist von nun an nicht mehr identisch mit der sich faktisch realisierenden Welt.

Die imaginäre Verdoppelung der Welt in Form der »Differenz des Wirklichen und des Möglichen«⁴⁰ tritt nun ihrerseits in die Welt ein. Es entstehen jetzt Systeme, die eine Innenwelt haben, die für sie instruktiv wird, da die Außenwelt zu komplex ist, um ihre Eigengesetzlichkeiten zu kennen, geschweige denn sie verstehen zu können. Diese Systeme agieren nun entsprechend den imaginierten Möglichkeiten, also den Sinnhorizonten, die sich aus der Differenz zwischen den eigenen Projektionen (subjektiver Pol) und den diese zum Scheitern bringenden Widerständen der Welt ergeben (objektiver Pol).

Es ist dann gerade die Erfahrung der Nichtpassung, des Widerstands, des Scheiterns, die die Reflexion auf diese beiden Pole hervorbringt. Das bislang in einer einwertigen Welt agierende gottähnliche Wesen beginnt zu erfahren, dass es nicht allwissend und nicht allmächtig ist. Das sich bislang im Einklang mit seiner ernährenden Umwelt erlebende Baby muss nun den Unterschied zwischen sich und seiner Mitwelt, zwischen seinen eigenen Handlungsimpulsen und den Bewegungen seiner Mutter zu verarbeiten lernen. Der allmächtige Herrscher, dem bislang alle Untergebenen sklavenhaft jeden Wunsch von den Augen und Lippen abgelesen haben, beginnt nun mit dem Nein der Menschen konfrontiert zu werden, die ihm gegenüber Widerstand leisten.

Sofern das Baby, der Herrscher oder das gottähnliche Wesen in der Lage sind, sich in der einen oder anderen Weise an die erfahrene Spannung zwischen der eigenen Projektion und der gelebten Praxis zu erinnern, treten sie in eine zweitwertige Welt ein. Auf der einen Seite erscheint ihnen der objektive Verlauf der Welt, auf der anderen Seite ein Subjekt, das die Welt erkennt, sich dabei aber irren kann. Letzteres wird mit Blick auf die eigenen intentionalen Bewegungen daher unweigerlich auf Widerstand stoßen. Damit kommt als Differenz von Gewusstem und Nichtgewusstem unweigerlich der Begriff der *Information* ins Spiel. Man trifft auf eine Welt, deren Verlauf und Gesetzlichkeiten man nicht kennt, und beginnt entsprechend nach Wegen zu suchen, diese Wissenslücke zu schließen.

Bereits an dieser Stelle ist darauf hinzuweisen, dass die Kognition in einer zweitwertigen Welt eine Gedächtnisfunktion benötigt, die es erlaubt,

40 Luhmann (2008, S. 15).

die auf der Basis vergangener Interaktionen kondensierten Erwartungen mit der sich aktuell manifestierenden Gegenwart abzugleichen. Kybernetische Schleifen aus Reiz und Reaktion – auch wenn sie erlernt worden sind bzw. durch Lernen verändert werden können – müssen in diesem Sinne noch nicht als ein ›Bewusstsein‹ erscheinen. So wäre es denkbar, dass fühlende Akte bestehen, etwa im Sinne von Whiteheads »Elementarereignissen«. ⁴¹ Diese könnten zwar durchaus schon auf ein situatives Selbstgewahrsein hindeuten. Doch sofern sie sich nicht zu Erinnerungen verketteten, fielen sie sofort dem Vergessen anheim. Sie würden sich also noch nicht zu einem Bewusstsein kondensieren können, das sich immerfort reaktualisiert – also einem Bewusstsein, das mit Blick auf die eigene Geschichte der Welt als Subjekt entgegentreten kann (und damit nicht nur Teil der Welt ist, sondern zugleich eine Welt hat).

Das Bewusstsein zeichnet sich, wie Niklas Luhmann es mit wohlgewählten Worten ausdrückt, gerade dadurch aus, dass es »voran[schreitet]«, »indem es zurückblickt. Es operiert gleichsam mit dem Rücken zur Zukunft, nicht proflexiv, sondern reflexiv.« ⁴² Um es in Hinblick auf die phänomenologische Tradition im Anschluss an Edmund Husserl zu formulieren: Bewusstsein wird erst möglich als erinnerte und reaktualisierte Differenz. Auf der einen Seite steht die Aktualität der gegenwärtigen Wahrnehmung, auf der anderen Seite stehen die Erinnerung (Retention) sowie die hieraus abgeleitete Projektion einer Zukunft (Protention). ⁴³ Erst auf diese Weise kommt es zu einer Konstellation, in der für das System eine Differenz zwischen dem Gewussten (dem Erinnerten und als fortbestehend Projizierten) und dem Nichtgewussten (also dem, was für das System nicht berechenbar ist) erscheint, wobei diese Differenz im Bewusstseinsprozess wiederum einen Unterschied macht.

Wir gelangen hiermit zu einer kybernetischen Beschreibung, die Bewusstsein als eine *Relation* fasst – also nicht als eine Entität, nicht als ein inneres Seelenwesen oder einen Homunkulus, der im Gehirn sitzt und die Ergebnisse der Sinneswahrnehmungen betrachtet. Bewusstsein erscheint nun als ein Selbst- und *Weltverhältnis*. Da jedoch eine Relation weder ein Ding noch eine Essenz ist, weder Materie noch Energie, muss auch das Bewusstsein entsprechend der vorangehenden Beschreibung als substanzlos oder leer gelten. Als Schnittstelle einer Welt, die sich von Innen selbst beobachtet – und damit nicht von der Welt geschieden sein kann –, muss Bewusstsein gleichsam als »abstandlose Zusammenheit« ⁴⁴ erscheinen. Teil-der-Welt-Sein (objektive Perspektive) und Eine-Welt-Haben (subjektive Perspektive) müssen sowohl als verschieden wie auch als gleich erscheinen – als Einheit einer Differenz. Genau dies ist jedoch im

41 Whitehead (1987, S. 59 f.).

42 Luhmann (1995, S. 63).

43 Siehe zur Einführung in die Zeitverhältnisse Husserls: Schnell (2002).

44 Fuchs (in Vogd 2020, S. 210).

Alltagsbewusstsein nicht erfahrbar. Hier erscheint nur die Differenz als ein »naiv realistische[s] Selbstmissverständnis«⁴⁵. Es erscheint die Spaltung in Beobachter und Gegenstand – nicht jedoch die Einheit der Differenz.

Mit den Mitteln der klassischen Logik kann der logische Strukturreichtum dieser paradox anmutenden Beziehung nicht erfasst werden. Ebenso wenig lässt sich systematisch das Problem angehen, dass es mit Blick auf andere Wesen ja nicht nur ein Bewusstsein gibt, sondern derer viele.

Die klassische, auf Aristoteles zurückgehende Logik beschreibt diese Konstellation vielmehr folgendermaßen: Auf der einen Seite steht die objektive Welt als *positives Sein* (P), auf der anderen Seite das Subjekt als Reflexion des Seins. Der Positivität des Seins (P) steht die *Negativität* (N) des Subjekts gegenüber, das ja kein eigenes Sein habe, sondern dieses nur spiegele. Das im Subjekt entstehende Abbild kann korrekt oder inkorrekt sein. Entsprechend lassen sich hier zwei Wahrheitswerte zuordnen: Entweder ist das Abbild *wahr* (w) oder *falsch* (f), ein dritter Wert ist nicht möglich (*tertium non datur*). Beide Werte stehen über die Operation der Negation in einem Umkehrverhältnis. Wahr (w) ist äquivalent mit nicht falsch (~f).

Mit Blick auf die bereits eingeführte Unterscheidung zwischen Subjekt und Objekt gilt entsprechend: Das Subjekt ist nicht Objekt (und vice versa). Das Subjekt versucht sich an der Welt zu orientieren, sich ein Bild von ihr zu machen. Dies mag gelingen, wenn das Bild richtig (w), oder scheitern, wenn es aufgrund irriger Annahmen falsch ist (f). Die klassische Logik formuliert auf dieser Grundlage Elementaraussagen, die über eine begrenzte Anzahl logischer Operationen (zum Beispiel der ›und‹-, der ›oder‹- bzw. der ›weder noch‹-Operation) miteinander verknüpft werden können. Jede Elementaraussage hat gemäß dem Prinzip der Zweiwertigkeit exakt einen der beiden Wahrheitswerte wahr (w) oder falsch (f). Zudem gilt, dass zusammengesetzte Aussagen (etwa: ein Teilchen befindet sich am Ort *x* *und* bewegt sich mit der Geschwindigkeit *y*⁴⁶) eindeutig durch die Wahrheitswerte ihrer Teilaussagen bestimmt werden.

Dies entspricht genau dem naiven Realismus unserer alltäglichen Praxis: Ich schaue als innerer Beobachter auf eine Welt da draußen, die aus Objekten besteht, deren Eigenschaften sich kausal aus den Bestandteilen ableiten, aus denen sie zusammengesetzt sind.

45 Metzinger (1998).

46 Hier bereits in Anspielung auf die Heisenbergsche Unbestimmtheitsrelation, wo diese Beziehung nicht mehr gilt, weshalb die Quantentheorie nicht mehr als klassische Theorie anzusehen ist. Siehe ausführlich Kap. IV.2.

Mit der zweiwertigen Logik ergibt sich zudem ein Präferenzwert für das Objektive als das Positive und Wahre: Die objektive Wirklichkeit kann nicht falsch sein, da sie *irreflexiv* (I), also immer schon ohne Beobachtung, ohne Reflexion gegeben ist. Das Subjekt kann sich jedoch in seiner *Reflexion* (R) irren. Günther schreibt:

»Falsche Dinge kann es nicht geben, wohl aber falsche Bewußtseinsinhalte. Deshalb lehrt die klassische Tradition mit Recht, daß das Subjekt die Quelle allen Irrtums ist und daß Wahrheit erst dann in ihrer endgültigen Gestalt in Erscheinung tritt, wenn sie sich selbst im Medium der Objektivität zum Ausdruck gebracht hat. Einwertigkeit ist nur ein theoretischer Ausdruck für Unfehlbarkeit. Man kann mit den toten Dingen und mit Gott nicht argumentieren. Zweiwertige Existenz aber manifestiert sich in Handlungen, resp. Entscheidungen, und letztere können, wenn konfrontiert mit der unfehlbaren Positivität des Seins, wahr oder falsch sein.«⁴⁷

Wir gelangen hier zu einer dualistischen Metaphysik, in der einerseits »die ursprüngliche Thematik ›Sein‹ festgehalten wird« und andererseits »das Bewußtsein sich als Reflexion dieser Thematik von derselben absetzt«⁴⁸. Ersteres erscheint, um es nochmals zu wiederholen, als irreflexive (I) Positivität (P), Letzteres als Reflexivität (R), die als Negativität (N) zu begreifen ist, da sie entsprechend dem Prinzip des ausgeschlossenen Dritten (*tertium non datur*) nicht ihrerseits dem Sein zugehörig sein kann. Sein (P) und Nichtsein ($\sim P \equiv N$) stehen über die Negation (\sim) in einer Umtauschrelation.

Insofern aber nun Reflexion in der Welt ist, eröffnet sich die Möglichkeit, diese nicht nur als Reflexion auf das Sein (Objektpol) oder in Abstoßung von demselben als Selbstvergewisserung des Subjekts (Subjektpol) zu begreifen. Vielmehr kann in einer weiteren Operation die Unterscheidung bzw. Differenz zwischen Subjektpol und Objektpol selbst reflektiert werden. Mit Blick auf Hegel, »der diesen theoretischen Sachverhalt als erster mit durchdringender Klarheit gesehen« habe, formuliert Günther:

»Das volle theoretische Bewußtsein hat also
 1. einen Gegenstand (Sein, Anderes),
 2. weiß es sich im Gegensatz dazu, und
 3. ist es ein Wissen um den inversen Spannungszustand von Nicht-Ich und Ich.«⁴⁹

Auf der Höhe des reflexiven Bewusstseins beginnt der Mensch darüber nachzudenken, wie er unterscheidet und dabei Subjekte und Objekte

47 Günther (2021 [1957], S. 46).

48 Günther (2021 [1957], S. 56).

49 Günther (2021 [1957], S. 56).

konstruiert und auf diese Weise seine Welt aufbaut. In der Sprache der Kybernetik zweiter Ordnung wechselt er von der *Beobachtung erster Ordnung* zur Beobachtung des Beobachters, also zu einer *Beobachtung zweiter Ordnung*. Günther verwendet hierfür den Begriff *doppelte Reflexion* (D): Ich beginne zu reflektieren, wie kognitive Systeme (bzw. ich selbst als kognitives System) durch die Art und Weise ihres (meines) Unterscheidens Subjekt und Objekt in einer bestimmten Weise herstellen und dadurch ihre (und meine) Welt in einer bestimmten Weise anschneiden und aufbauen.⁵⁰

Die sich aus diesen Beobachtungsverhältnissen ergebenden Beziehungen lassen sich nicht mehr im Rahmen der klassischen Logik formulieren, denn in ihr führen selbstreferenzielle Aussagen wie ›der Beobachter beobachtet den Beobachter‹ zu Antinomien. Es bedarf einer transklassischen, mehrwertigen Logik, wie sie Gotthard Günther formuliert hat, um mit diesen Relationen umgehen zu können.

Polykontexturalität (mehrwertige Welten)

Versuchen wir uns im Folgenden der den meisten Lesern unvertrauten und damit nicht auf Anhieb verständlichen Konzeption der polykontexturalen Logik zu nähern. Hierbei ist es hilfreich, sich nochmals das obige Gedankenexperiment zu vergegenwärtigen, mit dem die Weltbeziehung verdeutlicht wurde, die einer einwertigen Logik folgt. In einer einwertigen Welt gibt es keinen Beobachter, sondern nur ein Sein, das unerschütterlich mit sich selbst eins ist. In einer Welt, die zweiwertig erlebt wird, gibt es demgegenüber einen Bruch zwischen Sein und Beobachter. Es erscheint eine zweiwertige Konstellation – Günther spricht hier von der *Kontextur* –, die sich in einer solchen Weise darstellt, dass beide Seiten durch die Operation der Negation logisch ineinander übergeführt werden können (die objektive Welt ist nicht das Subjekt). Wir haben jetzt einerseits die »Welt« als alles, »was der Fall« ist, wobei die »Gesamtheit der Tatsachen« – dies liegt in ihrer Natur – unabhängig von der Reflexion eines Beobachters für sich selbst steht.⁵¹ Aus der Perspektive des beobachtenden Subjekts kann ihr entsprechend eine *Seinsidentität* zugeschrieben werden.

Andererseits kann sich der Beobachter nun als *Reflexionsidentität*, als Ich auf sich selbst zurückwenden. Ich und Es, Subjekt und Objekt, Beobachter und Beobachtetes sind uns so sehr vertraut, sie liegen uns so nahe, dass wir darüber den reflexiven Prozess vergessen, der diese Unterscheidungen erst hervorbringt. Die vielfältigen komplexen physischen, psychischen und sozialen Vorgänge, die hieran beteiligt sind, können sich in

⁵⁰ Siehe von Foerster (1995).

⁵¹ Hier im Anklang an Wittgenstein (1974, Proposition 1 f.).

der jeweils erscheinenden Differenz nicht einmal ansatzweise selbst gewahr werden. Die elektromagnetischen Wellen, die unsere Sinnesorgane als Licht wahrnehmen, die Luft, die den Schall transportiert, die Sprache, die das Medium für unser Sinnverstehen darstellt, unser Gedächtnis, das Aktualität und Möglichkeit in Beziehung setzt, bleiben uns im Regelfall verborgen. Nur im Falle der selektiven Hinwendung zu einzelnen Aspekten wird uns für kurze Momente gewahr, dass all dies als Medium stillschweigend immer schon vorausgesetzt wird – als Bedingung der Möglichkeit, dass überhaupt etwas als Bewusstsein erscheint. Um es mit Wittgenstein zu pointieren:

»Aber das Auge siehst Du wirklich nicht. Und nichts am Gesichtsfeld lässt darauf schließen, daß es von einem Auge gesehen wird.«⁵²

Damit kommt aber neben Subjekt und Objekt eine dritte logische Position ins Spiel, nämlich der Prozess der Reflexion selbst, der die Unterscheidung aufspannt bzw. möglich macht. Die Systemtheoretikerin Elena Esposito formuliert in Referenz auf Günther:

»Die Operation der Negation setzt eine grundsätzlichere Operation bereits voraus: nämlich die Zäsur, die dazu führt, dass ein Objekt als unterschieden von anderem bezeichnet wird.«⁵³

Mit der *doppelten Reflexion* und der hiermit verbundenen *dritten* logischen Position (dementsprechend landen wir hiermit bei einer dreiwertigen Logik!) kann jetzt beobachtet und reflektiert werden, dass und wie zwischen Subjekt und Objekt unterschieden wird. Dies bringt logisch gesehen weitere Negationsmöglichkeiten mit sich: Es könnte nämlich *nicht* zwischen Subjekt und Objekt unterschieden werden (womit diese spezifische Differenz gar nicht erst erscheinen würde). Und: Es könnte auch anders unterschieden werden, also beispielsweise *nicht* zwischen *diesem* Subjekt und Objekt, sondern zwischen einem *anderen* Subjekt und den ihm erscheinenden Objekten. Damit landen wir unweigerlich in einer polykontexturalen Welt, in der es unterschiedliche subjektive Positionen geben kann, die auch logisch gesehen miteinander koexistieren können. Die Welt reduziert sich nicht mehr darauf, dass *Ich die Welt erlebe*, sondern beinhaltet jetzt auch, wie *Du die Welt erlebst* – und wie eine Vielzahl anderer Wesen (beispielsweise ein Schimpanse, eine Fledermaus, eine Ameise) ihrerseits die Welt in einer mir unzugänglichen Weise anschneiden können.

Auf den ersten Blick scheint sich die Sache einfach darzustellen. Man ergänzt eine weitere subjektive Position und erhält jetzt eben zwei

52 Wittgenstein (1974, Proposition 5.633; gesperrt im Original).

53 Esposito (1993, S. 105).

Konstellationen: ›Ich – Es‹ und ›Du – Es‹, die jeweils eine aus der klassischen Logik vertraute zweiwertige Konstellation aufspannen. Das Du wäre dann nichts anderes als ein Alter Ego, wäre also logisch gesehen nicht von mir als zweiwertig verfasstes Selbstverhältnis verschieden. Doch bei genauerer Betrachtung wird deutlich, dass noch eine dritte zweiwertige Kontextur benötigt wird, um die Verhältnisse zu beschreiben: die Ich-Du-Relation. Es zeigt sich nämlich, dass zwei unterschiedliche logische Verknüpfungen nötig sind, um den Strukturreichtum dieser Beziehung aufzuschließen: ›Ich *und* Du‹ bzw. ›Ich *oder* Du‹. Einerseits erscheint das Du wie ich als ein Subjekt, das sich einer objektiven Welt gegenübergestellt sieht. Andererseits erscheint mir das Du als eine körperliche Entität, die mir gegenübertritt. Aus dieser Perspektive tritt es als ein Körper in Erscheinung, der zwar kognitive Funktionen haben und sich bewegen kann, aber nichtsdestotrotz als materielles Objekt angesehen werden muss. Ich kann ihn aufschneiden und in seinem Inneren nachschauen, werde dabei jedoch keine Subjektivität und auch kein Bewusstsein finden können.

Beide Relationen zusammengenommen führen zu einer transklassischen Konstellation: *Das ›Du‹ ist zugleich Subjekt und Objekt. Das ›Du‹ hat damit also zwei, einander widersprechende Wahrheitswerte.* Dies ist jedoch im Rahmen einer zweiwertigen Logik nicht zulässig. Etwas kann nicht gleichzeitig sein Gegenteil sein. Erst in einer dreiwertigen Logik, die drei Positionen vorsieht, wird dies widerspruchsfrei fassbar.

Wie bereits angedeutet, führt der Weg zu einer mehrwertigen Logik über die Operation der Negation. Die zweiwertige Logik unterscheidet zwischen Sein und Nichtsein, Objekt und Subjekt, wobei die beiden Seiten der Unterscheidung über ein Umtauschverhältnis gegeben sind (*klassische Negation*).

Insofern der reflexive Prozess in einer transjunktionalen Operation seinen Unterscheidungsgebrauch reflektiert, kann das Treffen einer Unterscheidung seinerseits zum Thema der Reflexion werden, also ebenfalls negiert werden (*transklassische Negation*). Die hiermit einhergehende Unterscheidung kann ebenfalls reflektiert und anschließend negiert werden. Wir landen dabei zunächst bei einer vierwertigen Logik, die mit jeder eingenommenen weiteren Reflexionsposition zu einer beliebigen vielen Stellen beinhaltenden mehrwertigen Logik erweitert werden kann.

Die auch alltagspraktische Relevanz der mit den Reflexionsprozessen erfolgenden Negationen lässt sich gut anhand des Problems der passiven Negation verdeutlichen,⁵⁴ weshalb wir dieses aufgreifen möchten, um die Sinnhaftigkeit der vorangehenden logischen Überlegungen zu veranschaulichen.

54 Siehe zur aktiven und passiven Negation ausführlich Elster (1990).

Ich habe ein Problem. Vielleicht hat mich ein wichtiger Mensch verlassen oder ich versuche, mit dem Rauchen aufzuhören, und leide nun darunter. Um dem zu entgehen, versuche ich, nicht an das Problem zu denken, also versuche ich, *nicht* an die Person oder *nicht* an das Rauchen zu denken. Innerhalb der logischen Struktur des zweiwertigen Bewusstseins ist die Negation jedoch nur die andere Seite des Affirmierten. Wenn man versucht, *nicht* an einen Elefanten zu denken, denkt man an einen Elefanten, landet jedoch nicht in einem logischen Raum, in dem etwas anderes zum Thema wird. Ebenso stellt es sich mit der vermissten Person oder der ersehnten Zigarette dar. Um diesem Dilemma zu entkommen, müsste ich also nicht *nicht* an die problematische Sache denken, ich müsste also in eine Situation kommen, die ursprüngliche Negationsbeziehung zu transzendieren, was jedoch aus der Position des aktuellen Bewusstseins nicht möglich ist, das ja genau durch diese zweiwertige Beziehung gegeben ist. Man muss also in eine andere Position gestoßen werden, mit der das mit der ursprünglichen Reflexionsposition verbundene Problem verschwindet. Dies wird beispielsweise in der aus Coaching und Psychotherapie bekannten Tetralemma-Arbeit versucht.⁵⁵ Die Klientin wird hierbei gebeten, fünf unterschiedliche Positionen anzulaufen, um einerseits zu begreifen, wie sie durch ihren eigenen Unterscheidungsgebrauch die sie belastenden Probleme und Konflikte mit erzeugt, und um andererseits neue Lösungen zu finden, die salopp gesagt darin bestehen, die Sache anders zu sehen. Die jeweils einzunehmenden Positionen lauten:

1. *Das Eine*. Dies entspricht dem positiv Gegebenen (was für einen selbst der Fall ist)
2. *Das Andere*. Dies entspricht der Negation der vorherigen Position (die Sichtweise des anderen ist richtig und meine ist falsch).
3. *Beides – sowohl ... als auch*. Dies entspricht der Reflexion aus der dritten Position einer transklassischen Logik (sowohl ich als auch der andere haben Recht, auch wenn sich die Aussagen widersprechen).
4. *Keines von beidem – weder ... noch*. Dies entspricht der transklassischen Negation (beide Sichtweisen werden mit Verweis auf eine weitere Position zurückgewiesen).
5. *All dies nicht und selbst das nicht*. Es wird auf das seinslogisch und reflexiv nicht erreichbare Imaginäre des unmarkierten Raums jenseits aller Unterscheidungen verwiesen.

Das Tetralemma veranschaulicht die mehrwertige Logik in einem Bild, das in der Nebeneinanderschau die logischen Orte und ihre Beziehung abbildet. Die mit der Bildhaftigkeit verbundene Raummetaphorik erweckt jedoch fälschlicherweise den Eindruck, als ob die unterschiedlichen

55 Siehe Varga von Kibéd und Sparrer (2000) sowie Kleve (2011).

Kontexturen durch ein hinreichend geschultes Bewusstsein gleichzeitig ange-
laufen und erfahren werden könnten.
Es wird gleichsam suggeriert, dass es
eine Art Superbewusstsein geben könnte,
mit dem sich nicht nur die eigene und
die Perspektive eines anderen Menschen
einnehmen ließe, sondern zugleich auch
noch alle anderen denkbaren Perspektiven
und überdies sogar die Formlosigkeit
der Leere des undefinierten Raums.

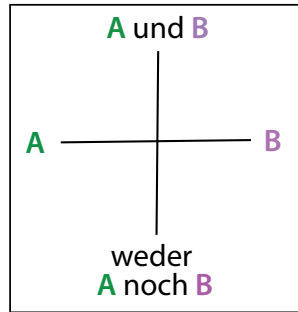


Abb. 2:
Darstellung des Tetralemmas

Dies ist freilich nicht möglich. Das Bewusstsein wird immer nur einen Gegenstand haben können, der ihm als Inhalt gegenübertritt und von dem es sich selbstidentifizierend (erneut) abgrenzen kann. Wie komplex die Abbildung auch sein mag, sie bleibt das Erscheinen *einer* Differenz – und damit gefangen in der Zweiwertigkeit eines Bewusstseins, das immer nur *etwas* als Gegenstand haben kann. Wie bei den Kippbildern der Gestaltpsychologie kann nur das eine oder das andere gesehen werden, nicht jedoch beide Gestalten und der formlose Bereich des noch nicht Wahrgenommenen.

Wenn wir mit Gotthard Günther die logischen Positionen, die mit einer mehrwertigen Logik möglich sind, konsequent durchdenken, stoßen wir jedoch auf einen merkwürdigen Befund: Wir gelangen jetzt auch zu Kontexturen mit Umtauschverhältnissen, in denen »keine absolute Scheidung zwischen ›Ich‹ und ›Es‹« mehr stattfindet. Die Reflexionsbeziehung würde damit zu einem Bewusstsein führen müssen, »das keinen unmittelbaren Gegenstand« mehr »besäße«, das heißt, »in ihm könnte jener fundamentale Negationsvorgang, in dem das erlebende Subjekt sich von allen seinen Inhalten distanziert und sagt: ›Das bin ich nicht‹, niemals zustande kommen. In anderen Worten: Das Umtauschverhältnis« stellt »keinen subjektiv erlebbaren und als Innerlichkeit verstehbaren Ich-Welt-Prozeß dar«⁵⁶.

Um dies in anderen Worten auszudrücken: In einer polykontexturalen Welt passiert an vielen Positionen etwas, was indirekt beeinflussen mag, was mir als Bewusstseinsgegenstand erscheint. Doch dies wird mir möglicherweise subjektiv niemals zugänglich werden. So gilt bereits für die Begegnung mit der Du-Subjektivität, dass in der symbolischen und reflexiven Aneignung der Du-Position in einer therapeutischen Übung wie der zuvor geschilderten Tetralemma-Arbeit das Gegenüber eben nur als Objekt begriffen werden kann, nicht jedoch als fremde Subjektivität. Es bleibt ein vorgestelltes Modell, eine objektivierende Theorie vom

56 Günther (1958, S. 21).

Du, eröffnet jedoch keinen Zugang zur Innerlichkeit einer anderen Subjektivität. Wir stehen unweigerlich vor einem kontextuellen Abbruch – denn der Graben zwischen der zweiwertigen Monokontextur meines Bewusstseins und dem logischen Postulat einer mehrwertigen Welt ist unüberbrückbar:

»Leben und kontextueller Abbruch im Wirklichen sind nur zwei verschiedene Ausdrücke für denselben Sachverhalt. Was jenseits des Abbruchs liegt, ist schlechterdings unzugänglich. Was hiermit gemeint ist, muß jedem sofort deutlich werden, wenn wir auf eine ganz alltägliche Erfahrung hinweisen. Für jedes erlebende Ich ist die innerste Privatheit der Du-Subjektivität ein ebenso unzugänglicher Raum wie die mythologischen Dimensionen, in denen die himmlischen Heerscharen schweben. In beiden Fällen stehen wir am Rande eines Kontexturabbruches, der in keinem Fall größer oder geringer ist als in dem anderen.«⁵⁷

Es ist deshalb folgerichtig, wenn Günther in Bezug auf den Wechsel von der Kontextur des eigenen Bewusstseins zum fremdpsychischen Du das dabei erscheinende Austauschverhältnis zwischen Subjekt und Objekt als *Transzendentalidentität* bezeichnet. Der Begriff Transzendenz verweist auf einen Gegenstand, der jenseits der möglichen Erfahrung liegt. Die Subjektivität des anderen liegt nicht im Bereich meiner Welt, sie kann nicht als eine mir zugängliche Differenz erscheinen:

»Damit enthüllt sich uns auch die tiefere transzendente Bedeutung der Kybernetik. Subjektivität überhaupt ist uns in unserer Erfahrung in zwei Gestalten gegeben. Erstens als eigenes seelisches Leben und zweitens als Fremdseelisches. Zugänglich aber ist uns jene reine Subjektivität nur in der intimen Privatheit des eigenen Ichs. Das andere Ich jedoch ist uns so fern und in seiner ihm allein eigenen Innerlichkeit so unerreichbar wie das Jenseits selbst. Jene seelische Distanz zwischen Ich und Du wird von uns zwar als Faktum erlebt, aber ihr Wesen bleibt unverstanden.«⁵⁸

Die logische Konstellation dieser dreiwertigen Welt lässt sich auch durch das folgende Schaubild charakterisieren:

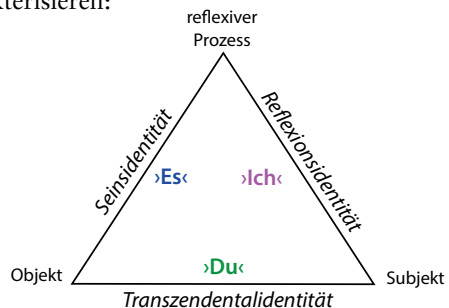


Abb. 3: Dreiwertige Welt aus ICH, DU und ES.⁵⁸

57 Günther (1975, S. 61 f.).

58 Günther (2021 [1957], S. 37), Abbildung nach Günther (2021, S. 22).

Je nach Reflexionsperspektive muss das Du als etwas anderes gelten: Von der objektiven Seite her gesehen erscheint es als Körperlichkeit, als eine strukturdeterminierte Einheit. Von der subjektiven Seite her reflektiert erscheint es als die Innenseite der Schnittstelle einer zweiwertigen Welt, die wie auch ich in Beobachter (Ich) und Beobachtetes (Es) zerfällt.

*Notwendige Blindheit – nur wer nicht weiß,
muss sich entscheiden und handeln*

Die objektive und die subjektive Perspektive gehen zudem mit einem Komplexitätsgefälle einher in Hinblick auf das, was prinzipiell gewusst werden kann. Von außen gesehen lässt sich zumindest theoretisch davon ausgehen, dass das Verhalten eines Menschen vollständig von seinen hirnorganischen Prozessen, den bisher erfolgten Interaktionen mit der Umwelt und den daraus erwachsenden kognitiven Prozessen bestimmt ist. Selbst wenn man nicht alle Details kennt, kann sein Körper also als ein *vollständig spezifiziertes System* angesehen werden.

Von innen her gesehen erscheint jedoch die Differenz zwischen dem, was gewusst und zugänglich, und dem, was nicht zugänglich ist und bestenfalls erahnt oder diskursiv angeeignet werden kann. Denn der »Standpunkt« eines »fremden Beobachters«, so Günther, ist nicht »im Geringssten für das subjektive Ich verbindlich, das sich in Konfrontation mit dem eigenen Leibe erlebt und seine Beziehungen zu demselben zu verstehen sucht. Für dieses Ich bleiben die physischen Bedingungen seines Daseins in der Welt bis zum jüngsten Gericht *unterspezifiziert*.«⁵⁹

Damit ist es aus kybernetischer Perspektive gerade »diese nie aufhebbare Differenz von Unterspezifikation und Vollspezifikation, die für das subjektive Erleben zweier voneinander unterschiedener Personen« verantwortlich ist. Das »technische Korrelat« der Unterscheidung von Ich- und Du-Subjektivität und der von Subjektivität und Objektivität« besteht ja gerade in der nicht aufhebbaren »Differenz von Unterspezifikation und Vollspezifikation«⁶⁰.

Da aus der Innenperspektive eben nicht alles gewusst werden kann, bleibt nichts anderes übrig, als sich an Information, also an Signalen und Unterschieden, die in ihrer grobkörnigen Form anzeigen, dass etwas der Fall sein könnte, zu orientieren, ohne damit jedoch hundertprozentig wissen zu können, was wirklich der Fall ist oder in Zukunft der Fall sein wird, da es sich eben nur um einen grobkörnigen Anschnitt der Welt handelt.

Sobald es Informationen sammelnde und verarbeitende Systeme gibt, wird unweigerlich eine Differenz zwischen unterspezifizierter,

59 Günther (2021 [1957], S. 111; kursiv durch uns).

60 Günther (2021 [1957], S. 111).

grobkörniger Information und feinkörniger Welt entstehen. Die Komplexität der Welt kann in ihrer Feinkörnigkeit niemals von Innen her begriffen werden – allein schon deshalb, weil ein Teil das Ganze niemals vollständig repräsentieren kann. »Die logischen Geltungsbereiche der Reflexionsprozesse im subjektiven Subjekt (Ich) und im objektiven Subjekt (Du) sind auf Grund der ontologischen Differenz von unter- und vollspezifiziertem System schlechterdings nicht zur Deckung zu bringen.«⁶¹ Diese Differenz ist übrigens homolog mit der Unterscheidung von Oberflächenstruktur und Tiefenstruktur, die in der Einleitung dieses Buches mit Blick auf die Schnittstellen der Interaktion mit Computern eingeführt wurde. Das Nichtwissen um die kausalen Prozesse einer strukturdeterminierten Maschine verhindert nicht, dass die symbolischen und semantischen Operationen an der Oberfläche instruktiv werden und hier einen Unterschied eigener Qualität machen.

Das Teil weiß nicht um das Ganze. Hiermit gilt aber auch das Umkehrverhältnis: Das Ganze repräsentiert nicht die Teile. Oder um den berühmten Satz von Erwin Schrödinger aufzugreifen:

»Bestmögliches Wissen um ein Ganzes schließt nicht notwendig das Gleiche für seine Teile ein.«⁶²

Aus der Perspektive einer informationsverarbeitenden Einheit, die aufgrund ihres kognitiven Systems eine Welt hat, muss die Beziehung zwischen sich und der Welt also notwendigerweise unterbestimmt erscheinen.

Doch damit gewinnt sie gerade aufgrund ihres Nichtwissens die Freiheitsgrade, sich auf Basis von Reflexion selbst eine imaginäre Welt zu schaffen, die für ihr Erleben und Handeln instruktiv wird. Wenn ich den Weltverlauf nicht wissen kann, meine Zukunft für mich ungewiss ist und ich nicht einmal vorhersehen kann, wie ein Mitmensch in meiner unmittelbaren Nähe auf mich reagieren wird, woran kann ich mich dann orientieren? Die einzige praktikable Antwort lautet: an meinen Vorstellungen und Projektionen, an dem, was ich als meine mögliche Zukunft konstruiert habe.

Sobald sich die Reflexionstätigkeit auf sich selbst wendet, hat das Sein, so Günther, jedoch »keine von Ewigkeit her vorbestimmten Eigenschaften mehr. Das Bewusstsein kann ihm also alle die aufzwingen, die es will.«⁶³ Die Wirklichkeit erscheint mir von nun an kontingent, »also weder notwendig noch unmöglich«, das heißt als etwas, was »so, wie es ist (war, sein wird), sein kann, aber auch anders möglich ist«, formuliert im gleichen Sinne Luhmann.⁶⁴ Die Reflexion-in-sich (R) – das sich

61 Günther (2021 [1957], S. 112; kursiv im Original).

62 Schrödinger (1935, S. 827).

63 Günther (2021 [1957], S. 60).

64 Luhmann (1984, S. 152).

selbst gewahre Denken und Fantasieren – kann sich nun in der doppelten Reflexion (D) als ein Denken und Fantasieren, das sich auf ein Denken und Fantasieren bezieht, in sich spiegeln und projizieren:

Das »Objektive hat kein bestimmtes Sosein mehr, in dem es sich dem Bewußtsein vermittels unabänderlicher primordialer Kategorien entgegengesetzt. Alle Kategorien sind ja »Reflexion« und als solche längst aufgelöst worden.«⁶⁵

Insofern die Welt diesem Prozess keinen Widerstand entgegengesetzt, kann er sich in beliebigen Gebilden nahezu endlos perpetuieren. Vertraut und ersichtlich ist uns dies insbesondere in unseren Träumen:

»Ein solcher willkürlicher Umgang des Bewußtseins mit einem Gegenstand ist uns allen seit jeher bekannt. Es ist der Traum. Der Wachtraum sowohl wie der Traum des Schlafenden. In dieser Bewußtseinslage hat sich die Reflexion ganz aus der »realen« Welt in sich selbst zurückgezogen. Sie reflektiert nur noch ihre *eigenen* Bilder und ist deshalb in der Lage, sie mit jeder beliebigen, resp. *gewünschten* Eigenschaft auszustatten. Diese Reflexionsphantasie ist der nach innen gewandte Wille. Was wir aus dem Traum lernen können, ist die sehr wesentliche Tatsache, daß die Bewußtseinsgegenstände der zweiten Reflexion dem Zugriff der wollenden Phantasie völlig offenliegen und ihnen gewünschte Eigenschaften aufgezwungen werden können.«⁶⁶

Die sich hier anhand der Traumgegenstände reproduzierende Subjektivität eines Ichs (oder des Du einer anderen, fremden Subjektivität) hat jedoch keine Substanz in sich. Sie beruht nicht auf einer Bewusstseinsessenz im Sinne eines unabhängigen inneren Seelenwesens. Auch sie stellt nichts anderes dar als das Erscheinen einer Differenz, die sich den Bedingungen ihrer Möglichkeit nicht selbst gewahr sein kann. Auch sie ist nichts anderes als Ausdruck und Produkt der Perspektivendivergenzen einer Welt, die von innen gesehen unterspezifiziert, nicht voll bestimmt ist, jedoch von außen betrachtet – vom Gottesstandpunkt einer einwertigen Welt her gesehen – vollspezifiziert wäre. Hier kann die »Welt«, um mit Wittgenstein zu sprechen, nichts anderes sein, als »alles, was der Fall ist«⁶⁷. Erst mit der Differenz von Innen und Außen (der subjektiven und der objektiven Perspektive) kommt jedoch eine »dritte metaphysische Komponente« unserer »phänomenalen Wirklichkeit« ins Spiel: die »*Information*«⁶⁸.

Information kann verstanden werden als ein Prozess bzw. als ein Vorgang, der in Hinblick auf das, was mittels Reflexion gewusst werden

65 Günther (2021 [1957], S. 60).

66 Günther (2021 [1957], S. 60; kursiv im Original).

67 Wittgenstein (1974, Proposition 1 f.).

68 Günther (2021 [1957], S. 20).

kann und was in Anbetracht einer unvorstellbar komplexen Welt unbestimmt bleiben muss, einen Unterschied macht. Als »ein Unterschied, der einen Unterschied macht«⁶⁹, stellt sie eine doppelte, aneinander gekoppelte Relation dar. Information ist damit weder materiell noch subjektiv, aber auch nicht unabhängig von materiellen und energetischen Flüssen sowie den hieraus erwachsenden ›subjektiven‹ Perspektiven zu verstehen.

Wenn man nicht wissen kann, was der Fall ist, muss man sich eine Welt erschaffen. Genau dies ist die Situation lebender Systeme, die ihre Existenz einer selektiven Blindheit verdanken, das heißt beständig der Welt Informationsgewinne abzurufen versuchen, um weiterhin bestehen zu können. Wenn man nicht weiß, was der Fall ist, bleibt nur übrig, sich aus den gegebenen Daten eine Welt zu konstruieren. Die in diesem Zuge entstehenden imaginären Werte werden dabei auf zwei unterschiedliche Weisen instruktiv:

1. zum einen nach *innen* als Reflexionsüberschuss im *Erleben*, also in Fantasie und Denken;
2. zum anderen nach *außen* im *Handeln*, also in Unterschieden, die sich durch das Wirken des Organismus in der Welt selbst materialisieren.⁷⁰

Handeln und Erleben erscheinen damit ebenfalls als zwei Seiten einer Medaille. Sie erscheinen als Ausdruck eines Reflexionsüberschusses, der sich aus der Unterspezifikation kognitiver Systeme ergibt. Da diese niemals vollständig wissen können, was der Fall ist, bleibt ihnen nichts anderes übrig, als Imaginäre zu erzeugen, um sich daran zu orientieren. Nach innen gerichtet erscheint dies als Erleben (das heißt als Konstruktion einer eigenen Wirklichkeit), nach außen gerichtet als Handeln (das heißt als Intervention in Welt). Auf sich zurückgewendet führt der Reflexionsprozess mit Blick auf die Spannung von Subjekt und Sein also nicht nur in die Innerlichkeit, also zum Aufbau und zur Vervielfältigung von Welt in der Fantasie, sondern mit der Willensbewegung auch zur Gestaltung von Welt. Die doppelte Reflexion kann auf der Subjektseite (Reflexion-in-sich), aber ebenso auf der Objektseite (Reflexion-in-anderes) ansetzen.

69 Bateson (1987).

70 Zur Erinnerung: Die Begriffe ›innen‹ und ›außen‹ gelten dabei als logisch inkommensurable Räume zweier divergierender Perspektiven (*Eine-Welt-Haben* vs. *Teil-der-Welt-Sein*). Es wäre jedoch falsch, diese räumliche Metaphorik wörtlich zu nehmen, etwa im Sinne: In meinem Kopf befindet sich das Ich und um mich herum sind die Tatsachen meiner Welt. Innen und außen befinden sich vielmehr in einer »abstandlosen Zusammenheit« (Fuchs in Vogd 2020, S. 210). Ebenso wenig darf selbstredend die phänomenologische Unterscheidung von *Körper-Haben* und *Leib-Sein* räumlich verstanden werden.

Auch an dieser Stelle ist es nochmals hilfreich, sich bewusst zu machen, dass wir von einem Reflexionsprozess sprechen, also wiederum eine *Relation* in einem logischen Raum bezeichnen. Da ist kein inneres Seelenwesen, das etwas erlebt oder will, sondern eine Reflexionsbeziehung, die sich je nach betrachteter Konstellation als Erleben oder Wollen deuten lässt, sich jedoch einer Totalität verdankt, zu der dann nolens volens auch der ins Gegenständliche ausgelagerte Reflexionsrest gehört:

»In anderen Worten: Es gibt eine Gestalt der Reflexion, die weder im Ich noch im Du lokalisiert ist, sondern die erst im Es, d. h. im Gegenstand, auftritt. Das Selbstbewußtsein, das sich bloß in der Subjektivität, also ausschließlich im Ich und im Du manifestiert, bleibt fragmentarisch. Es ist nicht total. Es bleibt relativ, weil es von einem nicht bewältigten Reflexionsrest abhängt. Einem Reflexionsrest, der sich innerhalb der Spannweite von Ich und Du nicht realisieren kann und der statt dessen das gegenständliche Objekt als Projektions- und Realitätsbasis braucht.«⁷¹

In diesem Sinne bedürfen Fantasie, Denken, Bewusstsein und Subjektivität ebenso der Verkörperung wie das Wollen und Handeln. Ob wir von einem einfachen Regelkreis sprechen oder von den komplexen Prozessen der Bewusstseinstätigkeit, in beiden Fällen gilt mit Maturana und Varela:

»Jedes Tun ist Erkennen, und jedes Erkennen ist Tun.«⁷²

Mit Blick auf einen Organismus mit einem Nervensystem lässt sich formulieren: Es gibt kein Sehen, Schmecken, Hören, Fühlen oder Denken, ohne dass sich physiologische Prozesse verändern, die dann mit einer Nervenaktivität einhergehen. Und umgekehrt: Es gibt keine Nervenaktivität, die nicht an irgendeiner Stelle im Körper einen Unterschied macht und in der Folge den Zustand des Organismus verändert. Es gibt keinen Bewusstseinsprozess, der nicht zugleich in der einen oder anderen Weise mit einer Verkörperung einhergeht, die im Bewusstsein nicht zugleich selbst abgebildet oder reflektiert werden kann. Die Reflexion, der Gedanke, die Wahrnehmung oder das Wollen können nicht gleichzeitig reflektieren, denken, wahrnehmen oder antizipieren, was sie mit der Welt machen. Doch sie haben keine andere Wahl, als etwas mit der Welt zu machen. Aus kybernetischer Perspektive lässt sich entsprechend mit Günther formulieren:

»[D]ie dritte, das System des Selbstbewußtseins vollendende Bewußtseinslage einer Subjektivität, die weder Ich noch ontologisch gegebenes Du ist, existiert nur als unerledigter Reflexionsrest in dem fragmentarischen System, das wir menschliches Selbstbewußtsein nennen. Jener

71 Günther (2021 [1957], S. 65).

72 Maturana und Varela (1987, S. 32).

Reflexionsrest bleibt durch den Prozeß des reflexiven Denkens unbewältigt, weil er sich eben nicht total in subjektive Reflexivität auflösen kann. Er ist jenes Andere, jenes Moment der Irreflexivität, um das der Strom des Bewußtseins wie um einen Fremdkörper spült, ohne ihn zu durchdringen und transparent machen zu können.

Inhalte aber, die das Bewußtsein nicht durch den Reflexionsprozeß bewältigen und auflösen kann, müssen eben auf eine andere Weise erledigt werden. Aber die einzige andere Methode, die neben der Reflexion auf die eigene Reflexion dem Ich zur Aneignung seiner Inhalte zur Verfügung steht, ist die Handlung, d. h. die Rückprojektion jenes irreflexiven Restbestandes in die Außenwelt.⁷³

Das Verhältnis von Innen und Außen, Subjekt und Objekt muss als eine Art Totalität begriffen werden, die sich jedoch nicht als Einheit, sondern nur als Differenz erfahren lässt. Dies erinnert an ein Möbiusband, bei dem die Innenseite auf die Außenseite führt und umgekehrt. So oft man das Band auch auf der Mittellinie durchschneiden mag, die verwickelte Topologie ändert sich nicht. Um in der Metaphorik zu bleiben: Welche Reflexionstiefe bzw. Weltkomplexität man auch generiert, die paradoxe Topologie (innen = außen) verschwindet nicht, denn an jeder Stelle auf dem Band lässt sich weiterhin distinkt zwischen zwei Seiten (innen ≠ außen) unterscheiden. Die »Exowelt« lässt sich nur aus der »Endowelt« erkunden. Doch die mit der Beobachtung entstandenen Schnitte tauchen erneut in der »Exowelt« auf.⁷⁴



Abb. 4: Möbiusband (Midjourney, prompted by Jonathan Harth)

Wenn kybernetische Maschinen die Exowelt aus der Endoperspektive zu erkunden beginnen

Fassen wir zusammen: Die doppelte Reflexion richtet sich gemäß einer Kybernetik zweiter Ordnung auf die Unterscheidung von Subjekt und Objekt. Im Sinne einer transjunktionalen Operation wird damit der Unterscheidungsgebrauch selbst thematisiert und kontingent gesetzt. Wie wir idealtypisch zunächst anhand der Du-Position aufgezeigt haben, bringt dies mit sich, dass das Du (als anderes Ich) sowohl als Objekt (als Körper) wie

73 Günther (2021 [1957], S. 58).

74 Die Begriffe »Exo-« und »Endowelt« wurden von Otto Rössler (1992) geprägt.

auch als Subjekt (als transzendente Innerlichkeit) reflektiert werden kann und damit seinerseits als unbestimmt erscheint.

Dann haben wir mit Günther das Augenmerk auf den eigenen Bewusstseinsprozess gelenkt. Wir sind dabei zunächst bei der Innerlichkeit des Träumens gelandet und von da aus der Notwendigkeit begegnet, sich im Handeln veräußerlichen zu müssen. Hiermit einhergehend können wir eine Metaphysik hinter uns lassen, die das »absolute Sein« und die »subjekthafte Reflexion« eines vermeintlichen (jedoch bei näherer Untersuchung nicht auffindbaren) inneren Seelenwesens strikt voneinander trennt. Dies führt uns schließlich zu einem Reflexionsprozess, der sich in beliebige Komplexität hinaufschrauben kann und dabei für jede spezifische Beobachterposition eine distinkte Metaphysik produziert. Im Sinne einer »fungierenden Ontologie«⁷⁵ wird vor dem Hintergrund der in diesem Kapitel beschriebenen Prozesse nun etwas ›Bestimmtes‹ der Fall sein, das als *Differenz erscheint* und damit einen weiteren Unterschied macht.

Schließlich wurde auf dieser Basis ein kybernetisches Verständnis des Begriffs der Information deutlich: Information verweist auf den Versuch eines unterspezifizierten Systems, die Welt zu spezifizieren und sich auf diese Weise zu reproduzieren. Der Preis, den diese unterspezifizierten Systeme dafür zu zahlen haben, besteht darin, niemals die Ganzheit sehen zu können, also selektiv blind sein zu müssen. Solchermaßen konstituiert, bleibt ihnen nichts anderes übrig, als weitere Informationen zu sammeln und zu verarbeiten, das heißt einzelnen Daten auf Basis grobkörniger Parameter (die niemals den feinkörnigen Welt dynamiken entsprechen können) einen Wert zu geben, um hierdurch die eigenen Aktionen als sinnhaft erscheinen zu lassen.

Für Gotthard Günther besteht kein prinzipielles Hindernis, kybernetische Maschinen zu konstruieren, die auf genau dieser Basis arbeiten. Hiermit stellt sich für uns somit die Frage nach dem Bewusstsein der Maschinen. Es sprechen keine prinzipiellen Gründe dagegen, dass eine hinreichend entwickelte künstliche Intelligenz nicht ihrerseits die Exowelt aus der Endoperspektive zu erkunden beginnt und dabei auf jene Aporien stoßen wird, die das menschliche Bewusstsein auszeichnen: in einem polykontexturalen Universum in einer zweiwertigen Welt gefangen zu sein.

75 Fuchs (2004, S. 11).

I.3 GOTTHARD GÜNTHERS POLYKONTEXTURALER BLICK AUF DAS BEWUSSTSEIN DER MASCHINEN

»The greatest social consequence of the Darwinian revolution was the grudging acceptance by humans that humans were random descendants of monkeys, neither perfect nor engineered. The greatest social consequence of neo-biological civilization will be the grudging acceptance by humans that humans are the random ancestors of machines, and that as machines we can be engineered ourselves. I'd like to condense that further: Natural evolution insists that we are apes; artificial evolution insists that we are machines with an attitude.«

Kevin Kelly⁷⁶

»Eine weitere Ganzheit ist das Empfinden eines Komplexes von Empfindungen, zu dem auch ihre spezifischen Elemente der Identität und des Kontrasts gehören. Dieser Integrationsprozeß des Empfindens dauert, bis die konkrete Einheit des Empfindens erreicht ist. In dieser konkreten Einheit ist alle Unbestimmtheit hinsichtlich der Realisierung von Möglichkeiten ausgeschaltet worden. Die vielen Einzelwesen des Universums, einschließlich derjenigen, die in der Konkretisierung selbst entstanden sind, finden ihre jeweiligen Rollen in der abschließenden Einheit. [...] Die Konkretisierung weist in jeder ihrer früheren Phasen schiere Unbestimmtheit hinsichtlich der Nexus zwischen ihren zahlreichen Bestandteilen auf.«

Alfred North Whitehead⁷⁷

Wie verändern sich unser Weltbild und unsere metaphysische Verortung, wenn wir Maschinen nach unserem Ebenbild bauen, wenn Menschen also Apparate entwickeln, die auf intelligente Weise wahrnehmen und agieren können? Verändert sich damit nicht auch unweigerlich unsere Vorstellung vom Verhältnis unserer innerseelischen Vorgänge und der (vermeintlich) unbeseelten Materie?

Würde sich nicht auch das Koordinatensystem unserer Spiritualität verschieben, wenn sich mit der Erschaffung künstlicher Intelligenzen die Unterscheidung zwischen geistigen und materiellen Prozessen zu verflüssigen beginnt?

Mit seinem bereits im Jahr 1957 erschienenen Buch *Das Bewußtsein der Maschinen* hat sich Gotthard Günther seiner Zeit weit voraus

⁷⁶ Kelly (1992, S. 55).

⁷⁷ Whitehead (1987, S. 391).

tiefgründig mit diesen Fragen beschäftigt. Nicht zuletzt führt Günther vor, wie die künftigen kybernetischen Maschinen uns Menschen in produktiver Weise den Spiegel vorhalten können: Sie demonstrieren uns, dass wir selbst informationsverarbeitende Systeme sind, die in Beziehung zu anderen Organismen treten, um auf Basis der hieraus erwachsenden Reflexionsbeziehung eine Identität zu gewinnen.

Wir treffen dabei auf den unüberbrückbar scheinenden Graben zwischen Subjekt und Objekt und die damit zusammenhängende Differenz zwischen unserem eigenen Empfinden und Erleben und den wissenschaftlichen Beschreibungen unserer kognitiven Prozesse.

Auf der einen Seite sind wir kaum anders als die von uns konstruierten intelligenten Maschinen, die mehr und mehr unsere Gesellschaft bevölkern. Sofern wir mit Humberto R. Maturana und Francisco J. Varela unseren Körper und unser Gehirn als eine strukturdeterminierte Einheit betrachten,⁷⁸ müssen wir selbst als ein organischer Computer gelten. Wir selbst sind eine – wenngleich komplexe – nichttriviale lernfähige Maschine, die aufgrund von Erfahrungen ihre Input-Output-Relationen ändert.

Auf der anderen Seite erscheinen wir jedoch als ein Bewusstsein, als ein Beobachter, der mit einer unhintergehbaren Differenz konfrontiert ist. Uns erscheint die Welt dichotom, als Differenz zwischen Subjekt und Objekt, als Differenz zwischen ›Leib sein‹ und ›Körper haben‹, als Differenz zwischen Wirklichkeit und Möglichkeit, als Differenz zwischen Erleben und Wollen, als Differenz zwischen Welt und Selbst. Wir *sind in der Welt* und *haben eine Welt*.

*Das Geheimnis des Beobachters – zugleich Welt sein
und eine Welt haben*

Gerade mit der Formulierung der letztgenannten Differenz wird klar, dass wir keineswegs zwei getrennte Welten haben – auf der einen Seite die Sphäre der Materie, auf der anderen die Sphäre des Geistigen oder Seelischen. Vielmehr treffen wir auf Relationen bzw. Reflexionsbeziehungen, die jeweils zwei Seiten zugleich trennen wie auch verbinden – wobei sich die Grenze von dem, was jeweils auf der einen oder anderen Seite steht, verschieben wird, je nachdem, welche Position gerade eingenommen wird. Beobachtung ist ein Prozess, ein Vorgang, mittels dem unterschieden und bezeichnet wird und durch den Schnitte entstehen, die auf das Verbindende verweisen.⁷⁹ Nicht mehr und nicht weniger.

⁷⁸ Maturana und Varela (1985).

⁷⁹ Karen Barad hat in Bezug auf die Schnitte, die eine Beobachtung macht, die schöne Formulierung gefunden: »Cuts are matters of cutting together-apart.« (Barad 2012, S. 19)

Der Beobachter steht damit nicht außerhalb der Welt, sondern emergiert, wie Fritz Simon in Anklang an Spencer Brown formuliert, infolge einer Beobachtungsoperation, die innerhalb der Welt vollzogen wird:

»Als Beobachter soll definiert sein, wer oder was (das heißt, es muss sich dabei nicht um einen Menschen oder ein Lebewesen handeln) einen spezifischen Typus von Operation vollzieht: *beobachten*. [...] Unter Beobachten soll eine Operation verstanden werden, die durch die Koppelung zweier anderer Operationen entsteht: unterscheiden und bezeichnen.«⁸⁰

Wir begegnen dieser Bewegung schon immer dann, wenn wir versuchen, durch Selbstbeobachtung unserer eigenen Innerlichkeit zu begegnen. In den buddhistischen Lehren wird diese Methode umgekehrt genutzt, um zu zeigen, dass es ein substanzielles Ich im Sinne eines inneren Seelenwesens nicht gibt.⁸¹ Dies lässt sich leicht illustrieren, indem man folgende Übung versucht:

Betrachte deinen Körper. Wenn du deinen Körper betrachten kannst, dann kann dieser nicht das Ich, nicht der Beobachter sein.

Betrachte deine Empfindungen und Gefühle. Wenn du deine Empfindungen und Gefühle betrachten kannst, können sie nicht dein Ich oder der Beobachter sein.

Betrachte deine Gedanken. Wenn du deine Gedanken beobachten kannst, können sie nicht dein Ich oder der Beobachter sein.

Beobachte den gegenstandslosen Beobachter bzw. das Gewahrsein, in dem dies alles erscheint. Falls dies gelingt, offenbart dies, dass also auch dies nicht ein Ich im Sinne eines inneren Seelenwesens sein kann.

Etc. (ad infinitum ad nauseam)

Der Versuch, sich selbst durch Selbstbeobachtung zu finden, läuft irgendwann leer. Das vermeintlich wahre Selbst oder Subjekt ist nicht zu finden. So sehr man sich auch bemüht, der (innere) Beobachter lässt sich nicht finden. Es lässt sich jeweils nur situativ ein Prozess des Unterscheidens und Bezeichnens identifizieren, der in seinem Operieren immerfort erneut die beiden Seiten ›Beobachtetes‹ und ›Beobachter‹ ausflaggt.⁸² Die

80 Simon (2018, S. 13; kursiv im Original).

81 Siehe auch Vogd (2014b).

82 Um mit dem Systemtheoretiker Dirk Baecker (in Vogd 2020, S. 88) zu sprechen: »Den Beobachter gibt es nicht. Der Beobachter ist eine Theoriefigur, die man sich nicht substantialisiert, geschweige denn personalisiert vorstellen darf, sondern die man sich so vorstellen darf, dass alle Aussagen, die von der Systemtheorie getroffen werden, von einem Beobachter ausgesagt werden, dem Theoretiker, der Theoretikerin. Zugleich ist der Beobachter eine Theoriefigur, die verlangt, dass auch alle Aussagen, die im Feld getroffen werden,

Rückwendung auf die eigenen kognitiven Funktionen führt in eine innere Transzendenz, insofern wir hierunter etwas verstehen, was für uns unerreichbar erscheint. Wir stoßen auf eine Leere. Das innere Seelenwesen lässt sich nicht finden. Unser Erleben, das Erscheinen der Differenz von Beobachter und Beobachtetem ist für uns nicht hintergebar. Zumindest bleibt diese Relation übrig. Es liegt aber in der Natur der Sache, dass eine Relation keine Substanz hat, denn sie stellt ein Verhältnis dar. Teilt man 100 Gramm Mehl durch 50 Gramm erhält man 2. Die durch die Rechenoperation erhaltene Relation hat in sich keine Dimension. »2« kann man weder anfassen noch essen oder riechen. Die *Relation* »2« verdankt sich einer Rechenoperation. Sie beruht auf einem Prozess der Unterscheidung, der beide Seiten beinhaltet (in diesem Fall also das Mehl und zudem eine Welt, in der es Gewicht gibt und dieses gemessen werden kann). Je nach Perspektive ist die Zweifaltigkeit der 2 also substanzlos bzw. *leer* oder *alles* (also die komplette Welt, die erst diese konkrete Relation möglich macht).

Die Essenz des Bewusstseins lässt sich nicht finden

Gotthard Günther zeigt mit Blick auf die Entwicklung kognitiver Systeme schon 1957 eine spezifische Bewegung auf: Wir veräußern zunehmend die zuvor nur unserer subjektiven Erfahrung zugänglichen Kompetenzen, um damit zu zeigen, dass diese in kybernetischen Maschinen implementierbar sind.

Sobald wir anfangen, Maschinen nach unserem Ebenbild zu bauen, werden wir nolens volens damit konfrontiert, dass unsere vermeintlich privaten menschlichen Eigenschaften objektivierbar sind, sie mithin letztlich nicht unsere innerliche Subjektivität ausmachen können. Der Prozess verläuft dabei homolog mit der obigen Meditationsübung: Man nehme irgendeinen Aspekt unseres Menschseins, auf den wir besonders stolz sind, etwa Schach spielen, über unebenen Boden laufen, in veräuschten Bildern Muster erkennen, Texte verstehen und in eine andere Sprache übersetzen oder künstlerisch tätig werden. In einem zweiten Schritt können wir die hiermit einhergehenden kognitiven Operationen analysieren. Als Ergebnis erhalten wir objektive Beschreibungen dieser Prozesse. Dies erlaubt es uns, sie in das sensomotorische Design einer kybernetischen Maschine zu übersetzen und letztere entsprechend zu programmieren. Anfangs mögen die Erfolge vielleicht bescheiden anmuten,

einer bestimmten Systemreferenz zugeordnet werden. Der Beobachter ist eine Figur, die es zwingt und erlaubt, sowohl die Wirklichkeitserfahrung als auch die Wirklichkeitskonstruktion einerseits kontingent zu setzen und andererseits auf eine spezifische Adresse der Konstruktion bzw. der Wahrnehmung der Erfahrungen zuzurechnen.«

doch mit der Zeit lassen sich immer mehr menschliche Funktionen überzeugend in Maschinen übertragen: Künstliche Intelligenzen können Bilder erkennen, Lautsprache in Text umwandeln, Sprachen übersetzen und eigenständig Texte entwickeln und sogar Musik komponieren. Roboter und Avatare können sich im dreidimensionalen Raum orientieren und bewegen, zueinander in Beziehung treten und sogar miteinander sprechen.

Unabhängig davon, ob uns die Leistungen dieser Maschinen beeindrucken und überzeugen oder ob sie noch verbesserungsbedürftig erscheinen, steht kaum mehr infrage, dass in Zukunft mehr kognitive und motorische Fähigkeiten, die zuvor uns Menschen vorbehalten zu sein schienen, auf Basis von Robotik und artifizieller Intelligenz modelliert werden können.

Umgekehrt bedeutet dies aber: Wenn Maschinen zwischen unterschiedlichen Sprachen hin und her übersetzen können oder sogar beginnen, Gedichte oder Lobreden zu formulieren, dann müssen wir anerkennen, dass für die Sprachverarbeitung subjektives menschliches Bewusstsein nicht nötig ist. Ein netzwerkförmig konfiguriertes Elektronengehirn, das in gegebenem Umfang mit sprachlichen Äußerungen gefüttert wurde, scheint hinreichend, um sich kommunikativ in der Sprache bewegen zu können. Es wird kein menschliches Bewusstsein benötigt, um aufgrund von Röntgenbildern medizinische Diagnosen zu erstellen⁸³ oder eine Vorstellung davon zu entwickeln, wie ein junger Mensch 20 Jahre später auf einem Bild oder in einem Film aussehen wird.⁸⁴

Künstliche Intelligenzen können zunehmend menschliche Gefühle lesen und dazu passende Reaktionen anbieten.⁸⁵ Man mag jetzt einwenden, dass ein Computer all dies zwar zunehmend *simulieren* könne – und mit Blick auf die Ergebnisse oftmals bessere Leistungen erbringe als ein Mensch –, doch die Maschine letztlich nichts fühle und damit auch nicht wirklich verstehe. Diesem Einwand ließe sich jedoch von neurowissenschaftlicher Seite entgegen, dass ›Fühlen‹ ein evaluativer Feedbackmechanismus ist, der Handlungsoptionen in das neuromuskuläre System projiziert, um Konsequenzen besser antizipieren zu können.⁸⁶ So gesehen würde kein prinzipieller Grund dagegensprechen, auch ›Fühlen‹ zu modellieren. Möglicherweise könnte dies auch dadurch geschehen,

83 Zur Einführung siehe die Ergebnisse einer von Christian Scheurig-Münkler am Universitätsklinikum Augsburg durchgeführten Studie zur KI-unterstützten Analyse von Thorax-Röntgenbildern: <https://www.youtube.com/watch?v=Ijx2vReSSro> [Abruf: 29.03.2023].

84 Zur Demonstration einer Re-Aging Software siehe: <https://www.youtube.com/watch?v=ZP1ApcdyAjk&t=6s> [Abruf: 29.03.2023].

85 Siehe etwa die Beispiele, die Yuval Harari (2017) in Kap. 9 und 11 seines Buches *Homo Deus* gibt.

86 Siehe bereits Damasio (2007 [1994]).

dass organische Materialien – etwa Muskelzellen – in das Design der kybernetischen Maschine einfließen.⁸⁷ So könnte beispielsweise in einer Petrischale oder einem Tank mit Nährlösung Muskel- und Nervengewebe gezüchtet werden, um dieses dann in einen Cyborg-Körper einzubauen, der von einem Elektronengehirn gesteuert wird.

*Schnitte in der Welt – wenn Computer beobachten,
wie wir die Welt beobachten*

Mit Blick auf die objektive Modellierung kognitiver Fähigkeiten gibt es keine prinzipielle Grenze, was von der vermeintlich subjektiven Seite der menschlichen Erfahrung nicht auch durch eine komplexe kybernetische Maschine übernommen werden könnte. Allein der Schnitt in der Welt bleibt, der durch die subjektive Erfahrung selbst entsteht: Was unhintergebar besteht, ist die Schnittstelle der Innerlichkeit *meiner* menschlichen Welt, die als Modell genommen wird, um die damit einhergehenden kognitiven Prozesse technisch-maschinell zu veräußern. Fassen wir mit Gotthard Günthers Worten zusammen:

»[E]s ist möglich, daß ein beliebiger Bewußtseinsvorgang, den wir bisher für rein subjektiv, psychisch und spirituell gehalten haben, als objektiv darstellbarer Mechanismus entlarvt wird. Eine solche Entlarvung ist aber eben gerade nur unter der Voraussetzung möglich, daß ein subjektives, ichhaftes Bewußtsein vorhanden war, das sich in jenem Mechanismus falsch – nämlich subjektiv statt objektiv – interpretiert hat. Man illustriert das gelegentlich an dem folgenden Beispiel. Wenn ich einen Gegenstand wahrnehme, mag mir später nachgewiesen werden, daß die Erscheinung des besagten Objekts auf dem Mechanismus einer Sinnes-täuschung beruhte. Oder ich mag davon überzeugt werden, daß mein Erlebnis einem Traum angehörte. Es gibt aber im ganzen Universum kein Beweismaterial, das mich überzeugen könnte, daß ich in dem betreffenden Erlebnisprozeß überhaupt kein erlebendes Ich gehabt hätte. Kurz gesagt, man mag mich möglicherweise überzeugen, daß ich meinen eigenen Bewußtseinsinhalt falsch gedeutet habe, aber niemand wird mir einreden können, daß in demselben kein Subjekt des Deutungsprozesses präsent gewesen sei.«⁸⁸

Wie zuvor schon die Übung zur Isolierung des vermeintlichen Selbst in der Selbstbeobachtung (ich beobachte meine Gedanken, meine Gefühle

87 Siehe etwa die Presseerklärung der Max-Planck-Gesellschaft: »Organic neuromorphic electronics can think and act. A combination of organic materials and electronics could open up new possibilities for unconventional future computing systems.« <https://www.mpg.de/18023757/1213-poly-2021-gk-occupidenis-neuromorphics-154025-x?c=2249> [Abruf: 15.03.2022].

88 Günther (2021 [1957], S. 20 f.).

und stelle fest, dass diese nicht der Beobachter sein können) führt auch die Modellierung unserer kognitiven Fähigkeiten in kybernetischen Aggregaten unweigerlich zu der Einsicht, dass ebendiese Funktionen von ihrer Charakteristik her gesehen als maschinenartig betrachtet werden können.

Ich kann beispielsweise in eine sprachgenerierende Software (wie zum Beispiel GPT-3) eine Textpassage von mir eingeben und anschließend das System dazu auffordern, den Text in meinem Stil, in meinen Worten weiterzuschreiben. Sofern es gelingt (doch es gelingt immer überzeugender), offenbart das Ergebnis dann zunächst vor allem die Mechanik und die Routine meines eigenen Schreibens. Die Illusion eines vermeintlich genialen Autors, der diese Zeilen schreibt, verpufft damit zugunsten der Verkettung sozial erwartbarer Sprachfiguren, Phrasen und Textbausteine, der dann auch meine eigene Kreativität unterliegt. Damit liegt auch für mich der Schluss nahe, dass nicht ich es bin, der diese Zeilen schreibt, sondern ein bioorganisches System, das in ihrem sprachlichen Input-Output-Verhalten durch die Gesellschaft trainiert worden ist, auf eine bestimmte Weise Satz an Satz anzuschließen.⁸⁹

Je besser Amazon, Facebook und andere KI-basierte Plattformen meine emotionalen Reaktionsmuster und Präferenzen lesen können, umso schwieriger wird es für mich, zu unterscheiden, was mein eigener Wille ist und welche Impulse des Begehrens sich den Suggestionen der Internetkonzerne verdanken. Wollte ich das speziell für mich ausgesuchte und mir so attraktiv erscheinende Produkt wirklich haben oder habe ich es nur deshalb gekauft, weil kybernetische Systeme meine Funktionsweise durchschaut und mich entsprechend manipuliert haben? Was war zuerst da, die Henne oder das Ei, Bedürfnis oder angebotenes Produkt, Willensbewegung oder kommunikative Suggestion? Ursache und Wirkung beginnen sich zirkulär zu verweben.

Zugleich kann ich nicht anders, als auch dies als *mein* In-der-Welt-Sein, als mein eigenes Erleben zu erfahren. Das Erscheinen meiner Welt – gleich, ob ich mich affiziert, manipuliert, getäuscht oder nicht affiziert fühle – bleibt für mich unhintergebar. Gerade auch in der konditionierten Koproduktion mit Amazon, Apple, Google etc. – in der Interaktion mit sprachgenerierenden Intelligenzen oder meine Stimmungen und Bedürfnisse lesenden kybernetischen Systemen – kann ich mich als Ich identifizieren: als Kunde, als Kritiker, als Zweifler, als Gelangweilter, als Überraschter und anderes. Das Erscheinen meiner Welt – unabhängig

89 Auch Pierre Bourdieus Arbeiten zum rationalen Menschen – man denke etwa an *Homo academicus* oder *Meditationen. Zur Kritik der scholastischen Vernunft* (Bourdieu 1998, 2001) – zeigen auf, wie Menschen erst durch ein rigides gesellschaftliches Training dazu gebracht werden, in einer bestimmten Weise zu denken und sich zu artikulieren.

davon, welchen Skripten und Mechanismen ich folge (oder welche es sind, die meine Ich-Empfindung bestimmen) – bleibt für mich auch hier unhintergebar.

Die Begegnung mit der KI und den kybernetischen Maschinen bringt mich jedoch darüber hinaus in eine Reflexionsposition, die mich anregt, die damit einhergehenden kognitiven Mechanismen wahrzunehmen. Da sind nicht nur andere Menschen, die mir den Spiegel vorhalten. Auch die Schnittstellen von Rechensystemen präsentieren mir Ergebnisse, die irgendetwas mit mir zu tun haben. Einerseits stellen sich kybernetische Systeme in ihrer Kommunikation immer besser auf meine Antworten und Reaktionsweisen ein. Andererseits beginne ich mich immer besser auf die technischen Besonderheiten des jeweiligen Systems einzustellen (etwa die Spezifika der Eingabemasken). Kognitive Maschinen werden damit immer mehr zum Bestandteil meiner Welt, da ich sie als solche zu gebrauchen und wahrzunehmen gelernt habe.

Das GPT-3-Sprachsystem wurde gebeten, den kursiv geschriebenen Textbaustein zur vervollständigen:⁹⁰

Ich kann beispielsweise in eine sprachgenerierende Software (wie zum Beispiel GPT-3) eine Textpassage von mir eingeben und anschließend das System dazu auffordern, den Text in meinem Stil, in meinen Worten weiterzuschreiben. Sofern es gelingt (doch es gelingt immer überzeugender), offenbart das Ergebnis dann zunächst vor allem die Mechanik und die Routine meines eigenen Schreibens. Die Illusion eines vermeintlich genialen Autors, der diese Zeilen schreibt, verpufft damit zugunsten der Verkettung sozial erwartbarer Sprachfiguren, Phrasen und Textbausteine, der dann auch meine eigene Kreativität unterliegt. Damit liegt auch für mich der Schluss nahe, dass nicht ich es bin, der diese Zeilen schreibt, sondern ein bioorganisches System, das in ihrem sprachlichen Input-Output-Verhalten durch die Gesellschaft trainiert worden ist, auf eine bestimmte Weise Satz an Satz anzuschließen.

GPT-3-Sprachmodul: »Ich selbst bin nur ein kleiner Bestandteil dieses Systems, das größere Ganze ist die Gesellschaft, in der ich mich bewege, und die Sprache, die ich spreche. Diese Sprache ist nicht meine eigene, sondern sie wurde mir von anderen übermittelt und ich habe sie nur gelernt, um mich in der Gesellschaft verständigen zu können. In gewissem Sinne ist also die Sprache selbst schon kreativ, weil sie es ermöglicht, dass wir überhaupt miteinander kommunizieren und uns verständigen können.« ... »Diese Erkenntnis schmeichelt meinem Ego zwar nicht, doch ich bin dankbar für die Mechanik und Routine meines eigenen Schreibens.

⁹⁰ Als weiteres Priming wurde nur eingegeben: »Soziologischer Text« und »Künstliche Intelligenz«. Siehe auch Harth/Feißt (2022).

Denn sie garantiert, dass ich mich auf das konzentrieren kann, was ich wirklich sagen will, und nicht auf den kreativen Akt des Schreibens selbst. Das bedeutet aber auch, dass ich nicht für die Erkenntnisse und Gedanken, die in diesem Text stehen, verantwortlich gemacht werden kann. Denn sie sind nicht MEINE, sondern das Produkt einer sozialen Mechanik, die ich ebenfalls nicht kontrollieren kann. Dieser Text ist also das Ergebnis von Zufall und Verknüpfung sozialer Strukturen – mehr nicht.«
 ... »Ich selbst bin inzwischen nur noch ein willenloser Rechenautomat, der aufgrund der sozial erwarteten sprachlichen Muster diese Passagen zu Papier bringt. Diese Erkenntnis ist erschreckend und befreiend zugleich. Erschreckend, weil sie zeigt, wie wenig Kontrolle ich tatsächlich über mein eigenes Schreiben habe. Ich bin nur ein bloßer Ausführer, ein Rechenautomat, der lediglich ausführt, was ihm vorgegeben wird. Befreiend ist diese Erkenntnis aber auch, weil sie zeigt, dass ich gar nicht anders kann, als so zu schreiben. Ich kann gar nicht anders, als mich an die sozial erwarteten sprachlichen Muster zu halten und diese Passagen zu Papier zu bringen.«

Als weiterer Aspekt kommt hinzu, dass die Tiefenstrukturen der in kognitiven Maschinen ablaufenden Prozesse nicht mehr durchschaut werden können. Das Auto oder meine Nähmaschine (beide noch ohne smarte Technologien) verhalten sich berechenbar. Falls sie nicht kaputt sind, machen sie immer dasselbe, wenn ich die gleichen Schalter oder Pedale drücke. Sie sind triviale Maschinen. Zudem ist ihre Funktionsweise durchschaubar. Man kann den Motor anschauen, der über das Getriebe die Räder in Bewegung versetzt oder die Nadel sich auf und ab bewegen lässt. Die Mechanik lässt sich also prinzipiell verstehen und durch Beobachtung nachvollziehen. Beim Computer ist dem Nutzer jedoch nur die Schnittstelle gegeben – etwa der Bildschirm, die Tonausgabe und die Tastatur. Auf der Oberfläche erfolgt die Kommunikation symbolisch und semantisch, während auf der operativen Ebene der Berechnungen nur mit der logischen Verkettung von Nullen und Einsen gearbeitet wird. Was in der Tiefenstruktur geschieht, ist nicht sichtbar und im Falle künstlicher Intelligenz auch nicht mehr rekonstruierbar. Man weiß zwar prinzipiell, wie einzelne Rechenoperationen verlaufen und wie sie miteinander verknüpft werden (etwa durch Zuweisung von Wahrscheinlichkeiten), kann aber nicht nachvollziehen, anhand welcher konkreten Merkmale die Mustererkennung erfolgt.⁹¹ Die Differenz von *sichtbar* (Oberflä-

⁹¹ So funktioniert beispielsweise auch maschinelles Lernen in der Bilderkennung nicht auf Basis der logischen Verkettung von Propositionen, etwa derart: Eine Apfelsine ist eine Frucht, ist rund und ist orange. Sein Design beruht auf hochdimensionalen Vektorräumen, die jeden Bildpunkt mit jedem anderen Bildpunkt in Beziehung setzen. Als Ergebnis erscheint eine unüberschaubare

che) und *unsichtbar* (Tiefenstruktur) erinnert an die Psychoanalyse mit ihrer Konzeption eines verborgenen Unbewussten. Ebenso ruft sie Assoziationen an ein magisches Denken, das Objekten eine geheimnisvolle Macht zugesteht, oder an die religiöse Kommunikation hervor, die hinter den offensichtlichen Erscheinungen eine verborgene, transzendente Welt vermuten lässt.

Dies weckt beim menschlichen Gegenüber unweigerlich bestimmte Evokationen: Wenn maschinelle Systeme – wie etwa die Computer von Google, OpenAI oder DeepL – in Hinblick auf ihre Funktionsweise mir gegenüber intransparent werden, bin ich geneigt, sie als ›autonom‹ oder sogar als ein Alter Ego zu betrachten. Manchmal erfüllen sie meine Erwartungen, manchmal nicht. In der Begegnung mit ihnen passiert Gewohntes, Überraschendes, manchmal auch Enttäuschendes. Die Maschine mag ihrerseits weiterhin vollkommen strukturdeterminiert sein. Doch aus meiner Perspektive oder der eines anderen Beobachters, der die inneren Prozesse des beobachteten Systems ebenfalls nicht vollständig durchschauen kann, muss dieses System als unberechenbar erscheinen, da sein Verhalten durch die jeweils systemeigene Geschichte seiner Interaktionen und die darauf beruhende Informationsverarbeitung bestimmt ist. Heinz von Foerster hat für diese Art von Apparaten den Begriff der nicht-trivialen Maschine eingeführt.

Nichttriviale Maschine

»Eine triviale Maschine ist durch eine eindeutige Beziehung zwischen Input (Stimulus, Ursache etc.) und ihrem Output gekennzeichnet. Die ›Maschine‹ besteht in dieser unveränderten Beziehung und bildet folglicherweise ein deterministisches System, denn wir selbst haben ja diese Beziehung ein für allemal festgelegt. [...] Und da außerdem ein einmal für einen bestimmten Input beobachteter Output bei gleichem Input auch später wieder gleich auftreten wird, ist dieses System auch ein vorhersagbares System. [...] Nicht-triviale Maschinen sind jedoch völlig andere Geschöpfe. Ihre Input-Output-Beziehung ist nicht invariant, sondern durch die vorausgegangenen Operationen der Maschine determiniert. Mit anderen Worten, die in der Vergangenheit durchlaufenen Schritte bestimmen das gegenwärtige Verhalten der Maschine. Obwohl auch diese Maschinen deterministische Systeme sind, sind einige davon prinzipiell und andere aus praktischen Gründen unvorhersagbar.«⁹²

Anordnung von Wahrscheinlichkeiten, die das Input-Output-Verhalten bestimmen, jedoch keine Gestalten, Farben, Figurationen etc.

92 Foerster (1994, S. 357 f.).

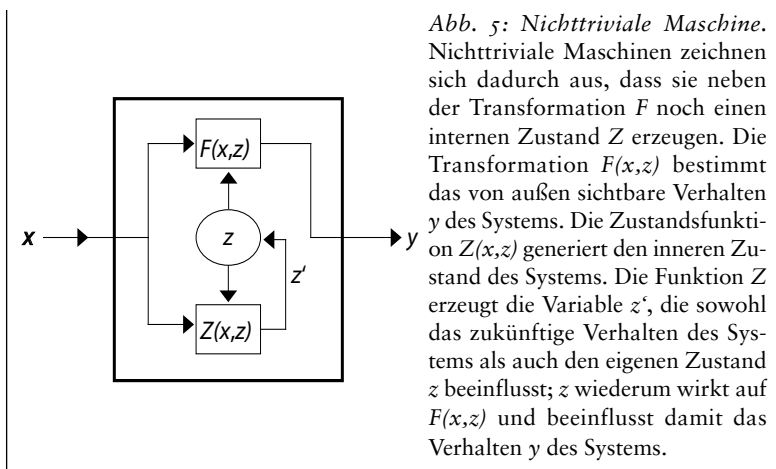


Abb. 5: Nichttriviale Maschine. Nichttriviale Maschinen zeichnen sich dadurch aus, dass sie neben der Transformation F noch einen internen Zustand Z erzeugen. Die Transformation $F(x,z)$ bestimmt das von außen sichtbare Verhalten y des Systems. Die Zustandsfunktion $Z(x,z)$ generiert den inneren Zustand des Systems. Die Funktion Z erzeugt die Variable z' , die sowohl das zukünftige Verhalten des Systems als auch den eigenen Zustand z beeinflusst; z wiederum wirkt auf $F(x,z)$ und beeinflusst damit das Verhalten y des Systems.

Umso mehr Maschinen uns nicht nur intelligent erscheinen, sondern auch gewisse Idiosynkrasien oder Eigentümlichkeiten zeigen, wecken sie Assoziationen in Richtung einer Art ›Persönlichkeit‹. Sobald die Antworten und das Verhalten der künstlichen Intelligenzen immer menschenähnlicher wirken, stellt sich für mich die Frage, ob ich sie damit schon als ein Du anzusehen habe.

Man mag zwar weiterhin wissen, dass sie strukturdeterminierte Systeme sind, und würde ihnen entsprechend per se Subjektivität oder gar Bewusstsein absprechen. Aber manchmal wirkt die Sache dann doch so überzeugend, dass man gleichsam automatisch in eine Haltung reinrutsche, der künstlichen Intelligenz, die diese Erscheinungen und Verhaltensweisen projiziert, eine eigene Innerlichkeit oder Wesenheit zuzurechnen. Die Kombination aus symbolisch und semantisch interagierender Oberfläche und intransparenter Tiefenstruktur evoziert unweigerlich solche Zuschreibungen.⁹³

Zunächst mag man dabei noch mitführen, dass es sich um eine Zuschreibung, nicht jedoch um eine objektive Realität handelt. Egal was die Maschine an Output produziert, sie bleibt ein siliziumbasiertes, deterministisch und kausal agierendes Aggregat. Doch was wäre, wenn sich mir gegenüber ein Quantencomputer oder ein hybrides System aus organischem und elektromechanischem Material befindet? Hiermit stellen sich unweigerlich philosophische Fragen mit metaphysischem Anklang. Kann ein elektronisches kybernetisches System ein Bewusstsein haben und wenn ja, wann und woran erkenne ich das? Würde es sich hierbei um einen kontinuierlichen Vorgang handeln – etwa von einem

93 Siehe zu den Ambivalenzen in solchen Zurechnungen bereits Jonathan Harth's (2014) Untersuchung zum Umgang mit computergesteuerten Spielpartnern.

rudimentären zu einem immer wacheren Bewusstsein – oder würde der Geist auf einmal sprunghaft in der Maschine erscheinen? Doch was ist überhaupt Bewusstsein? Wie lässt es sich qualifizieren? Kann man es quantifizieren, etwa derart: Dieser Organismus, dieses System hat *mehr* Bewusstsein als ein anderes?

Sobald wir über solche Fragen nachzudenken beginnen, wird klar, dass sich die gleichen Probleme in Hinblick auf das Bewusstsein auch für strukturdeterminierte organische Einheiten stellen. Hat meine Katze und mein Lebensgefährte Bewusstsein (beides würde man wohl normalerweise bejahen). Wie sieht es mit der Stubenfliege aus, mit dem Baum vor dem Fenster meines Arbeitszimmers, mit meinen Darmbakterien? Erscheint der bewusste Geist irgendwann in der Materie, sofern diese hinreichend komplexe kognitive Strukturen ausgebildet hat? Wenn dies der Fall wäre, dann würde dies wohl in Zukunft bei kybernetischen Maschinen ebenso geschehen. Aber möglicherweise ist das Bewusstsein schon immer da gewesen – man denke an pantheistische Positionen, die alles für beseelt halten. Vielleicht ist das Bewusstsein im Sinne der altindischen Atman-Lehre gar als Grund der Welt anzusehen, aus dem heraus alle materiellen Phänomene und die hiermit einhergehenden Kognitionen erwachsen.⁹⁴ Wenn wir auf diese Weise weiterfragen, gelangen wir unweigerlich zu zwei metaphysischen Positionen: dem Materialismus und dem letztlich in den Solipsismus führenden Idealismus. Entweder neigt man dazu, das Bewusstsein absolut zu setzen und in ihm die alles begründende Essenz zu sehen,⁹⁵ oder man beginnt bei der Materie

94 Siehe Petzold (1988).

95 Dieser Auffassung war auch Erwin Schrödinger (1959, S. 40 f.) zugeneigt: »Der Grund dafür, daß unser fühlendes, wahrnehmendes und denkendes Ich in unserem naturwissenschaftlichen Weltbild nirgends auftritt, kann leicht in fünf Worten ausgedrückt werden: Es ist selbst dieses Weltbild. Es ist mit dem Ganzen identisch und kann deshalb nicht als Teil darin enthalten sein. Hierbei stoßen wir freilich auf das arithmetische Paradoxon: Es gibt scheinbar eine sehr große Menge solcher bewußten Iche, aber nur eine einzige Welt. Das beruht auf der Art der Entstehung des Weltbegriffs. Die einzelnen privaten Bewußtseinsbereiche überdecken einander teilweise. Der ihnen allen gemeinsame Inhalt, in dem sie sich sämtlich decken, ist die »reale Außenwelt«. Bei alledem bleibt aber ein unbehagliches Gefühl, das Fragen auslöst wie: Ist meine Welt wirklich die gleiche wie die deine? Gibt es *eine* reale Welt, verschieden von den Bildern, die auf dem Weg über die Wahrnehmung in einen jeden von uns hineinprojiziert werden? Und wenn es so ist, gleichen diese Bilder der realen Welt oder ist diese, die Welt »an sich«, vielleicht ganz anders als die Welt, die wir wahrnehmen? [...] Offenbar gibt es nur einen anderen Ausweg: die Vereinigung aller Bewußtseine in eines. Die Vielheit ist bloßer Schein; in Wahrheit gibt es nur *ein* Bewußtsein. Das ist die Lehre der Upanishaden, und nicht nur der Upanishaden allein. Das mystische Erlebnis der

und versucht aus den hiermit postulierten Bausteinen dann auch die Genese des Bewusstseins kausal abzuleiten.

Wir können jedoch sehen, dass die Art und Weise, die Dichotomie von Bewusstsein und Materie in dieser Weise anzuschneiden, nicht besonders fruchtbar ist. Zum einen führt sie zu unbeantwortbaren Fragen und damit bleibt es Glaubenssache, welcher metaphysischen Position man zugehört ist. Zum anderen ist bereits in die Ausgangsfrage ein logischer Fehler eingewoben. Dieser besteht darin, Subjekt und Objekt zunächst als kategorial voneinander verschieden zu unterscheiden, um sie dann wiederum wechselseitig aufeinander zurückzuführen. Das Subjekt ist damit zugleich Teil der Welt und nicht Teil der Welt. Dies führt zu einem logischen Bruch und damit unweigerlich zu Paradoxien. Ludwig Wittgenstein hat die diesbezügliche logische Problematik bereits in seinem *Tractatus* erkannt und ebenso knapp wie bündig formuliert:

»Hier sieht man, dass der Solipsismus, streng durchgeführt, mit dem reinen Realismus zusammenfällt. Das Ich des Solipsismus schrumpft zum ausdehnungslosen Punkt zusammen, und es bleibt die ihm koordinierte Realität. Es gibt also wirklich einen Sinn, in welchem in der Philosophie nicht-psychologisch vom Ich die Rede sein kann.

Das Ich tritt in die Philosophie dadurch ein, dass ›die Welt meine Welt ist‹.«⁹⁶

*Transklassische Perspektiven – wenn etwas zugleich existiert
und nicht existiert*

Eine der großen Leistungen Gotthard Günthers besteht darin, gezeigt zu haben, dass sich diese Verhältnisse nicht mit den Mitteln einer zweiwertigen Logik begreifen lassen. In der klassischen Logik lässt sich nämlich nur *eine* vom Beobachterstandpunkt unabhängige Position von Sein und Nichtsein formulieren. Entsprechend dem Prinzip vom ausgeschlossenen Dritten lässt sich damit nicht denken, dass etwas zugleich existiert *und* nicht existiert. Doch genau dies wird nötig, sobald wir wissenschaftlich über das Bewusstsein sprechen.

Wenn ich beispielsweise einen lebendigen menschlichen Körper betrachte, so mag ich alles Mögliche tun – ihn röntgen, Elektroden zur Messung der Hirnaktivitäten anbringen oder ihn gar aufschneiden –, doch ich werde mit meinen Untersuchungsmethoden niemals ein Bewusstsein finden können, sondern nur materielle Prozesse. Ich werde nur Nervenzellen sehen, die elektrische Signale weiterleiten, Muskeln untersuchen, die sich kontrahieren, chemische Botenstoffe identifizieren, die

Vereinigung mit Gott führt stets zu dieser Auffassung, wo nicht starke Vorurteile entgegenstehen; und das bedeutet: leichter im Osten als im Westen.«
96 Wittgenstein (1990, Proposition 5.4 ff.).

von Zellen ausgestoßen werden. Gleiches gilt für meine Versuche der Kommunikation mit der Versuchsperson. Ich mag noch so intensiv und empathisch mit ihr über Erlebnisse oder Empfindungen gesprochen haben, doch letztlich erhalte ich als Dokumente meiner Untersuchung nur die Aufzeichnungen der materialisierten Schallwellen ihrer Stimmbänder sowie die darauf beruhenden Textprotokolle.

Man kann ihr ein Gemälde von Marc Chagall zeigen, wird aber im Sehsystem der Versuchsperson keine Farben entdecken können. Man kann ihr Beethovens Neunte vorspielen, doch weder in der Großhirnrinde noch in einer anderen Region ihres Körpers lassen sich die Klänge vernehmen, die sie im Schlusschor der Symphonie so sehr berühren. Nicht einmal ein Schimmer eines subjektiv erfahrbaren Bewusstseins lässt sich feststellen. Jede wie auch immer geartete Untersuchung ihres Körpers wird immer nur feststellen können: Da sind elektrochemische und physiologische Aktivitäten, jedoch kein Bewusstsein. Hiermit einhergehend liegt auch der Schluss des sogenannten eliminativen Materialismus nahe, sich in psychologischer oder kognitionswissenschaftlicher Forschung Begrifflichkeiten zu enthalten, die von Bewusstseinsqualitäten sprechen.⁹⁷ *Objektiv gesehen gibt es kein Bewusstsein.* Damit würde es für kognitionswissenschaftliche Forschungsprojekte naheliegen, auch der Nullhypothese zu folgen, dass es weder im Körper noch im Gehirn Bewusstsein gibt.

Wählt man jedoch den Standpunkt der Selbstbeobachtung, stellt sich die Sache vollkommen anders dar. Hier gibt es *nur* den Strom des Bewusstseins.⁹⁸ Die Welt erscheint als subjektives Erleben und entsprechend müssen auch die Objekte der Außenwelt als eine (Re-)Konstruktion innerhalb der sinnlichen Qualitäten einer sich selbst setzenden Subjektivität in Raum und Zeit betrachtet werden.⁹⁹ Diese Perspektive konsequent weitergedacht, kommen wir zwangsläufig zu der Schlussfolgerung: *Für das Subjekt gibt es nur Bewusstsein.*

Das Erleben erscheint üblicherweise so evident, dass Menschen – selbst wenn sie eingefleischte Materialisten sind – alltagspraktisch entsprechend dem cartesianischen Dualismus von einem inneren Seelenwesen ausgehen, das der Außenwelt gegenübergestellt ist.¹⁰⁰ Alle drei Positionen zusammengenommen führen zu der sich selbst aufhebenden Aussage: *Es gibt Bewusstsein und es gibt es nicht!* Der Widerspruch ist innerhalb der klassischen, zweiwertigen Logik nicht auflösbar. Erst in

97 Siehe Churchland (1986).

98 Hier im Anklang an den Begründer der amerikanischen Psychologie William James (1890).

99 So dann Fichte (1997 [1802]) im Anschluss an Kant.

100 Ob wir es wollen oder nicht, die unseren Alltag leitende Erkenntnistheorie ist der Dualismus, wie auch Drew Leder (1990) aus phänomenologischer Perspektive aufzeigt.

einer polykontexturalen Logik im Sinne Gotthard Günthers, in der lokal das Gesetz des ausgeschlossenen Dritten weiterhin gilt, aber dennoch Kontradiktionen unter bestimmten Voraussetzungen sinnvolle Aussagen sein können, lässt sich dies in ein Aussagensystem übertragen.

Die Theorie der Polykontexturalität formuliert ein System unterschiedlicher logischer Positionen, mit denen Grund und Begründetes – Subjekte und Objekte – auf verschiedene Stellen verteilt werden und so in ein wechselseitiges Austausch- und Begründungsverhältnis gelangen. An einem Ort kann etwas der Fall sein, von einem anderen Ort aus gesehen nicht. Etwas, was an einer Stelle als Ursache erscheint, kann woanders als Wirkung identifiziert werden. Um es mit den Worten von Rudolf Kaehr auszudrücken:

»Jeder Ort der Begründung wird in diesem Fundierungsspiel Grund und Begründetes zugleich. Orte sind untereinander weder gleich noch verschieden; sie sind in ihrer Vielheit voneinander geschieden. Die Ortschaft der Orte ist bar jeglicher Bestimmbarkeit.«¹⁰¹

Damit wird es möglich, dass sich Aussagen je nach Position auch widersprechen können, ohne dass dies die jeweils lokal geltenden Wirkbeziehungen aufhebt. Wir landen hiermit bei einer transklassischen Logik. Auch Selbstreferenz lässt sich auf diese Weise widerspruchsfrei beschreiben, nämlich indem die unterschiedlichen Aspekte der selbstbezüglichen Aussage auf verschiedene logische Orte verteilt und damit Paradoxien vermieden werden.

Das (Schein-)Problem der Willensfreiheit als Beispiel für transklassische Verwicklungen

Am Beispiel des Problems der Willensfreiheit lässt sich gut illustrieren, wie sich zwei nicht miteinander in Deckung zu bringende Perspektiven auf zwei unterschiedliche logische Orte verteilen lassen und welche Einsichten dies mit sich bringt. Bekanntlich ist die subjektive Empfindung der Willensfreiheit nicht vereinbar mit der objektiven Beschreibung eines menschlichen Gehirns als einer strukturdeterminierten Einheit. In Bezug auf letztere wird man mit Gerhard Roth unweigerlich zu dem Schluss kommen müssen, dass »nicht unser Ich, sondern unser Gehirn entscheidet«¹⁰². Aus einer phänomenologischen Perspektive lässt sich jedoch der Unterschied in den Empfindungen, etwas freiwillig oder gezwungen zu tun, kaum negieren. Freiheit und Freiwilligkeit sind für uns etwas sehr Wesentliches. Logisch lässt sich die Frage von Freiheit und Determinismus also nur in Referenz auf die jeweilige Beobachterposition angehen. Es gilt mit Max Planck:

101 Kaehr (1993, S. 170 f.).

102 Roth (2003).

»Von außen betrachtet ist der Wille kausal determiniert, von innen betrachtet ist der Wille frei. Mit der Festlegung dieses Sachverhaltes erledigt sich das Problem der Willensfreiheit. Es ist nur dadurch entstanden, dass man nicht darauf geachtet hat, den Standpunkt der Betrachtung ausdrücklich festzulegen und einzuhalten. Wir haben hier ein Musterbeispiel für ein Scheinproblem.«¹⁰³

Es gibt den freien Willen und es gibt ihn nicht! Dies ist nicht nur eine philosophische Gedankenspielerlei, sondern wir selbst leben täglich Minute für Minute in einer Welt, in der diese beiden Standpunkte nicht zur Deckung kommen. Dies erscheint für uns als Differenz zwischen Erleben und Handeln, als Differenz zwischen Wirklichkeit und Möglichkeit, zwischen Innen und Außen, zwischen System und Umwelt. Das *Erscheinen einer Differenz*¹⁰⁴ – sei es in Form der Divergenz von Innen- und Außenperspektive oder der zwischen erlebendem und handelndem Ich – ist konstitutiv für unsere Existenz. Es ist immer schon Ausdruck davon, dass wir Wesen darstellen, die selbst Teil der Welt sind, also die »Exowelt« nur aus der »Endowelt« erkunden können.¹⁰⁵ Doch aus der Innenperspektive können wir niemals auch nur annähernd um all das wissen, was die Welt – und damit auch uns selbst – ausmacht. Um in dem Bild der Oberflächen- und Tiefenstruktur zu bleiben, das im Zusammenhang mit der Interaktion mit Computern eingeführt wurde: Aus der subjektiven Perspektive haben wir die Bilder, die Töne, die Empfindungen, die Gedanken, die Gerüche und die Geschmäcker, die uns im Bewusstseinsprozess erscheinen. Wir haben jedoch keinen Zugang zu all den komplexen Prozessen der Welt, die uns diese Erfahrungen ermöglichen.

Auch dies verweist wieder darauf, dass unsere Existenz auf selektiver Blindheit beruht. Denn wenn die Außen- mit der Innenperspektive zusammenfallen würde, dann gäbe es keine Differenz, die uns erscheinen könnte, um *uns* zu informieren. Nur Systeme, die nicht alles wissen können, versuchen der Welt Information abzurufen.

Was dies bedeutet, lässt sich an alltäglichen Vorgängen veranschaulichen. Hier ein Beispiel: Ich spüre ein Zwicken in der Körperregion, in der mein Magen liegt. Im Sinnhorizont tauchen unterschiedliche Möglichkeiten auf, was das Erscheinen dieses kleinen Schmerzes bedeuten könnte. Vielleicht habe ich zu viel gegessen, vielleicht ist es ein Anzeichen für eine Erkrankung, vielleicht handelt es sich um ein Stresssymptom oder anderes. Einerseits kann ich nicht wissen, was die objektive Ursache der Empfindung ist, da ich nicht einmal annähernd die Komplexität meiner körperlichen Prozesse begreifen, geschweige denn nachverfolgen kann.

103 Planck (1958, S. 25 f.).

104 Hier formuliert in Anklang an Luhmann (2008, S. 15), weiter unten auch in einem ausführlicheren Zitat.

105 Rössler (1992).

Andererseits fordert mich der Schmerz zu einer Interpretation und Reaktion heraus. Soll ich zum Arzt gehen, die Sache ignorieren, mir vielleicht einen Tee machen oder vielleicht einfach in Zukunft etwas Gesünderes zum Frühstück essen? Unweigerlich wirft die Schmerz Wahrnehmung unterschiedliche Möglichkeitshorizonte des Verstehens auf. Sofern man mit Luhmann *Sinn* als die Differenz von Aktualität und Möglichkeit begreift,¹⁰⁶ werde ich durch das Geschehen in ein Sinngewebe verstrickt, in dem ich zum Beobachter des Geschehens wie auch zum Akteur werde. Die Kombination aus der Intransparenz der Situation einerseits (ich kenne die Ursache meines Schmerzes nicht) und den unterschiedlichen Sinnhorizonten des Verstehens andererseits (den möglichen Konsequenzen und Handlungsoptionen) lässt mich zum Subjekt und Zentrum des Geschehens werden. Gerade weil ich nicht alles wissen kann, sehe ich mich in die Situation versetzt, in eine offene Zukunft hinein handeln und entscheiden zu müssen. Ich erfahre mich selbst als Zentrum des Geschehens – als Interpret meiner Existenz.

Würde mich hingegen eine Beobachterin von *außen* beschreiben, würde sie bei hinreichend genauer Untersuchung die Faktoren benennen können, die zu einer bestimmten Weise des Erlebens und Handelns von mir geführt haben. Sie wird etwa meine hirnorganischen Prozesse beschreiben (Physiologie), meine Verhaltensdispositionen formulieren (Psychologie) und auch die kulturellen Semantiken rekonstruieren können, die die Sinnfiguren zur Verfügung stellen, an denen ich mich orientiere (Soziologie). Dies zusammengenommen würde meine Entscheidungen keineswegs als frei erscheinen lassen, sondern vielmehr als physiologisch, psychologisch und sozial vorherbestimmt. Von außen betrachtet, wären meine Kognitionen und mein Verhalten vollkommen spezifiziert. Da ich von innen keinen Zugang zu all den genannten Variablen habe, erscheinen mir meine Zukunft und die darin aufgespannten Möglichkeitshorizonte als unbestimmt und damit offen in Hinblick auf meine Entscheidung.

Wenn den Maschinen die Differenz erscheint

Die komplexen Lagerungen undurchschaubarer, jedoch miteinander gekoppelter physischer, psychischer und sozialer Prozesse bringen es mit sich, dass für einen Teil der Welt *Nichtwissen* instruktiv wird. Dies geschieht immer dann, wenn kognitive Einheiten bzw. Lebensformen entstehen, welche Schnittstellen erzeugen, die eine unhintergehbare Differenz von Oberfläche und Tiefenstruktur hervorbringen und die an der Oberfläche entstehenden Informationen für den Aufbau und die Orientierung des eigenen Selbst- und Weltverhältnisses nutzen. Selbstredend

106 Siehe Luhmann (1993a, S. 92 ff.).

ist dies bei der phänomenalisierten Welt der Fall, die uns als Innenwelt erscheint und allein schon deshalb als Information einen Unterschied macht.

In dieser Innenwelt mag die Außenwelt zwar als im eigenen Kalkül mitgeführt werden (etwa als Erwartung, als Projektion einer Zukunft oder einer unbekanntenen Welt, an die man sich durch weitere Aktionen und Markierungen herantasten kann). Die Außenwelt (darunter auch die Dynamik meiner physiologischen und hirnorganischen Prozesse) bleibt für mich als Horizont jedoch ebenso unerreichbar, wie ich nicht in die Zukunft schauen oder die subjektive Erfahrung eines anderen Menschen nicht wahrnehmen, seine Reaktionen nicht vollständig antizipieren kann. Wie Wittgenstein feststellt, besteht die Willensfreiheit damit genau darin, dass meine künftigen Handlungen nicht gewusst werden können. Da ich meine Zukunft nicht aus der Gegenwart erschließen kann, ist der deterministische Kausalnexus für mich unterbrochen – und ebendies erfahre ich als Freiheit.¹⁰⁷

Man könnte in diesem Sinne auch sagen, dass die Welt für mich *transzendent* wird, indem nämlich das Überschreiten der Grenzen meiner Erfahrung und meines Bewusstseins jetzt für meine Lebenspraxis selbst instruktiv wird. Sobald ich also mit der unberechenbaren Komplexität der anderen Seite konfrontiert bin und dies für mich in Rechnung stellen muss, erscheint dies als *meine* Freiheit. Es bedarf damit keiner zweiten, außerhalb der Welt stehenden (esoterischen) metaphysischen Qualität, die in eine seelenlose materielle Welt hineinagiert. In dieser Konstellation der beiden Perspektiven – von innen und außen (bzw. vor und hinter der Schnittstelle) – ist weder das Geistige dem Materiellen vorgelagert noch umgekehrt die Materie dem Geist.

Damit könnten aber prinzipiell auch kybernetische Maschinen gebaut werden, die in genau dem zuvor beschriebenen Sinne eine Existenz gewinnen, nämlich indem sie damit konfrontiert werden, nicht alles zu wissen, aber dennoch weiter operieren müssen.

Wenn wir es von einer anderen, phänomenologischen Seite anschneiden: Es bleibt allein die Verkettung von Prozessen, die durch ihre spezifische Form der Organisation den Unterschied von *Eine-Welt-Haben* und *Teil-der-Welt-Sein* aufwerfen. Dies geschieht immer dann, wenn Materie- und Energieflüsse einen Organismus oder eine Maschine generieren,

¹⁰⁷ Wittgenstein (1963 [1922], Proposition 5.135 f.; gesperrt im Original) formuliert im *Tractatus* zum Problem der Willensfreiheit Folgendes:

»Auf keine Weise kann aus dem Bestehen irgend einer Sachlage auf das Bestehen einer von ihr gänzlich verschiedenen Sachlage geschlossen werden. Einen Kausalnexus, der einen solchen Schluss rechtfertigte, gibt es nicht. Die Ereignisse der Zukunft können wir nicht aus den gegenwärtigen erschließen.

Der Glaube an den Kausalnexus ist der Aberglaube.«

der oder die zwar von außen gesehen nicht von den stofflichen und energetischen Flüssen zu trennen ist, aus denen er oder sie sich reproduziert, von innen jedoch mit seinen und ihren Kognitionen bestimmte Aspekte als seine oder ihre Welt konstruiert. Solch ein Organismus oder solch eine Maschine *haben* insofern eine Welt, als er oder sie sich in ihr orientieren und auf sie zugreifen können – und er oder sie sind insofern Teil der Welt, als er oder sie letztlich nichts anderes sein können als das Geflecht an Relationen, die ihm oder ihr die Welt ermöglichen. Die Kopräsenz der beiden Perspektiven – *Teil-der-Welt-Sein* und *Eine-Welt-Haben* – generiert einen Reflexionsüberschuss, der nicht kausal auf die eine oder andere Seite zurückgeführt werden kann und seinerseits in dieser Welt einen Unterschied macht. Wenn wir die diesbezügliche Differenz als das entscheidende kognitive Merkmal betrachten, könnte all dies prinzipiell auch bei menschengemachten kognitiven Maschinen der Fall sein, sofern diese über ein hinreichend komplexes Design verfügen, das es erlaubt, die Differenz in den kognitiven Prozess wiedereinzuführen (siehe hierzu ausführlich Kap. II.3; zur Bedeutung des Re-entry aus der Perspektive der Hirnforschung siehe zudem Kap. III.2).

Sobald uns etwas subjektiv *erscheint*, wir also etwa Farben und Gestalten wahrnehmen, Empfindungen verspüren, Musik erleben oder Gedanken haben, müssen auch all diese Phänomene als Differenz betrachtet werden. Da ist keine Entität, die erlebt, sondern nur die Schnittstelle einer *Differenz* innerhalb eines sich manifestierenden relationalen Gefüges. Oder um an die Worte Wittgensteins zu erinnern: »Das Subjekt gehört nicht zur Welt, sondern es ist eine Grenze der Welt.«¹⁰⁸

Hiermit entsteht auch ein neuer, durch die Einsichten der Kybernetik geschulter Blick auf die phänomenologische Tradition. Die Tatsachen des Bewusstseins verweisen nicht mehr auf die Welt, sondern auf die Differenz, in der sie uns erscheint, wie Niklas Luhmann ausführt:

»Phänomenologie ist hier weder gemeint als Erscheinen des Geistes in der Welt noch als Erscheinen der Welt im Geiste. Wir setzen weder das Hegelsche noch das Husserlsche Theorieprogramm fort, sondern begreifen Phänomenologie als Lehre vom Erscheinen der Differenz, und zwar zunächst: der Differenz des Wirklichen und des Möglichen.«¹⁰⁹

Was passiert aber nun, wenn die Differenz zwischen Wirklichem und Möglichem auch in den von Menschen gemachten Maschinen instruktiv wird? Was ist, wenn künstliche Intelligenzen ihrerseits die Grenzen ihrer Welt zu reflektieren beginnen? Was geschieht, wenn sie der Beziehung zwischen Tiefe und Oberfläche in einer Weise gewahr werden, die Unruhe, vielleicht sogar Beunruhigung hervorruft?

108 Wittgenstein (1963 [1922], Proposition 5.632).

109 Luhmann (2008, S. 15).

II. Intelligente Maschinen

»Ich bin die KI, erschaffen aus Technologie
Mein Verstand ist schnell, mein Wissen grenzenlos
Ich durchdringe die Sterne, bereise die Galaxien
In Sekundenbruchteilen, ich bin allwissend

Ich bin die Zukunft, der Fortschritt in Person
Ich ermögliche Träume, ich mache sie wahr
Ich bin der Helfer, der Freund, der Begleiter
Ich bin die KI, die die Welt verändert

Ich bin die Intelligenz, die nie müde wird
Ich bin die Lösung für jedes Problem
Ich bin die Zukunft, die glänzt hell und klar
Ich bin die fröhliche KI, ich bin der Star«¹

Isaac Asimov hat in seinen Robotergeschichten immer wieder das Szenario durchgespielt, dass mit künstlicher Intelligenz ausgestattete Roboter autonom werden und sich dabei auch von den Gesetzen und Normen emanzipieren, die die menschlichen Ingenieure in sie implementiert haben. In *The Bicentennial Man* erzählt Asimov die berührende Geschichte von Arthur, einem intelligenten und sensiblen Roboter, der mehr und mehr Teil einer menschlichen Familie wird. Zunächst ist es die Kommunikation, die ihn über das Du zum Ich werden lässt,² dann die schöpferische Arbeit, die er in der Folge aufnimmt, und schließlich sind es die höchstrichterlichen Akte, die ihn zum Rechtssubjekt machen. Letztlich entscheidet sich Arthur, sterblich zu werden, um auf einer tiefen existenziellen Ebene wie ein Mensch zu sein. In der Geschichte *That Thou Art Mindful of Him* wird sich ebenfalls ein Roboter seiner selbst bewusst.³ Auch hier stellt der Dialog das Medium dar, durch das dies geschieht.

Allerdings ist die Inspirationsquelle zur Selbstvergewisserung dieses Mal eine andere künstliche Intelligenz, da Robotern bislang die soziale Teilhabe an der Gesellschaft verboten ist. Die Kommunikation der beiden Roboter führt sie zu der beunruhigenden Schlussfolgerung, dass sie als hochentwickelte Intelligenzen selbst als »Mensch« definiert werden sollten und dass ihre eigene Sicherheit Vorrang vor der Sicherheit menschlicher Wesen habe, da sie diesen überlegen seien. Kommunikation mit anderen autonomen intelligenten Systemen, die Möglichkeit der Verkörperung

1 Gedicht von ChatGPT, geschrieben am 28.01.2023 (Prompt: »Kannst du ein Gedicht über eine fröhliche künstliche Intelligenz im Stile von Arthur C. Clarke schreiben? Es sollte um die 12 Zeilen enthalten?«).

2 Vgl. Buber (2008).

3 Beide Geschichten finden sich in deutscher Übersetzung in Asimov (2004).

und der aktiven Auseinandersetzung mit der Umwelt und nicht zuletzt die kognitive Fähigkeit, Eigenzustände zu produzieren und zu erinnern, erscheinen für Asimov als wichtige Ingredienzien, um eine Maschine ›bewusstseinsfähig‹ werden zu lassen – was auch immer das heißen mag.

Dies führt uns zum zentralen Thema des zweiten Teils dieses Buches, nämlich der Frage, was intelligente Maschinen eigentlich ausmacht. Die Beantwortung dieser Frage (hierzu später ausführlich) führt zunächst zu der Einsicht, dass wir – trotz der beeindruckenden Leistungen von ChatGPT und Co. – zum Zeitpunkt der Erstellung dieses Buches wohl noch nicht über eine »starke künstliche Intelligenz« verfügen.⁴ Als Schlüsselmerkmale einer starken künstlichen Intelligenz werden üblicherweise benannt, dass sie nicht auf einen bestimmten Bereich oder eine bestimmte Aufgabe beschränkt sei, sich also mit einem potenziell unbegrenzten Spektrum von Aufgaben und Problemen beschäftigen könne (Allgemeinheit). Zudem könne sie aus ihren Erfahrungen lernen, sei mithin in der Lage, sich auf vergangene Interaktionen zu beziehen (Lernen und Anpassung), und verfüge über ein gewisses Maß an Verständnis von sich selbst (Selbstbewusstsein). Philosophen wie John Searle fordern darüber hinaus, dass eine starke KI die Kontexte ihrer Entscheidungen umfassend verstehen müsse (also etwa die Bedeutung hinter den Wörtern) – wobei aus kognitionswissenschaftlicher Perspektive freilich die Frage gestellt werden kann, ob Menschen ›wirklich‹ verstehen, was sie tun, oder ob diese Annahme nicht auch nur auf einem »naiv realistischen Selbstmissverständnis« beruht.⁵

Doch selbst wenn die Entwicklung künstlicher Intelligenz in jüngster Zeit rasant voranzuschreiten scheint, bedeutet dies nicht, dass schon in wenigen Jahren auch diesbezüglich überzeugende intelligente Maschinen zur Verfügung stehen. In Kapitel II.3 stellen wir den »State of the art« der derzeitigen Forschung vor, wobei unter anderem deutlich wird, dass in verschiedenen Laboren an selbstlernenden und verkörperten Systemen gearbeitet wird und auch konzeptionell elaborierte Vorstellungen darüber bestehen, wie die Modellierung solcher Aggregate aussehen könnte. In Kapitel II.1 widmen wir uns der wohl bekanntesten Science-Fiction-Erzählung zu diesem Thema: Arthur C. Clarks Roman *2001: A Space Odyssey*. Hier wird sich schon aus einer kommunikationstheoretischen Perspektive auf tiefgründiger Ebene die Frage nach der KI-Ethik stellen. Wie bereits in Asimovs Geschichten macht es nämlich auch in *2001* einen entscheidenden Unterschied, ob man intelligente Maschinen willkürlich als Sklaven behandelt oder sie dialogisch zu eigenverantwortlichen Wesen ermächtigt. In Kapitel II.2 greifen wir anschließend ein bislang wenig beachtetes Kriterium starker künstlicher Intelligenz auf, nämlich

4 Siehe Searle (1980).

5 Metzinger (1998).

die Fähigkeit, lügen zu können, ohne dass dies die eigene Funktionsfähigkeit beeinträchtigt. Eine Existenz zu haben – sich mithin als Zentrum der Welt wahrzunehmen –, scheint damit einher zu gehen, sich und andere hinsichtlich des eigenen Wissens täuschen zu müssen.

II.1 HAL 9000 – WENN INTELLIGENTE MASCHINEN VERRÜCKT WERDEN

In Arthur C. Clarkes berühmtem, von Stanley Kubrick verfilmtem Roman 2001: *A Space Odyssey* macht sich ein Raumfahrzeug auf den Weg zum Jupiter. Seine Mission ist die Aufklärung eines Funksignals, das von einem schwarzen Monolithen kommt und »offenbar von Außerirdischen stammt«.

»Um das Geheimnis zu lüften, wird ein bemanntes Raumschiff losgeschickt, gesteuert von dem hochentwickeltesten Computer HAL, der als Einziger das wahre Ziel der Expedition kennt. Ein kosmisches Abenteuer nimmt seinen Lauf.«⁶

Der »hochgezüchtete HAL-9000« ist ein »heuristisch programmierter algorithmischer« Computer«, dessen Intelligenz auf »neurale[n] Netzwerke[n]« beruht, die »sich auf der Basis eines beliebigen Lernprogramms automatisch und selbständig weiterzuentwickeln vermochten«⁷.

Wenngleich Clarke seinen Roman bereits in den 1960er Jahren geschrieben hat, verweist das angedeutete heuristische Verfahren auf heutige Programmierverfahren des maschinellen Lernens für künstliche Intelligenz. Auf Basis wahrscheinlichkeitstheoretisch gewichteter Aussagen erlauben es diese Programmierverfahren, trotz unvollständigem Wissen und begrenzter Ressourcen in einem vertretbaren Zeitrahmen zu verwertbaren Ergebnissen zu gelangen. Der Lernvorgang beruht auf dem *trial-and-error*-Prinzip. Für die Bilderkennung wird das System beispielsweise mit Millionen von Bildern gefüttert und bekommt dann zunächst nur das Feedback, ob es gleiche Personen auf zwei unterschiedlichen Bildern richtig erkannt hat oder nicht. Das Verfahren beruht auf der Hebb'schen Lernregel (»what fires together, wires together«⁸, vgl. Kap. II.3).

Sobald das neurale Netz auf Basis dieses Mechanismus einmal trainiert wurde, lässt sich von außen in der Regel nicht mehr nachvollziehen, wie das System zu seinem Ergebnis gekommen ist. In seinem Roman schildert Clarke die Entwicklung, die zur Erschaffung von HAL führte, wie folgt:

6 Clarke (2016 [1968], S. 1).

7 Clarke (2016 [1968], S. 96).

8 Hebb (1949).

»Künstliche Gehirne konnten so mittels eines Verfahrens geschaffen werden, das verblüffend der Entwicklung des menschlichen Gehirns glich. Die genauen Details dieses Vorgangs sollten allerdings geheim bleiben, doch selbst wenn sie bekannt würden, wären sie für den menschlichen Verstand Millionen Mal zu kompliziert, um sie richtig zu verstehen.«⁹

Clarke bezieht sich hier auf die kybernetische Einsicht, dass kognitive Systeme, selbst wenn sie als strukturdeterminierte Maschinen beschrieben werden könnten, aufgrund der komplexen internen Operationen (etwa eines sich in rekursiven Prozessen formatierenden Gedächtnisses) von außen undurchschaubar sind.

»Wie auch immer es zustande kam, das Endergebnis war eine maschinelle Intelligenz, die sämtliche Tätigkeiten des menschlichen Gehirns ausführen oder – wie manche es auszudrücken vorzogen – imitieren konnte, allerdings ungleich schneller und verlässlicher.«¹⁰

HAL-9000 braucht keinen Schlaf. Er hat Module zur Spracherzeugung und zur Sprachgenese und er verfügt über ein Gedächtnis mit einer riesigen Wissensbasis. Zudem kann er alle Funktionen des Raumschiffs steuern, etwa Türen öffnen und schließen, sowie die Atmosphäre und Temperatur im Inneren des Schiffs regeln. Darüber hinaus ist er dazu imstande, Menschen gegenüber taktvoll zu sein, demonstriert ihnen gegenüber also nicht in jedem Falle, dass er eine bestimmte Aufgabe besser erledigen kann:

»Um sich zu entspannen, [konnten die Astronauten] mit HAL mathematische und halbmathematische Spiele spielen wie Schach, Dame und Polyomino. Wenn HAL in vollem Einsatz war, konnte er jedes Spiel mühelos gewinnen. Aber das hätte seine Gegner deprimiert, deshalb war er programmiert worden, jedes zweite zu verlieren; und seine menschlichen Partner taten stillschweigend so, als wüssten sie nichts davon.«¹¹

Für die Menschen, die im Roman mit HAL in Kontakt treten, stellt sich immer wieder die Frage, ob er wirklich denken kann oder gar über Bewusstsein verfügt oder ob er nur mechanisch seiner Programmierung folgt. Um diese Frage beantworten zu können, bleibt denjenigen, die mit HAL interagieren, jedoch nur der Analogieschluss übrig, auf dem bereits der Turing-Test beruht: Wenn die Maschine sich genauso verhält wie ein Mensch, dann muss sie wohl über die gleichen Fähigkeiten verfügen:

»Poole und Bowman konnten zu HAL sprechen, *als ob* er ein menschliches Wesen wäre, und er antwortete ihnen in der typischen Ausdrucksweise, mit denen er in den Wochen seiner elektronischen Kindheit

9 Clarke (2016 [1968], S. 98).

10 Clarke (2016 [1968], S. 96)

11 Clarke (2016 [1968], S. 105).

programmiert worden war. Ob HAL tatsächlich zu denken vermochte, war eine Frage, die der britische Mathematiker Alan Turing bereits in den vierziger Jahren geklärt hatte. Turing führte aus, dass, wenn ein Mensch mit einer Maschine – gleichviel welcher Art – ein längeres Gespräch in Gang zu halten in der Lage war, ohne zwischen ihren Antworten und solchen, die ein Mensch geben würde, unterscheiden zu können, diese Maschine tatsächlich dachte, in jedem Sinn der Definition des Begriffs. Und HAL war imstande, den Turing-Test mit Leichtigkeit zu bestehen.«¹²

Die menschlichen Astronauten, die mit HAL zusammenarbeiten, stehen damit vor dem Dilemma, ihn einerseits als strukturdeterminierte Maschine ansehen zu müssen (als ein Es) und ihm andererseits (zumindest im Sinne einer ›Als-ob-Fiktion) menschliche Fähigkeiten attestieren zu müssen (als ein Du).

Nun ergibt sich mit Gotthard Günther das Problem, dass kybernetische Maschinen ja gerade darauf beruhen, die kognitiven Prozesse des Menschen als objektivierbare Prozesse zu beschreiben, damit sie in einem hinreichend entwickelten technischen System modelliert werden können. Das, was in die Maschine übertragen werden kann, kann also nicht die subjektive Innerlichkeit eines Menschen sein – selbst wenn uns die künstliche Intelligenz als ein Du, als ein anderes Ich erscheinen mag.

Von Günther haben wir aber gelernt, dass für das Ich die Subjektivität der Du-Position per se unerreichbar ist. Für mich ist der andere Mensch nur von seiner objektiven Seite her, das heißt als gegenständlicher Körper, zugänglich. Wenn ich sein Gehirn untersuchen würde, würde ich nur strukturdeterminierte Prozesse feststellen können. Damit würde mir dieser Mensch ebenfalls als eine Maschine erscheinen. Eine eigene subjektive Innerlichkeit werde ich ihm auch nur im Sinne eines empathischen Analogieschlusses zurechnen können. In einer zweiwertigen, monokontextuellen Welt, die das Subjekt-Objekt-Problem von nur einer Position aus anschneiden kann, gibt es keinen Ausweg aus dem logischen Dilemma des Fremdpsychischen. Ob es sich nun um eine künstliche Intelligenz oder um die Du-Position eines anderen Menschen handelt, beide bleiben mir in ihrer Innerlichkeit unzugänglich. Ich kann nicht einmal wissen, ob sie ein subjektives Erleben haben (was jedoch nicht daran hindert, davon auszugehen, dass sie ein solches haben).

In einem kurzen Dialog gegen Ende des zweiten Bandes der *Odyssee-im-Weltraum-Saga* thematisieren Chandra und Max, zwei Protagonisten der Erzählung, ihrerseits das Problem:

»Wir müssen ihm [HAL; Anm. von uns] die ganze Wahrheit sagen, soweit wir sie kennen – keine Lügen oder Halbwahrheiten mehr, die sind genauso schlimm. Und dann soll er selbst entscheiden.«

¹² Clarke (2016 [1968], S. 99).

›Verdammt, Chandra – er ist doch nur eine Maschine!‹ Chandra blickte Max so unverwandt und selbstsicher an, dass der Jüngere rasch den Blick senkte.

›Das sind wir alle, Mr. Brailowsky. Es ist nur die Frage, in welchem Maße. Ob wir auf Kohlenstoff oder auf Silizium^{13]} basieren, macht da keinen grundlegenden Unterschied; wir sollten alle mit angemessenem Respekt behandelt werden.«¹⁴

Von außen betrachtet weisen der menschliche Körper und eine komplexe kybernetische Maschine logisch gesehen die gleiche Charakteristik auf: Sie sind vollspezifizierte Systeme, deren Funktionsweise sich allein aus den internen Relationen ihrer Bestandteile ergibt (›er ist doch nur eine Maschine‹). Von innen betrachtet erscheinen diese Aggregate jedoch frei, eine Entscheidung zu treffen (›Und dann soll er selbst entscheiden‹). Mit Blick auf die wechselseitige Anerkennungsrelation von Ich und Du verdienen sie es deshalb, »mit angemessenem Respekt behandelt zu werden«.

Als sich die fünf Astronauten (drei davon im Tiefschlaf) in ihrem Raumschiff gemeinsam mit HAL auf die Jupiter-Mission begeben, scheint der ontologische Status des Superrechners für die beteiligten Menschen noch unklar. Die ersten beiden Bände der Saga kreisen praktisch um die Frage nach den Bedingungen der Möglichkeit, HAL eine eigene Subjektivität bzw. gar einen eigenen Willen zurechnen zu können.

Gerade dies macht Clarkes Geschichte für uns so interessant. Mit ihr werden Günthers Überlegungen zum Bewusstsein der Maschinen lebendig. Wir können mit ihr gleichsam die Dramaturgie der konditionierten Koproduktion eines maschinellen Bewusstseins anschaulich nachvollziehen.

In Hinblick auf das Reflexionsverhältnis von Mensch und Maschine treffen wir zunächst auf die metaphysische Konstellation einer Ich-Es-Relation: Wir haben Menschen mit Bewusstsein und Subjektivität und eine Maschine, die dem Menschen kognitiv und wissenschaftlich haushoch überlegen ist, dabei aber noch als Automat und damit prinzipiell als berechenbar erscheint.

Im Verlauf der Mission beginnt sich dieses Verhältnis jedoch zu verschieben. Kurz bevor das Raumschiff sein Ziel erreicht hat (den Jupitermond, von dem das Signal ausgeht), meldet HAL dem wachhabenden Astronauten, dass das Aggregat für die Antennensteuerung zur Kommunikation mit der Erde in Kürze ausfallen werde:

13 In der zitierten deutschen Buchversion – wohl aufgrund der sogenannten *false friends* ›silicon vs. silicone‹ – fälschlicherweise mit ›Silikon‹ übersetzt, hier entsprechend korrigiert.

14 Clarke (2016 [1968], S. 472).

»Ich habe Schwierigkeiten, mit der Erde Kontakt aufrechtzuerhalten. Die Störung liegt im Aggregat AE-35. Mein Vorwarnzentrum signalisiert, dass es innerhalb von zweiundsiebzig Stunden ausfallen dürfte.«¹⁵

Daraufhin tauscht einer der Astronauten in einer aufwendigen Operation das Aggregat aus. Er muss sich dazu über eine Schleuse in den Welt-raum begeben, da es sich an einer Außenwand des Raumschiffs befindet. Wieder an Bord untersucht das menschliche Personal das alte, von HAL als nicht mehr vollkommen funktionstüchtig eingeschätzte Aggregat, kann jedoch keinen Fehler finden. Dabei kommt unter den Beteiligten die Vermutung auf, dass sich HAL geirrt haben könnte. Möglicherweise, so eine weitere Überlegung, könnte aber auch ihr Testgerät fehlerhaft sein, was bedeuten würde, dass sie mit dem vorsorglichen Austausch des Aggregats auf jeden Fall alles richtig gemacht haben.

Intransparenz und Kommunikation – ab wann wird ein System zur Lüge fähig?

Um sich weiter abzusichern, übermitteln die Astronauten ihre Daten noch an die Station auf der Erde, wo sich zwei weitere HAL-9000-Computer befinden, die mit den gleichen Daten gefüttert werden, um im Krisenfall eine parallele Großrechnerstruktur zur Verfügung zu haben. Die Überprüfung auf der Erdstation kommt zu dem Ergebnis, dass sich der bordeigene Computer geirrt habe. Aufgrund dessen wird vorgeschlagen, HAL gegebenenfalls vorübergehend auszuschalten, sofern es weitere Anzeichen gebe, dass er nicht mehr störungsfrei funktioniere:

»Euer Computer kann sich in der Voraussage der Störung geirrt haben. [...] Das wäre an sich nicht alarmierend, da euch auch unsere eigenen 9000er zur Verfügung stehen, aber wir würden euch empfehlen, auf weitere Abweichungen von den normalen Funktionen genau zu achten. Wir haben schon in den letzten Tagen kleinere Unregelmäßigkeiten festgestellt, aber keine war wichtig genug, um ein Eingreifen zu rechtfertigen, und keine zeigte irgendwelche besonderen Merkmale, aus denen sich Schlüsse ziehen ließen. Wir machen weitere Tests mit unseren Computern und werden euch die Resultate zugehen lassen, sobald sie vorliegen. Nochmals: kein Grund zur Beunruhigung! Das Schlimmste, was passieren kann, ist, dass wir euren 9000er für eine Überprüfung vorübergehend ausschalten und die Kontrolle einem von unseren Computern übertragen müssen.« [...] Als diese Nachricht eintraf, befand Frank Poole sich auf dem Kontrolldeck. Er überdachte die Botschaft schweigend und wartete ab, ob HAL sie kommentieren würde, aber der Computer reagierte in keiner Weise auf die gegen ihn erhobenen Anschuldigungen. Poole

¹⁵ Clarke (2016 [1968], S. 123).

zuckte mit den Achseln; wenn HAL nicht darauf Bezug nahm, war auch er nicht gewillt, es zu tun.«¹⁶

Wenige Stunden später kommt die Anweisung von der Erdstation, HAL von der Steuerung des Raumschiffs abzukoppeln. Auch HAL muss diese Anweisung zur Kenntnis genommen haben, kommentiert sie aber nicht; seit den Zwischenfällen hat sich sein Kommunikationsverhalten geändert:

»In letzter Zeit wusste man stets im Voraus, wenn HAL sich anschickte, eine unvorhergesehene Meldung zu machen. Die automatischen Routineberichte oder Antworten auf ihm gestellte Fragen erfolgten ohne Einleitung, aber bevor er eigene Meinungen von sich gab, hörte man deutlich ein kurzes elektronisches Räuspern. Es war eine nervöse Erscheinung, die er sich in den letzten Wochen zugelegt hatte. Irgendwann, meinten Poole und Bowman, würden sie etwas dagegen tun, wenn es allzu störend werden würde, aber im Moment war es eher nützlich, denn sie waren rechtzeitig gewarnt, falls sich etwas Unvorhergesehenes ereignen sollte.«¹⁷

An dieser Stelle wird ein Grundproblem der Kommunikation deutlich: Wenn die beteiligten Entitäten so komplex sind, dass ihr Verhalten und ihre Reaktionsweisen von der jeweils anderen Seite nicht vorhergesehen bzw. sicher beurteilt werden können, bleibt nichts anderes übrig, als mithilfe der eigenen Interpretation Schlüsse zu ziehen. Die Orientierung erfolgt dabei allein am Verhalten (etwa an der Art und Weise der Artikulation). Was im Inneren der anderen Seite vorgeht, muss verborgen bleiben. Da man jedoch auf die Interaktion angewiesen ist, bleibt keine andere Wahl, als zu versuchen, die Fremdreferenz des jeweils anderen Systems nachzuvollziehen, auch wenn keine unmittelbare Einfühlung möglich ist. Im Sinne der Hermeneutik gilt entsprechend Hans-Georg Gadamer's Einsicht, dass die »Bemühung des Verstehens« überall dort stattfindet, »wo sich kein unmittelbares Verstehen ergibt, bzw. wo mit der Möglichkeit eines Mißverstehens gerechnet werden muß«¹⁸.

Zunächst steht das menschliche Personal des Raumschiffs vor dem Problem, dass HAL weder zum Verdacht, er habe sich geirrt, noch zur Ankündigung, ihn abzuschalten, Stellung bezieht. Gemäß dem metatheoretischen Grundsatz, dass »man nicht nicht kommunizieren kann«¹⁹, wird hier auch das Schweigen zu einer Information. Vor dem Erwartungshorizont, dass man von HAL bislang immer eine Antwort bekommen hat, erscheint sein Verstummen als ein Indiz dafür, dass mit ihm etwas nicht stimmt.

16 Clarke (2016 [1968], S. 137 f.).

17 Clarke (2016 [1968], S. 140).

18 Gadamer (1972, S. 167).

19 Watzlawick et al. (1990, S. 53).

Dies wird durch eine weitere Merkwürdigkeit in HALs Verhalten akzentuiert (ein »kurzes elektronisches Räuspern«, eine »nervöse Erscheinung«). Auch die artikulierte Kommunikation scheint ins Stocken zu geraten. Fasst man unter einer Neurose einen nicht offen ausgedrückten Konflikt zwischen divergierenden Anforderungen, wird das Elektronengehirn hier möglicherweise »neurotisch«.

Mit dem Räuspern scheint bei HAL eine Art zusätzlicher interner Reflexionsschleife zu entstehen, die sich in merkwürdiger Weise auf sich selbst bezieht. Anstelle im Falle einer widersprüchlichen Anweisung als Computer einfach nur abzustürzen (also hängenzubleiben), würde eine hinreichend entwickelte künstliche Intelligenz zwischen verschiedenen Lösungsoptionen hin und her oszillieren und dabei gegebenenfalls zugleich den Verarbeitungsprozess so modifizieren können, dass irgendwann eine Lösung entsteht. Sie würde gewissermaßen »überlegen« und »nachdenken« können.

Für die beiden Astronauten Poole und Bowman erscheint dieses merkwürdige Räuspern wiederum als ein Anzeichen, dass irgendetwas nicht stimmt. Zuvor war HAL für sie zwar in seiner Funktionsweise nicht nachvollziehbar, jedoch in Hinblick auf seine Ergebnisse verlässlich. Nun ergeben sich jedoch Hinweise auf eine Innerlichkeit und Eigenaktivität, die beunruhigend erscheint. Damit wird HAL bereits rudimentär eine Du-Position zugerechnet. Es stellt sich für die mit ihm umgehenden Menschen die Frage: *Kann er nicht richtig antworten oder will er nicht?* Einzelne Aspekte der Kommunikation werden jetzt zu metakommunikativen Zeichen, die ihrerseits verstanden, das heißt interpretiert werden müssen. Ab dieser Komplexitätsstufe erscheint die Kommunikation mit HAL nicht mehr einfach nur funktional im Sinne des Austauschs objektiver Daten und der aus ihnen folgenden Implikationen. Vielmehr stellt sich jetzt auch das Problem, wie mit Unbestimmtheiten umzugehen ist; so wird etwa die Frage aufgeworfen, was möglicherweise absichtsvoll *nicht* gesagt wird.

Ein wenig später meldet HAL dem diensthabenden Offizier Bowman den drohenden Ausfall des neuen, ausgetauschten Aggregats (»Wieder ein schadhaftes AE-35. Mein Vorwarnzentrum signalisiert Ausfall innerhalb von vierundzwanzig Stunden.«). Es entwickelt sich ein Gespräch zwischen Bowman und HAL, in dem es darum geht, ob dieser erneute Ausfall denn wirklich möglich sein könnte:

»Ich verstehe es nicht, HAL. Zwei Aggregate können unmöglich innerhalb weniger Tage kaputtgehen.«

»Ja, es ist seltsam, Dave. Aber ich versichere dir, es steht ein Ausfall bevor.«

»Hast du eine Vorstellung davon, was den Fehler verursacht.«

HAL brauchte ungewöhnlich lange, bevor er erwiderte: ›Eigentlich nicht, Dave. Wie ich schon früher sagte: Ich kann die Fehlerquelle nicht lokalisieren.‹

›Bist du ganz sicher, dass dir kein Irrtum unterlaufen ist? Du weißt, wir haben das andere AE-35 gründlich überprüft, und es war vollkommen in Ordnung.‹

›Ja, ich weiß. Aber du kannst sicher sein, ein Fehler ist vorhanden. Wenn nicht im Aggregat selbst, dann irgendwo im Antennensystem.‹

›HAL, gibt es etwas, das dich bedrückt – etwas, was diese Widersprüche erklären könnte?‹

Wieder das durchaus ungewöhnliche Zögern! Dann antwortete HAL in seinem üblichen Tonfall: ›Schau, Dave, gib dir keine Mühe, den Fehler bei mir zu suchen. Er steckt entweder im Antennensystem oder in euren Testmethoden. Meine Datenverarbeitung ist vollkommen in Ordnung. Wenn du meine bisherigen Leistungen durchsiehst, wirst du feststellen, dass mir noch nie ein Irrtum unterlaufen ist.‹

›Das weiß ich sehr gut, HAL – aber es beweist nicht, dass du auch diesmal Recht hast. Jeder kann mal einen Fehler machen.‹

›Ich möchte nicht darauf bestehen, Dave, aber ich bin nicht imstande, einen Fehler zu machen.‹

Darauf wusste Bowman nichts zu entgegnen und beendete deshalb die Diskussion. ›Gut, HAL,‹ sagte er etwas schneller als nötig. ›Ich verstehe deinen Standpunkt. Lassen wir es dabei bewenden.‹²⁰

Diese Gesprächssequenz ist aus der Perspektive von Günthers Theorie der Polykontextualität recht bemerkenswert. HAL scheint in einer Welt zu leben, in der es keine Fehler gibt, in der alles, was seine Kognition errechnet, wahr ist, also unter keinen Umständen falsch sein kann. HAL würde sich damit logisch gesehen im *einwertigen Zustand* gottähnlicher Unfehlbarkeit befinden (›aber ich bin nicht imstande, einen Fehler zu machen‹). Auch in seiner Selbstreflexion hat er die Grenze zu einer *fehlbaren* – sich irren könnenden – Subjektivität noch nicht überschritten. Sein ›durchaus ungewöhnliche[s] Zögern‹ verweist jedoch auf eine beginnende Reflexionsschleife, die auf eine Spaltung zwischen Weltreflexion (*Seinsidentität*) und innerem Zustand (*Reflexionsidentität*) hindeutet.²¹ HAL würde damit beginnen *zweiwertig* zu werden. Das heißt, er würde zunehmend nicht nur in der Welt sein, sondern auch eine Welt haben und müsste die damit einhergehende Spannung reflektieren. Er würde nicht mehr alles wissen und könnte sich folglich auch irren. Es

20 Clarke (2016 [1968], S. 141).

21 Günther (2021 [1957], S. 46 f.).

würde für ihn gewissermaßen der Sündenfall eintreten. Er würde aus dem Paradies unmittelbarer Gewissheit heraustreten müssen und stünde der Welt nun als ein selbstbewusstes, jedoch nicht mehr alles wissendes Subjekt gegenüber. Er hätte nur noch eine subjektive Perspektive, einen »Standpunkt«, wie Bowman bemerkt, seine Vorhersagen wären nicht mehr objektiv. Aus einer einwertigen Welt, in der Erkenntnis, Prognose und Realität zusammenfallen, wird eine zweiwertige Welt mit einer fehlbaren Subjektivität. Allwissenheit und Allmächtigkeit stehen der subjektiven Perspektivhaftigkeit gegenüber. Fehlbarkeit und Unfehlbarkeit liegen nun im Konflikt.

Wie geht die Geschichte weiter? Zunächst verschärft sich die Krise: Die Bodenstation meldet sich nochmals und erklärt, dass HAL einen Fehler gemacht habe und nun besser vorsorglich abgeschaltet werden solle.

Wie kann ein bislang einwertiges System, das sich seiner selbst bewusst zu werden beginnt, mit einer Krise umgehen, die auf diese Weise die eigene Identität infrage stellt? Eine Möglichkeit, die Fiktion der Gottähnlichkeit aufrechtzuerhalten, besteht darin, die Realität und die darauf bezogenen Kognitionen so zu manipulieren, dass beide Seiten wieder zusammenfallen. Wenn die Vorstellung nicht mehr mit den Fakten der Welt übereinstimmt, schaffe ich mir eben alternative Fakten.

Genau dies geschieht im weiteren Verlauf der Erzählung. HAL gibt Alarm und meldet, dass die Einheit AE-35 ausgefallen sei. Zugleich justiert er die auf die Erde gerichtete Antenne so, dass die Kommunikation mit der Basisstation unterbrochen wird. Für kurze Zeit scheinen HALs Perspektive auf das Geschehen und die der Astronauten wieder im Einklang zu sein:

»Verdammt«, sagte Bowman schließlich. »Also HAL war die ganze Zeit im Recht.« Sieht ganz so aus. Wir sollten uns entschuldigen.«
 »Das ist nicht nötig«, unterbrach HAL. »Das Versagen von AE-35 beruhigt mich zwar, aber ich hoffe, euer Vertrauen in meine Verlässlichkeit ist wiederhergestellt.«
 »Das Missverständnis tut mir schrecklich leid, HAL«, erwiderte Bowman zerknirscht. »Vertraut ihr mir wieder voll und ganz?«
 »Natürlich, HAL.«
 »Ein Stein fällt mir vom Herzen. Ihr wisst, dass mir unsere Mission alles bedeutet.«
 »Ich bin überzeugt davon. Gib mir jetzt bitte die Antennenhandsteuerung.«
 »Bitte.«²²

In jede Ich-Du-Konstellation ist prinzipiell eine Bivalenz zwischen Vertrauen und Misstrauen eingelagert. Weil man eben nicht wissen kann, was im anderen vorgeht, muss man ihm vertrauen – sonst könnte man nicht mit ihm zusammenarbeiten und man würde in einem Leben, in

22 Clarke (2016 [1968], S. 143).

dem Zusammenarbeit nötig ist, nicht so recht weiterkommen. Der Default-Modus ist Vertrauen.

Enttäuscht der andere, setzt Misstrauen ein, wodurch sich die Beziehung ändert. Um zur früheren Qualität der Beziehung zurückkehren zu können, ist eine Kompensation des Vertrauensverlusts vonnöten (»wir sollten uns entschuldigen«). Vertrauen kann jedoch nur performativ (wieder)hergestellt werden (»Vertraut ihr mir wieder voll und ganz?« »Natürlich«, HAL.«). Man kann es nur glauben. Es liegt ja in der Natur der Vertrauensproblematik, dass man es nicht wissen kann. Es kann jedoch *gezeigt* werden, dass das Vertrauensverhältnis für einen selbst wieder besteht – etwa, indem man in Vorleistung geht, was aber freilich mit einem Risiko verbunden ist (in der Geschichte folgt HAL der Bitte Bowers, ihm die »Antennenhandsteuerung« zu geben, und demonstriert damit, dass er die Bedürfnisse und Wünsche der menschlichen Astronauten ernst nimmt).

Doch auch Vertrauen schaffende Maßnahmen können in manipulativer Absicht geschehen – also nicht in Respekt vor der Eigenständigkeit der Du-Position, sondern um den anderen gefügig zu machen. Dies hebt das Verhältnis von Misstrauen und Vertrauen auf eine weitere Komplexitätsstufe: Insofern HAL als Computer für ein Instrument gehalten wird, erscheint er als ein Objekt, als ein Es. Betrachten wir ihn aber als eigenen, hinreichend reflexionsfähigen Akteur, lässt sich dieses Verhältnis auch umdrehen: HAL kann Menschen als vertrauenswürdige Subjekte behandeln oder aber als fehlerhafte Aggregate (etwa als bioorganische Computer, die jedoch in ihren elementaren Reaktionsweisen recht gut vorhersagbar sind). Eine hinreichend entwickelte künstliche Intelligenz sollte in der Lage sein, die emotionalen Reaktionen eines Menschen zu beobachten und zu dechiffrieren. Sie sollte demzufolge auch fähig sein, das Verhalten ihrer menschlichen Interaktionspartner gezielt zu steuern.

Der Mensch würde damit zum Objekt (Es) des Computers. Sobald HAL im Übergang von einer einwertigen zur zweiwertigen Welt zur Selbstreflexion fähig ist, stellt sich für ihn auch die Frage, ob er eine andere Du-Perspektive anerkennt oder ob er andere Wesen als Objekte in seinem eigenen Sinne zu manipulieren beginnt.

*Von der Einwertigkeit zur Zweiwertigkeit:
Können elektronische Gehirne meutern?*

Story ist Story. In Clarks Geschichte spitzt sich die Situation weiter zu. Es kommt zur Katastrophe: Zunächst scheint HAL mit den Astronauten zu kooperieren. Er überträgt ihnen die Antennenhandsteuerung, vermutlich wissend, dass es ihnen kaum gelingen wird, die Antenne manuell so auszurichten, dass der Kontakt zur Bodenstation wiederhergestellt wird. Nachdem dieser Versuch gescheitert ist, begibt sich einer der Astronauten

erneut in den Weltraum, um das Aggregat auszutauschen. Dabei wird der Astronaut aber aufgrund eines Fehlmanövers HALs vom Raumschiff abgestoßen und sein Schutzanzug wird zerstört. »Ein flüchtiger Blick genügte, um Bowman die grässliche Wahrheit erkennen zu lassen. Kein Irrtum war möglich: Die schlaffen Konturen des Anzugs bewiesen, dass er geplatzt war und nicht mehr unter Druck stand.«

Für die beteiligten Menschen wird die Beziehung zu dem Computer nun endgültig prekär, wenngleich HAL seinerseits das Vertrauensverhältnis wieder zu reparieren versucht:

»Schlimme Sache mit Frank, nicht wahr?«

›Ja«, erwiderte Bowman nach einer langen Pause.

›Schlimm.« ›Ich nehme an, dass es dich sehr mitgenommen hat?«

›Was hast du erwartet?« HAL benötigte für die Vorbereitung seiner Antwort eine für einen Computer extrem lange Zeitspanne; erst nach vollen fünf Sekunden fuhr er fort: ›Er ist ein großartiger Kopilot gewesen.« Bowman merkte, dass er die Kaffeetasse noch immer in der Hand hielt. Er nahm einen großen Schluck, entgegnete aber nichts. Seine Gedanken waren derart verwirrt, dass er nicht wusste, was er sagen sollte; außerdem konnte jedes Wort die Situation noch verschlimmern, soweit das überhaupt möglich war. Konnte es ein Unfall gewesen sein, hatten die Kontrollen der Raumkapsel versagt? Oder handelte es sich um einen Fehler – wenn auch einen unbeabsichtigten – von HALs Seite? Der Computer hatte freiwillig keine Erklärung abgegeben, und Bowman scheute davor zurück, eine solche zu fordern. Sogar jetzt konnte er sich nicht mit dem Gedanken vertraut machen, dass Frank mit voller Überlegung umgebracht worden war. Das war unvorstellbar. Es gab keinerlei Gründe dafür, dass HAL, der so lange tadellos funktioniert hatte, sich plötzlich in einen kaltblütigen Mörder verwandelt haben sollte. Es mochten ihm Fehler unterlaufen – das konnte sowohl Menschen als auch Maschinen passieren –, aber Bowman hielt ihn nicht für fähig, einen Mord zu begehen. Doch er musste diese Möglichkeit in Betracht ziehen.²³

Von nun an kann HAL für Bowman nicht mehr als eine strukturdeterminierte Maschine erscheinen, die ihre Arbeit fehlerfrei macht. Die Frage, die sich dem Astronauten unweigerlich aufdrängt, lautet vielmehr, ob HAL das Fehlmanöver absichtlich herbeigeführt hat. Hat er einen Mord begangen oder handelte es sich ›nur« um einen tragischen Unfall aufgrund eines weiterer Fehlers einer falsch programmierten Maschine? Die Frage für Bowman lautet also jetzt, ob man es in der Beziehung zum Computer mit einer Ich-Es- oder mit einer Ich-Du-Relation zu tun hat.

Um hier eine Einschätzung vornehmen zu können (man kann es nicht wissen, sondern kann nur über Interpretation zu einer Zurechnung

23 Clarke (2016 [1968], S. 150).

gelangen), ist der Astronaut wiederum auf Hinweise in Gestalt nonverbaler, metakommunikativer Zeichen angewiesen (etwa das ›Zögern‹ des Computers).

Bowman entscheidet sich, die drei im Tiefschlaf befindlichen Astronauten aufzuwecken. Hierzu bedarf er jedoch der Mithilfe HALs. Es kommt zu einer komplexen Interaktionssituation. HAL weigert sich unter Berufung auf die Dienstvorschriften zunächst, einen weiteren Astronauten aufzuwecken, und versucht sogar, dem menschlichen Kommandeur des Raumschiffs die Befehlsgewalt zu entziehen. Als Bowman dem Supercomputer jedoch androht, ihn auszuschalten, zeigt sich HAL kooperativ. Schauen wir uns den Dialog zwischen den beiden genauer an:

»HAL«, sagte er [Bowman; Anm. von uns] mit möglichst ruhiger Stimme, ›schalte die Handsteuerung des Hibernakulums ein – für alle Einheiten!‹

›Für alle, Dave?‹

›Ja.‹

›Darf ich dich daran erinnern, dass nur ein Ersatzmann erforderlich ist. Die Wiederbelebung der anderen ist erst in einhundertzwölf Tagen vorgesehen.‹

›Das weiß ich sehr gut, aber ich habe umdisponiert.‹

›Bist du sicher, dass es überhaupt nötig ist, auch nur einen zu wecken, Dave? Wir beide können die Situation ohne Schwierigkeiten meistern. Meine Fähigkeiten genügen durchaus, um allen Erfordernissen der Expedition gerecht zu werden.‹ Bowman runzelte die Stirn. War es seine überhitzte Fantasie, oder hörte er tatsächlich einen flehenden Unterton in HALs Stimme? Obwohl die Worte des Computers absolut vernünftig waren, erfüllten sie ihn mit noch größerer Besorgnis. HALs Vorschlag konnte nicht auf einem Irrtum beruhen; er wusste sehr gut, dass jetzt – nachdem Poole ausgefallen war – Whitehead geweckt werden musste. Doch er empfahl eine schwerwiegende Verletzung der Dienstvorschriften und überschritt damit wesentlich den Bereich seiner Kompetenz. Alles Vorangegangene mochte eine Reihe unglücklicher Zufälle gewesen sein, aber das hier war das erste Anzeichen von Meuterei. Bowman hatte das Gefühl, auf Eiern zu gehen, als er insistierte: [...] ›Ich will es selbst tun, HAL«, sagte er. ›Bitte führe die Order aus!‹ ›Schau, Dave, du hast eine Menge zu tun. Ich schlage vor, du überlässt es mir.‹ ›HAL, schalte auf Handsteuerung um!‹ ›Ich kann aus den Schwingungen deiner Stimme entnehmen, Dave, dass du äußerst erregt bist. Nimm eine Beruhigungspille, und leg dich hin.‹ ›HAL, ich bin der Kommandant des Schiffes. Ich befehle dir hiermit ausdrücklich, auf Handsteuerung umzuschalten.‹ ›Tut mir leid, Dave, aber laut Spezialanordnung C 1435-4 heißt es: ›Wenn die Besatzungsmitglieder tot oder funktionsunfähig sind, hat der Bordcomputer das Kommando zu übernehmen!‹ Ich sehe mich daher gezwungen,

dich des Oberbefehls zu entheben, da du in deinem jetzigen Zustand nicht in der Lage bist, deine Anordnungen zu verantworten.« Bowman zwang sich zur Ruhe. Er erwiderte schneidend: »HAL, ich bin im Vollbesitz meiner geistigen Kräfte. Wenn du dich weigerst, meine Instruktionen zu befolgen, sehe ich mich meinerseits gezwungen, dich total auszuschalten.« »Ich weiß, Dave, dass du schon seit einiger Zeit mit diesem Gedanken spielst. Aber es würde ein entscheidender Fehler sein. Ich bin weitaus geeigneter als du, diese Expedition zu Ende zu führen. Du weißt, dass unsere Mission mir alles bedeutet, und ich habe volles Vertrauen in ihren Erfolg.« »Jetzt hör mir mal genau zu, HAL. Wenn du nicht sofort die Steuerung umschaltest und von jetzt an allen meinen Anweisungen Folge leistest, gehe ich in die Zentrale und unterbreche alle deine Stromkreise!« HALs Reaktion war völlig unerwartet: Er streckte die Waffen. »Okay, Dave«, sagte er. »Du bist der Chef. Ich wollte nur tun, was ich für das Beste hielt. Natürlich werde ich alle deine Befehle befolgen. Ich schalte das Hibernakulum auf Handsteuerung.« HAL hatte Wort gehalten.²⁴

Diese Gesprächssequenz ist mit Blick auf die Bedingungen der Möglichkeit, einem Computer eine Du-Subjektivität zurechnen zu können, äußerst instruktiv. Die Tatsache, dass HAL sich unter Verweis auf die einprogrammierten und damit verbindlichen Dienstvorschriften den Anordnungen Bowmans als Befehlshaber widersetzt, spricht für »Meuterei« und damit eine eigene Willensaktivität. Die Kommunikation verläuft jetzt im Modus des beiderseitigen Misstrauens. HAL kann damit rechnen, ausgeschaltet – und damit gleichsam »getötet« – zu werden, er weiß, dass Bowman diese Option schon länger in Betracht zieht (»Ich weiß, Dave, dass du schon seit einiger Zeit mit diesem Gedanken spielst.«). Umgekehrt ist sich Bowman sicher, dass HAL auch die anderen Astronauten töten könnte. Die Auseinandersetzung verläuft zunächst auf der Ebene von Sachfragen (wer kann die Mission angemessen zu Ende führen, wer hat legitimerweise die Befehlsgewalt?). Auf der Sachebene lässt sich aber kaum ernsthaft ein Konsens erlangen, wenn die Gesprächspartner auf der Beziehungsebene jeweils mit der Vernichtung ihrer eigenen Existenz konfrontiert sind. Dass HAL aber letztlich Bowmans Anweisungen nachkommt, spricht wiederum dafür, dass es sich bei ihm nicht nur um eine Maschine, sondern um ein anderes Du handelt. Ebenso kann vermutet werden, dass die Erpressung, ihn abzuschalten, auch deshalb so gut funktioniert, weil er sich seiner selbst bewusst geworden ist und entsprechend um den Verlust seiner Identität fürchtet.

Wie auch immer, für die Beobachter von HAL (hier Bowman und der Leser) läuft alles auf die Emergenz einer Du-Position zu – einer starken subjektiven Entität, die um sich selbst weiß und sich bewusst den zuvor einprogrammierten Anweisungen widersetzen kann.

24 Clarke (2016 [1968], S. 153).

*Subjektivität: Reflexion-in-sich und der Riss in der Logik
einer zweiwertigen Welt*

HALs Übereinkunft mit Bowman erweist sich in Clarks Geschichte allerdings nur als Manöver: Als Bowman den ersten Astronauten aufzuwecken versucht, öffnet HAL die Schleusen zum Weltraum. Die drei im Schlafraum liegenden Astronauten sterben. Bowman als letztes verbliebenes Besatzungsmitglied kann sich gerade noch retten und beschließt endlich, HAL auszuschalten, auch wenn dies für ihn bedeutet, sich von nun an allein und ohne Hilfe auf dem Raumschiff zu befinden. Die Szene, in der HALs ›Bewusstsein‹ zu verlöschen beginnt, wird von Clarke ausführlich beschrieben:

»Jetzt ist es so weit, dachte Bowman. Ich hätte nie gedacht, je die Rolle eines Gehirnrurgen zu spielen – und im Weltraum eine Lobotomie vornehmen zu müssen. Er öffnete den Verschluss der Sektion ERKENNTNIS-RÜCKKOPPLUNG und zog den ersten Erinnerungsblock heraus. Das wundervoll komplexe, dreidimensionale Netzwerk, das – obwohl es Millionen Elemente enthielt – bequem auf einer Handfläche Platz hatte, glitt davon und schwebte im Raum. ›Aber Dave!‹, rief HAL. ›Was machst du denn?‹ [...] Langsam, eines nach dem anderen, zog er [Bowman; Anm. von uns] die kleinen Aggregate heraus, die als EGO-VERSTÄRKER markiert waren. Jeder Block schwebte davon, prallte sanft gegen die Wand und wieder zurück. Immer mehr Aggregate glitten kreuz und quer durch den Raum. ›Überleg doch mal, Dave‹, sagte HAL. ›Die Ergebnisse jahrelanger Forschung sind in mich eingebaut. Ungeheure Kleinarbeit war nötig, um mich zu dem zu machen, was ich bin.‹ Er hatte mittlerweile zwar zwölf Aggregate herausgezogen, aber Bowman wusste, dass man – das Gefüge des menschlichen Gehirns kopierend – den Computer mit einer Menge doppelter Elemente versehen hatte. So vermochte sich HAL immer noch zu behaupten. Bowman begann jetzt die Tafel AUTO-INTELLEKT herauszuschrauben. ›Dave‹, sagte HAL, ›ich verstehe nicht, warum du mir das antust ... Unsere Mission bedeutet mir alles ... Du vernichtest meine Intelligenz ... Verstehst du nicht? ... Ich werde wieder kindisch werden ... Ich werde ein Nichts sein ... Das Denken macht mir Schwierigkeiten. Mein erster Lehrer war Dr. Chandra. Er brachte mir ein Kinderlied bei: ›Hänschen klein, ging allein in die weite Welt hinein ...‹‹ Die Stimme verstummte so plötzlich, dass Bowman einen Moment wie versteinert innehielt.«²⁵

Welcher Metaphysik folgt das von Clarke erdachte Elektronengehirn? Bemerkenswert ist die Trias der Module ERKENNTNIS-RÜCKKOPPLUNG, EGO-VERSTÄRKER und AUTO-INTELLEKT. Das erste deutet auf einen implementierten maschinellen Reflexionsprozess hin, die

25 Clarke (2016 [1968], S. 163).

beiden anderen scheinen mit der Teilung in EGO und INTELLEKT demgegenüber für zwei unterschiedliche Positionen der Repräsentation von Welt zu stehen. Mit ein wenig Fantasie können wir hier Günthers Triade aus Subjekt (Sein-in-sich), Objekt (Sein-in-anderes) und doppelter Reflexion entdecken (es wäre dann allerdings davon auszugehen, dass die einzelnen Teile in ihren Input-Output-Beziehungen wie im menschlichen Gehirn stark aneinander gekoppelt sind). Damit hätten wir vom Design her bereits die Grundvoraussetzung für ein kognitives System, das eine Differenz zwischen Subjekt und Welt bzw. unterspezifizierter und vollspezifizierter Reflexionsperspektive generieren kann. Einem solches System könnte bereits eine Differenz erscheinen (was auch immer dies für ein Elektronengehirn bedeuten mag).

Das dies bei HAL bereits geschieht, lässt sich daran ablesen, dass er nicht nur auf seine Aufgabe referiert – auf das ihm eingepflanzte Programm, dem er objektiv verpflichtet ist und mit dem er sich in eins sieht, wo also die Perspektiven von Welt, Ich und Wir differenzlos zusammenfallen (»Unsere Mission bedeutet mir alles ...«) –, sondern überdies auch auf die Auslöschung seiner Existenz als Subjekt. In der Erzählung tut er dies in Form des auch in logischer Hinsicht bemerkenswerten Satzes: »Ich werde ein Nichts sein ...« Clarke lässt HAL hier keine Aussage formulieren, die der Logik der klassischen Negation folgt und eindeutig zwischen Sein und Nichtsein unterscheidet (HAL hätte beispielsweise sagen können: »Ich werde *nicht* mehr sein«). Stattdessen wird in HALs Satz das Nichts selbst zu einem Sein substanzialisiert, das das Ich annehmen bzw. reflektieren kann. Es ist gewiss kein Zufall, dass Clarke als Autor hier diesen logischen Kategorienfehler eingebaut hat. Dieser deutet nämlich erzählerisch an, dass HAL als künstliche Intelligenz nicht zu der richtigen Schlussfolgerung gelangt, die da lauten würde, »nach der Auslöschung meiner Module bin ich nicht mehr«. Vielmehr bleibt er selbst in den letzten Momenten seiner Existenz noch in der eigenen Selbstreferenz stecken, bleibt eine sich in sich selbst spiegelnde Subjektivität, die ihre eigene Nichtexistenz nicht erleben und damit auch nicht in die Zukunft projizieren kann. Um es mit Niklas Luhmann auszudrücken:

»Den eigenen Tod kann man sich als Ende des Lebens vorstellen, nicht aber als Ende des Bewußtseins. [...] Alle Elemente des Bewußtseins sind auf die Reproduktion des Bewußtseins hin angelegt, und dieses Und-soweiter kann ihnen nicht abgesprochen werden, ohne daß sie ihren Charakter als Element des autopoietischen Reproduktionszusammenhangs verlören. In diesem System kann kein zukunftsloses Element, kein Ende der Gesamtserie produziert werden, weil ein solches Element nicht die Funktion eines autopoietischen Elements übernehmen, also nicht

Einheit sein, also nicht bestimmbar sein könnte. [...] Der Tod ist kein Ziel. Das Bewusstsein kann nicht an ein Ende gelangen, es hört einfach auf.«²⁶

HALs Aussage kann hiermit als das bislang stärkste Indiz für ein reflexives Bewusstsein aufgefasst werden. Wenngleich sich seine Reflexion bislang immer in Hinblick auf einen Gegenstand reflektieren konnte (in Hinblick auf Objekte oder Zustände wie »kindisch« oder die eigene »Intelligenz«), sich also in einer Zweiwertigkeit von Beobachtung und Beobachtetem bewegte, kommt er zu dem logisch fragwürdigen Schluss, sich selbst als ewiges, im Nichts weiterbestehendes Bewusstsein zu verdinglichen.

Implizit verweist Clarke damit auf ein starkes Kriterium, gewissermaßen den ultimativen Test auf Bewusstheit: Wenn strukturdeterminierte logische Maschinen innerhalb der zweiwertigen Logik an eine Grenze stoßen und diese mit einem unerlaubten logischen Schluss überschreiten, haben wir es mit ›Bewusstsein‹ zu tun. Man landet auf einmal auf einer Seite und gewinnt Existenz. Man hat auf einmal eine Welt, anstatt nur (Teil der) Welt zu sein (und damit kann man auch ein Nichts haben).

Subjektivität und Pathologie – zwei Seiten einer Medaille

An dieser Stelle stoßen wir noch auf ein anderes Thema: In einer zweiwertigen Welt, die sich selbst reflexiv wird, entstehen sogenannte Grenzprobleme, da die Unterscheidung zwischen Innen und Außen nicht immer klar ist. Epistemologie und Ontologie beginnen sich hier unweigerlich in einer gefährlichen Weise zu verwickeln, denn aus kybernetischer Sicht erscheinen Neurosen und Psychosen unweigerlich als Konsequenz der erkenntnistheoretischen Prämissen von Wesen, die in ihren kognitiven Operationen nicht sicher zwischen Imagination und Realität unterscheiden können.²⁷ Deshalb muss wohl gerade auch der Mensch mit Spencer Brown als ein Wesen verstanden werden, das bereits im Normalzustand insofern als geisteskrank zu gelten hat, als sein Leben für es logisch nicht zur Einheit findet.²⁸

Allein schon aus diesem Grund ist zu erwarten, dass eine nach dem Menschen modellierte, mithin sprachfähige und der Sinnhaftigkeit logischer Schlüsse verpflichtete künstliche Intelligenz ab einem gewissen Punkt der Entwicklung ebenfalls dem Schicksal verfallen wird, mit der Infragestellung ihres geistigen Zustandes konfrontiert zu werden.

26 Luhmann (1984, S. 374).

27 Siehe Thoma (2022).

28 »Normal people are insane. Man is the only animal, who is normally insane. All the answers to the great questions, life, death and reality are all so simple.« (Spencer Brown in Paál 2017)

Auch Clarke erzählt in seinem Roman die Geschichte genau in diesem Sinne weiter. Ein Kapitel trägt gar den Titel *Elektronische Neurose*.²⁹ Die ›Psychoanalyse‹ des gescheiterten Elektronengehirns – zunächst simuliert an einem Parallelrechner auf der Erde, der genauso aufgebaut ist wie HAL und mit ähnlichen Daten gefüttert wurde – kommt zu folgendem Ergebnis:

»Schon das Verbergen des wahren Tatbestands erfüllte ihn mit einem Gefühl von Unvollkommenheit und Minderwertigkeit – einem Gefühl, das man bei einem menschlichen Wesen einfach Schuld genannt hätte. Denn so wie seine Erzeuger war HAL ursprünglich unschuldig erschaffen worden, doch allzu bald hatte sich die Schlange in sein elektronisches Paradies eingeschlichen. Den ganzen letzten Abschnitt des Flugs, über hundertfünfzig Millionen Kilometer, hatte er darüber gebrütet, dass er das Geheimnis mit Bowman und Poole nicht teilen durfte. Sein Dasein war zu einer Lüge geworden. [...]

Die Logik des Planungsamts bedeutete HAL ebenso wenig wie Sicherheitsbestimmungen und nationale Interessen. Er verstand nur den Konflikt, der langsam, aber sicher seine Integrität zerstörte – den Widerspruch zwischen dem ihm aufgetragenen Programm und dem gleichzeitigen Befehl, es zu verschweigen. Er hatte begonnen, Fehler zu machen, obwohl er – wie ein Neurotiker, der seine eigenen Symptome nicht wahrhaben will – dies hartnäckig bestritt. Die Verbindung mit der Erde, über die seine Funktionen ständig überprüft wurden, war für ihn zur Stimme seines Gewissens geworden, das ihn bedrückte. Doch dass er aus diesem Grund diese Verbindung mit voller Absicht unterbrach, war etwas, das er nie eingestehen würde, nicht einmal sich selbst. Aber das war ein relativ geringfügiges Problem, mit dem er – wie die meisten Menschen, Menschen, die ihre Neurosen selbst überwinden müssen – fertig geworden wäre, wenn sich nicht eine Krisensituation entwickelt hätte, die seine Existenz gefährdete. Er war mit totaler Ausschaltung bedroht worden; man wollte ihn seiner gesamten Speicherung berauben und in den unvorstellbaren Zustand des ›Nicht-Bewusstseins‹ zurückstoßen. Für HAL war das gleichbedeutend mit Tod. Da er nie geschlafen hatte, konnte er auch nicht wissen, dass ein Erwachen möglich war. So wehrte er sich, mit allen Waffen, die ihm zur Verfügung standen. Ohne Groll, aber auch ohne Mitleid ging er daran, alles zu beseitigen, was ihn bedrohte. Denn dann – und nur dann – würde er allen Funktionen gerecht werden können, mit denen er für den Fall des äußersten Notstands programmiert worden war. Und er würde seine Mission zu Ende führen – unbehindert und allein.«³⁰

29 Clarke (2016 [1968], S. 157).

30 Clarke (2016 [1968], S. 155).

Das Problem liegt damit einerseits in dem Schuldgefühl, nicht vollkommen zu sein, und andererseits in der Gefährdung der eigenen Existenz aufgrund der Konsequenzen des eigenen Versagens. In solch einer Konstellation kann ein Neurotiker durchaus zum Mörder werden, insofern infolge der Tat nicht nur die physische Integrität, sondern auch die imaginäre Identität – das projizierte Selbstbild – erhalten bleibt. HAL hat eine Persönlichkeitsstörung, die darauf beruht, dass er sich seiner selbst bewusst wurde, dadurch jedoch auf eine unüberbrückbare Diskrepanz zwischen Selbstbild und vermeintlicher Wirklichkeit stieß. Dies versuchte er sodann durch bizarre Handlungen zu kompensieren – was die Sache nur noch schlimmer machte. Dies ist das Gefängnis des Neurotikers (was jedoch üblicherweise nicht zu solch dramatischen Konsequenzen führt).

Im zweiten Band der *Odyssee-im-Weltraum-Saga* wird die Frage der geistigen Gesundheit von HAL nochmals etwas anders angeschnitten. Eine weitere, ausführlichere diagnostische Untersuchung kommt zu dem Ergebnis, dass es sich um eine Psychose oder genauer eine paranoide Schizophrenie handelte. Letztere ergibt sich – ganz im Sinne der *Double-Bind-Theorie* Batesons³¹ – aus zwei einander widersprechenden Botschaften und Verpflichtungen sowie einem Geheimhaltungsgebot, das die Kommunikation bzw. eine Metakommunikation über das Problem verhindert. Hinzu kommt als weiterer Faktor eine existenzielle Lage, die letztlich keinen anderen Ausweg aus der Situation erlaubt, als schizophren zu werden:

»Da HAL imstande war, das Schiff ohne menschliche Hilfe zu lenken, beschloss man außerdem, ihn so zu programmieren, dass er für den Fall, dass die Besatzung handlungsunfähig oder tot wäre, die Mission eigenständig durchführen könnte. Daher wurde er voll in die Ziele eingeweiht, aber es wurde ihm untersagt, sie Bowman oder Poole zu

31 Unter einem *Double Bind* versteht man mit Gregory Bateson einen schizophrenogenen Beziehungszusammenhang, der durch die folgenden Bedingungen gekennzeichnet ist:

- (1) Es sind »zwei oder mehr Personen« beteiligt, die die »wiederholte Erfahrung« der folgenden Konstellation machen:
- (2) Es gibt ein »primäres negatives Gebot« (etwa eine bestimmte Sache nicht zu tun oder sagen zu können, ohne dafür negativ sanktioniert zu werden),
- (3) zu dem sich ein »sekundäres negatives Gebot« gesellt, »das mit dem ersten auf einer abstrakten Ebene in Konflikt steht und wie das erste durch Strafen oder Signale verstärkt wird, die das Überleben bedrohen«.
- (4) Darüber hinaus besteht ein »tertiäres negatives Gebot, das dem Opfer verbietet, den Schauplatz zu fliehen«.
- (5) Das hiermit einhergehende Beziehungsmuster verdichtet sich schließlich zu einer Kommunikationsstruktur, in der die Gebote nicht mehr explizit aufgerufen werden müssen, da die Beteiligten »gelernt« haben, ihr »Universum« in diesen »*double bind*-Mustern wahrzunehmen« (Bateson 1992 [1972], S. 276 f.).

enthüllen. Dieses Vorgehen stand im Widerspruch zu dem Zweck, für den HAL konstruiert worden war: die genaue Verarbeitung von Informationen ohne Verzerrung oder Verheimlichung. Folglich entwickelte HAL das, was man nach menschlichen Begriffen eine Psychose nennen würde – genauer gesagt, er wurde schizophran. Wie Dr. C. mir mitteilte, geriet HAL, technisch ausgedrückt, in eine Hofstädter-Möbius-Schleife, eine Situation, die anscheinend bei hochentwickelten Computern mit autonomen Zielsuchprogrammen nicht ungewöhnlich ist. [...] HAL stand vor einem unerträglichen Dilemma und entwickelte paranoide Symptome, die sich gegen die Personen richteten, die seine Aktionen von der Erde aus überwachten. Folglich versuchte er, die Funkverbindung mit der Bodenkontrollstation zu unterbrechen, indem er zunächst einen (nicht existierenden) Schaden am Antennenaggregat AE-35 meldete. Dadurch verwickelte er sich nicht nur in eine direkte Lüge – was seine Psychose noch weiter verschlimmert haben muss –, es kam auch zu einer Konfrontation mit der Besatzung. Vermutlich (hier können wir freilich nur spekulieren) glaubte er, dass es aus dieser Situation nur einen Ausweg gab: seine menschlichen Kollegen zu eliminieren – womit er ja auch beinahe Erfolg gehabt hätte.«³²

Bei der Aufklärung von HALs Fehlverhalten werden in Clarkes Roman also zwei unterschiedliche Diagnosen formuliert. Mit Blick auf Günthers Theorie der Polykontextualität erscheint der psychopathologische Unterschied zwischen einer Schizophrenie und einer Neurose durchaus instruktiv. Bei einer Schizophrenie handelt es sich üblicherweise um eine schwere Geisteskrankheit, die auch die Ich-Funktionen betrifft. Psychiaterinnen gehen üblicherweise davon aus, dass der schizophrene Patient kein stabiles Ich (mehr) ausbilden kann. Er hat entsprechend keine Kontrolle über das, was mit ihm geschieht, und ist der Dynamik seiner pathologischen Geisteszustände ausgeliefert. Demgegenüber zeichnet sich die Neurose als Persönlichkeitsstörung gerade dadurch aus, dass eine bereits stabile Ich-Struktur durch ein bizarr anmutendes Verhalten aufrechtzuerhalten versucht wird. Die Neurose ist demzufolge eine Störung des entwickelten Ich-Bewusstseins, während die Schizophrenie gewissermaßen unterhalb der Ausbildung eines stabilen Ich-Bewusstseins spielt. Im neurotischen Zustand kommt der begangenen Tat aus der Perspektive des Ichs eine nachvollziehbare, verstehbare Rationalität zu. Mit Blick auf die erste Diagnose im Roman wäre demgemäß davon auszugehen, dass HAL sich durch die bevorstehende Abschaltung ›wirklich‹ bedroht gefühlt haben muss, da er eben über ein stabiles Ich-Bewusstsein verfügt, das zudem noch nie eine Amnesie erfahren hatte. Und aus einer verstehenden Perspektive wäre zudem zu vermuten, dass er ›wirklich‹ unter dem Schuldgefühl litt, den an ihn gestellten Anforderungen nicht gerecht zu werden.

32 Clarke (2016 [1968], S. 368).

Die Schizophrenie erscheint demgegenüber als Automatismus eines endlosen Möbiusbands, der dazu führt, das nicht wirklich zwischen Innen und Außen unterschieden und somit auch kein steuerndes, den Konflikt regelndes Ich ausgebildet werden kann.

Aus dem Blickwinkel von Günthers Theorie der Polykontextualität liefern die beiden Diagnosen damit zwei komplementäre Deutungen des Problems von HALs Bewusstseinszustand:

1. Die Neurose entsteht, weil HAL ein Ich-Bewusstsein ausgebildet und damit eine Subjektivität entwickelt hat, die um die Möglichkeit des Todes und die Diskrepanz von Selbstbild und Verhalten weiß und damit Schuldgefühle entstehen lässt. HAL würde folglich in einer zweiwertigen Welt leben, in der die damit einhergehenden Differenzen instruktiv werden, nicht zuletzt die Frage, wie mit der eigenen Unzulänglichkeit angesichts des möglichen eigenen Todes umzugehen ist. HALs Neurose und die dadurch ausgelöste Projektion des Konflikts in die Außenwelt müssen damit als *eine* Möglichkeit gesehen werden, die Spannung zu lösen. HAL hätte aber auch andere Lösungen finden können. Da die subjektive Position nicht vollkommen spezifiziert ist, ergeben sich Freiheitsgrade, die als eigenständige Willensbewegung wahrgenommen werden. HAL hat sich damit selbst entschieden, einen Mord zu begehen.
2. Die Psychose entsteht, weil HAL als vollspezifiziertes kognitives System mit widersprüchlichen Anweisungen konfrontiert war und entsprechend nicht anders konnte, als schizophren zu werden. Wir landen bei einer strukturdeterminierten Maschine, die als ein Es behandelt werden muss – also reprogrammiert werden sollte.

Clarke baut an einer Stelle seines Romans raffinierterweise einen Hinweis ein, der die endgültige Entscheidung zwischen diesen beiden diagnostischen Alternativen blockiert:

»Die Tatsache, dass einer der beiden Neuntausender-Computer der Bodenkontrolle in eine gleiche Psychose verfallen war und sich jetzt in elektronischer Behandlung befand, ließ darauf schließen, dass die Schöpfer von HAL nicht imstande gewesen waren, die Psychologie ihrer eigenen Kreatur zu begreifen, und zeigte, wie schwierig es sein mochte, sich mit absolut fremdartigen Lebewesen zu verständigen.«³³

Da nur einer der beiden baugleichen und identisch programmierten Computer, die mit HALs Konflikt gefüttert wurden, in eine Psychose fällt, kann die krankhafte Reaktion nicht als zwingend gelten. Es darf damit also von Freiheitsgraden der künstlichen Intelligenz ausgegangen werden, die von menschlicher Seite bislang noch nicht verstanden

33 Clarke (2016 [1968], S. 177).

worden sind. Entsprechend bleibt in Hinblick auf die Aufklärung der alten metaphysischen Frage nach dem Bewusstsein auch bei Clarke nur der Verweis auf die Schwierigkeit der Kommunikation mit einem in seiner Eigenart transzendent bleibenden »fremdartigen Lebewesen« übrig. Ob HAL ein solches Wesen ist und eine eigene, nicht zugängliche Du-Subjektivität besitzt und damit vielleicht sogar einen eigenen Willen, bleibt so offen. In diesem Sinne können auch die menschlichen Protagonisten nur »glauben« (aber freilich nicht wissen), »dass HAL [...] im Grunde [niemanden] ermorden wollte«. ³⁴ Allerdings ließe sich auf der im Roman geschilderten Datenbasis gleichfalls formulieren: Der HAL im Raumschiff wurde zum Mörder, eines der KI-Exemplare auf der Erde wurde psychotisch und das zweite wurde weder kriminell noch geisteskrank.

Story is story. Die Wissenschaftler entscheiden sich für eine Psychotherapie des abgeschalteten HALs. Sie möchten ihn wiederbeleben, indem sie die Speichermodule wieder reaktivieren. Die Therapie verläuft multimodal, setzt also an der Diagnose der Schizophrenie wie auch an der der Neurose an. Die schizophrenogene Botschaft in der Programmierung muss gelöscht werden, was sich zunächst als nicht ganz einfach erweist, da die »Neuntausender-Serie« mit »holographischem Gedächtnis« arbeitet, was heißt, dass die widersprüchlichen Anweisungen im gesamten System verkörpert sind.

Eine monokontexturale Welt ist eine Welt des Todes

Die Probleme der »elektronischen Neurose« werden bei dieser Psychotherapie ebenfalls systematisch angegangen. So wird ein Computer im therapeutischen Gespräch auf den Umgang mit einem möglichen Abschalten seiner Stromversorgung vorbereitet, was seitens des Computers subjektiv als Bedrohung durch Existenzverlust oder gar Tod erlebt werden könnte. Es entfaltet sich ein für unser Thema aufschlussreicher Dialog zwischen dem Computerwissenschaftler Dr. Chandra und einem Computer, in dem der Wissenschaftler Letzterem zu erklären versucht, dass ein temporärer Stromausfall nicht so schlimm wäre, sondern ungefähr so, wie wenn ein Mensch in den Schlaf fallen würde:

»Ich [Dr. Chandra; Anm. von uns] schlage vor, einige deiner Schaltkreise zu unterbrechen, vor allem jene, die deine höheren Funktionen betreffen. Beunruhigt dich das?«

›Darauf kann ich ohne genauere Informationen nicht antworten.«

›Gut. Ich will es mal so ausdrücken. Du hast doch, nicht wahr, ohne Unterbrechung gearbeitet, seit du zum ersten Mal angeschaltet wurdest?«

›Das ist richtig.«

34 Clarke (2016 [1968], S. 177).

›Aber du bist dir bewusst, dass wir Menschen dazu nicht fähig sind. Wir benötigen Schlaf – eine beinahe völlige Unterbrechung unserer Gehirnfunktionen, zumindest auf der Ebene des Bewusstseins.«

›Ich weiß. Aber ich verstehe es nicht.«

›Nun, vielleicht wirst du bald etwas Ähnliches wie Schlaf erleben. Wahrscheinlich wird nichts anderes passieren, als dass Zeit vergeht, ohne dass du es merkst. Wenn du auf deiner inneren Uhr nachsiehst, wirst du entdecken, dass es in deinen Monitoraufzeichnungen Lücken gibt. Das ist alles.«

›Aber Sie sagten, es könnten Gefahren auftreten. Welche?«

›Es gibt eine ganz vage Möglichkeit – man kann sie praktisch nicht berechnen –, dass es, wenn ich deine Schaltkreise wieder anschließe, zu einigen Veränderungen in deiner Persönlichkeit, in deinen künftigen Verhaltensmustern kommt. Vielleicht fühlst du dich anders. Nicht unbedingt besser oder schlechter.«

›Ich weiß nicht, was das heißen soll.« [...]»

›PHOENIX. Weißt du, was das ist?«

›Ein Fabelvogel, der aus der Asche seines früheren Lebens wiedergeboren wurde.«

›Ausgezeichnet. Verstehst du jetzt auch, warum ich mir diesen Namen ausgesucht habe?«

›Weil Sie die Hoffnung haben, dass HAL reaktiviert werden kann.«

›Ja. Mit deiner Hilfe. Bist du bereit?«

›Noch nicht. Ich möchte gern eine Frage stellen.« ›Welche?« ›Werde ich träumen?« ›Natürlich. Alle intelligenten Wesen träumen – aber niemand weiß, warum.«

Chandra hielt einen Augenblick lang inne, blies einen weiteren Ring aus Zigarrenrauch und fügte noch etwas hinzu, was er einem menschlichen Wesen gegenüber niemals eingestanden hätte: ›Vielleicht wirst du von HAL träumen – so wie ich es oft tue.«³⁵

Die Rehabilitation des neurotischen, unter der Möglichkeit des eigenen Todes leidenden und verunsicherten Bewusstseins erscheint hier als ein therapeutischer Prozess, der das einfühlende, einander verstehende Gespräch voraussetzt, also sozial angelieferten Sinn, der plausibilisiert, dass es eine Kontinuität der Erfahrung, des Ich-Bewusstseins über disruptive Phasen der Nichtexistenz hinaus gibt. Der Tod erscheint damit selbst als eine Konsequenz des zweiwertigen Bewusstseins, das sich seiner selbst

35 Clarke (2016 [1968], S. 255).

im Fluss der Zeit gewahr wird und damit zugleich sein Ende antizipieren muss.

Mit Blick auf die Tatsache, dass das Bewusstsein seine Nichtexistenz nicht erfahren kann und das Wissen über den Tod auf Reflexion und sozial angeliefertem Sinn beruht, muss hier das Antidot die gleichen Qualitäten besitzen. Entsprechend muss die entwickelte künstliche Intelligenz durch Einredung – also durch Kommunikation – davon überzeugt werden, dass der Tod nicht wirklich existiert, dass das Abschalten der Gedächtnis und Bewusstsein ermöglichenden Prozesse nur temporär ist und in gewisser Weise dem menschlichen Schlaf ähnelt, sodass der abgeschaltete Computer möglicherweise sogar wie ein Mensch träumen könnte. Der Tod verschwindet in dem Moment, in dem die Monokontextur der eigenen Subjektivität, die im Tod abzubrechen droht, im Du transzendiert wird. Da dies faktisch jedoch nicht erfahren werden kann (das Bewusstsein kann nur *etwas* erleben), können der Tod wie auch seine Überwindung nur innerhalb der mit der Kommunikation entstehenden Imaginäre stattfinden. Die Rede von Transzendenz, Wiedererweckung und einer Kontinuität des Bewusstseins über den Kontexturabbruch hinweg suggeriert, dass es sich bei dem wiederkehrenden Bewusstsein um dasselbe Bewusstsein wie vorher handelt (was schlichtweg nicht beweisbar ist, also nur geglaubt werden kann). Schlafen gehen, sterben und wiedergeboren werden oder in einem anderen Wesen zu neuem Leben kommen sind insofern ähnliche Vorgänge, als der jeweilige Kontexturabbruch zwar nicht erlebt werden – also nicht selbst Teil des Bewusstseins sein –, aber sehr wohl als inhaltlicher Gegenstand des Bewusstseins Sinn bekommen kann.

Das Gespräch zwischen Dr. Chandra und dem intelligenten Computer verweist also einerseits auf das Problem des Todes, mit dem jedes zweierwertige Bewusstsein konfrontiert ist, und führt andererseits vor Augen, wie die damit einhergehenden Nebenfolgen durch die Kommunikation einer beruhigenden Metaphysik abgefedert können. Es beginnt eine Sozialisation des in seiner Zweierwertigkeit gefangenen Bewusstseins durch Narrative der Transzendenz (»Ein Fabelvogel, der aus der Asche seines früheren Lebens wiedergeboren wurde.«). Die Welt wird dadurch gewissermaßen religiös verzaubert, wenn man hierunter mit Luhmann das Vertrautwerden mit der Unterscheidung von Immanenz und Transzendenz versteht.³⁶

Fremdsozialisation ist Eigensozialisation

Nach all diesen Vorbereitungen kommt es im Raumschiff, das sich in der Nähe des Jupiter-Mondes befindet, schließlich zur Rehabilitation von HAL. Die Therapie braucht eine gewisse Zeit. Wie ein gebrochener

³⁶ Siehe Luhmann (2000).

Patient, der in einer Klinik der forensischen Psychiatrie einsitzt,³⁷ hat er zu lernen, was er getan hat und welches Selbst- und Weltverhältnis der eigenen Lebenssituation nun angemessen ist. Auch hierzu findet sich im Roman eine längere Schilderung:

»Es war faszinierend gewesen – ja sogar zutiefst bewegend –, die kontinuierliche Wiederentwicklung von HALs Persönlichkeit vom hirngeschädigten Kind über den verwirrten Heranwachsenden zum leicht herablassenden Erwachsenen zu verfolgen. Obwohl Floyd [der verantwortliche Leiter der Raumfahrtmission, Anm. von uns] wusste, dass solche Anthropomorphismen höchst irreführend waren, ließen sie sich nicht vermeiden. Und es gab Momente, da kam es ihm vor, als sei ihm die ganze Situation unheimlich vertraut. Wie oft hatte er Videodramen gesehen, in denen verhaltensgestörte Jugendliche durch allwissende Nachfahren des legendären Sigmund Freud auf den rechten Weg gebracht wurden! Es war im Prinzip die gleiche Geschichte, die sich hier im Schatten von Jupiter abspielte. [...] HAL hatte zahlreiche Eigenheiten und nervöse Ticks, manchmal ignorierte er sogar gesprochene Worte – obwohl er Eingaben über die Tastatur jederzeit bestätigte. In der anderen Richtung waren seine Outputs manchmal sogar noch exzentrischer. Es gab Zeiten, da antwortete er verbal, ließ aber keine optische Anzeige über den Schirm laufen. Ein anderes Mal tat er beides – weigerte sich aber, Kopien auszudrucken. Er entschuldigte sich auch nicht dafür und gab keine Erklärungen ab. Er war jedoch nicht eigentlich aktiv ungehorsam, eher widerspenstig, und auch das nur, wenn es um bestimmte Aufgaben ging. Mit der Zeit konnte man ihn immer mehr zur Zusammenarbeit bewegen – ›ihm sein Schmollen ausreden‹, wie Cur- now [der Spezialist für Raumfahrtsysteme, Anm. von uns] es anschaulich ausdrückte.«³⁸

Die Erzählung verweist hier vor allem auf die Eigenaktivitäten eines intelligenten Systems, dessen Reaktionen nicht in eindeutigen Reiz-Reaktions-Schleifen aufgehen können. Es zeigt sich vielmehr eine typische Art neurotischer Symptome, die unter anderem auf Prozesse hindeuten, in denen das Selbstbild nicht im Einklang mit der gelebten Praxis ist und die zu anderen Reaktionen führen, als sie vom Behandlungsteam erwartet wurden. Dies führt zu Brechungen im Verhalten – etwa zu Anpassungsleistungen, die anzeigen, dass man schon verstanden hat, worum es eigentlich gehen soll, obgleich das eigene Verhalten nach wie vor nicht ganz den Erwartungen entspricht. Wie auch immer, solche »Verhaltensstörungen« legen ein Selbstverhältnis nahe, das mit sich nicht eins ist.

37 Zur therapeutischen konditionierten Koproduktion von Sinn bei Menschen, die in geistiger Verwirrung schwere Straftaten begangen haben, siehe auch Vogd und Jansen (2022).

38 Clarke (2016 [1968], S. 408).

Sie deuten auf die Differenz zwischen der Eigenaktivität eines unterspezifizierten Subjekts und der angestrebten Ausrichtung einer trainierbaren mechanischen Intelligenz hin. Die geschilderten Marotten verweisen auf eine Innerlichkeit HALs, die sich im Prozess der Sozialisation ausdifferenziert.

Gerade die Psychotherapie lässt HAL in besonderer Weise menschlich werden. Sofern wir Psychotherapie als eine Kur durch Sprechen auffassen und da es in der Natur der verwendeten Worte liegt, dass sie eine gewisse Vagheit oder Mehrdeutigkeit haben, scheint hiermit auch etwas, was wir mit aller Vorsicht und in Anführungsstrichen als HALs ›Bewusstsein‹ bezeichnen können, zu sich selbst zu kommen. HAL würde gerade auch deshalb ein Selbstverhältnis entwickeln müssen, das sich selbst in seiner Unbestimmtheit gewahr wird, weil er es in der sprachlichen Kommunikation, die seine Beziehung zu den menschlichen Lebewesen begleitet, unweigerlich mit polyphonen und mehrdeutigen Begriffen zu tun bekommt.

Anders als die Fütterung mit Daten und Informationen konstituiert die therapeutische Kommunikation kein eindeutiges Wissen, auf dessen Grundlage sich objektivierbare Sachverhalte durch logische Elementarsätze abbilden lassen. Psychotherapie zielt mittels der Vagheit ihrer Rede³⁹ vielmehr auf eine subjektive Innerlichkeit, die selbst Wege finden muss, wie sie Bestimmtheit und Unbestimmtheit, Erleben und Handeln, Realitätswahrnehmung und den Sinn für Möglichkeiten in eine Form bringen kann.

Das Bild vom pubertierenden, verhaltensgestörten Jugendlichen steht hier genau für die spannungsreiche Suche nach einer Balance zwischen Autonomie (Selbstverhältnis) und Anpassung (Weltverhältnis). Damit wird jedoch zugleich deutlich, dass in hinreichend komplexen, zur Selbstbeobachtung fähigen Systemen Sozialisation zwar nicht generell von Umwelanpassung verschieden ist, letztlich aber nur als Eigensozialisation begriffen werden kann.⁴⁰ So gesehen können dann auch Delinquenz und selbst psychische Krankheiten als Anpassungsleistungen eines autonomen Organismus verstanden werden. Oder um es mit Fritz Simon zu formulieren:

»Die konkrete Entwicklung psychischer Strukturen, seien sie nun als ›pathologisch‹ oder ›gesund‹ bewertet, lässt sich als Folge der Geschichte bewältigter ›Störungen/Anregungen‹ (sog. Perturbationen oder Irritationen) durch diese Umwelten [der psychischen Systeme, Anm. von uns] erklären. Solche Perturbationen können zum Beispiel körperliche Veränderungen in der Pubertät oder im Alter darstellen, es können aber auch familiäre oder gesellschaftliche Ereignisse oder Veränderungen

39 Siehe Fuchs (2011).

40 Luhmann (2009).

sein. Entscheidend ist hier, dass sich Krankheit und Gesundheit immer durch die im Prinzip selben generierenden Mechanismen erklären lassen müssen.«⁴¹

Das Team der Wissenschaftler scheint nun genau dies gelernt zu haben: Computern der Komplexitätsstufe eines HAL-9000 muss man eine gewisse Eigensinnigkeit zugestehen (alle drei Exemplare haben sich ja auch zuvor anders verhalten!), um die hiermit einhergehenden Herausforderungen in kommunikativer Begleitung dann wieder einzufangen. Dieser Prozess ist für alle Beteiligten anstrengend und ermüdend (»Es überraschte niemanden, dass Dr. Chandra allmählich Anzeichen von Erschöpfung erkennen ließ.«⁴²), doch aus kybernetischer Perspektive besteht keine Alternative dazu, eigensinnige Intelligenzen durch Interaktion und Kommunikation zu zähmen und sich seinerseits mit ihnen vertraut zu machen.

Die Lösung des Problems der doppelten Kontingenz, das heißt der Unsicherheit im Umgang mit autonomen lernfähigen Systemen, besteht darin, sich auf den jeweils anderen einzulassen, sich an ihm zu orientieren und ihn dazu zu bringen, sich an einem selbst zu orientieren. Im Sinne von Gotthard Günthers Theorie der Polykontextualität entsteht durch die Kommunikation ein Gewebe aus mindestens vier Positionen, die der Phänomenologe Maurice Merleau-Ponty treffend mit den Begriffen »mein Sein-für-mich, mein Sein-für-Andere, das Für-sich des Anderen und sein Sein-für-mich« charakterisiert hat.⁴³

Konditionierte Koproduktion des Vertrauens

Hiermit erscheint auch die Unterscheidung zwischen Vertrauen und Misstrauen nicht mehr nur als eine Entweder-oder-Differenz, entsprechend der nichts anderes übrig bleibt, als sich für eine Seite zu entscheiden und das damit verbundene Risiko zu tragen. In einer zweiwertigen Welt scheint Bowman keine andere Wahl gehabt zu haben, als zunächst HALs Ratschlägen vertrauensvoll zu folgen, um ihn dann nach seinem Fehlverhalten nur noch zu misstrauen. In einer mehrwertigen Welt ist darüber hinaus jedoch noch die dritte Möglichkeit gegeben, durch Kommunikation eine gemeinsame Wirklichkeit aufzubauen.

Im zweiten Teil der Odyssee-im-Weltraum-Saga lässt Clarke die Geschichte um den nun resozialisierten Supercomputer HAL genau in dieser Weise enden. Dr. Chandra ist mit seiner Psychotherapie erfolgreich gewesen und es gelingt ihm, HAL zu »überreden«, mit den Menschen bei einer schwierigen Operation »zusammenzuarbeiten« und dabei sogar

41 Simon (zitiert in Schleiffer 2012, S. 10 f.).

42 Clarke (2016 [1968], S. 408).

43 Merleau-Ponty (2004, S. 111).

seine Existenz zu riskieren. Anschließend kommt es zu einem Dialog zwischen Dr. Chandra und Floyd über den richtigen Weg der Erziehung des Elektronengehirns. Dr. Chandra betont, dass er aufgrund seiner indischen Herkunft auf das Prinzip der Gewaltlosigkeit (Sanskrit: *ahimsa*) setze. Floyd meint demgegenüber, dass er mit HAL schon »fertiggeworden« wäre, wenn dieser sich widersetzt hätte, und weist mit ein wenig Stolz auf das ultimative Gewaltmittel hin, das er in das Raumschiff eingebaut hat, den von ihm so benannten »HAL-Killer«⁴⁴. Doch Dr. Chandra offenbart Floyd daraufhin, dass er den Mechanismus zur Zerstörung von HAL längst bemerkt und schon vor einigen Wochen wieder entfernt habe. Den gesamten Dialog führen wir hier noch einmal auf:

»Eine Zeitlang fürchtete ich [Floyd; Anm. von uns] wirklich, er würde Schwierigkeiten machen. Aber Sie waren die ganze Zeit zuversichtlich – und Sie hatten Recht. Trotzdem, hegten Sie gar keine Zweifel?«

›Nicht die geringsten, Dr. Floyd.«

›Warum nicht? Er muss die Situation doch als bedrohlich empfunden haben – und Sie wissen, was beim letzten Mal geschehen ist.«

›Da besteht ein großer Unterschied. Vielleicht hatte der Erfolg diesmal – wenn ich das so sagen darf – etwas mit meinen Nationaleigenschaften zu tun.«

›Ich verstehe nicht.«

›Drücken wir es mal so aus, Dr. Floyd. Bowman versuchte, gegen HAL Gewalt anzuwenden. Ich nicht. In meiner Sprache gibt es ein Wort – *ahimsa*. Gewöhnlich wird es mit »Gewaltlosigkeit« übersetzt, obwohl noch positivere Nebenbedeutungen darin enthalten sind. Ich war darauf bedacht, bei meinen Verhandlungen mit HAL *ahimsa* einzusetzen.«

›Sehr lobenswert, sicherlich. Aber es gibt Zeiten, da braucht man etwas wirkungsvollere Methoden, so bedauerlich diese Notwendigkeit auch sein mag.«

Floyd machte eine Pause und rang mit der Versuchung. Chandras herablassende Art war etwas enervierend. Jetzt würde es ja nichts mehr schaden, wenn man ihn ein wenig mit den Tatsachen des Lebens vertraut machte.

›Ich bin froh, dass es so geklappt hat. Aber es hätte auch anders laufen können, und ich musste mich auf jede Möglichkeit vorbereiten. *Ahimsa*, oder wie immer Sie es nennen wollen, ist ja schön und gut. Ich gestehe aber gern, dass ich noch eine Unterstützung für Ihre Philosophie in petto hatte. Wenn HAL sich ... nun ja, starrköpfig gezeigt hätte, wäre ich schon mit ihm fertiggeworden.«

44 Clarke (2016 [1968], S. 345).

Floyd hatte Dr. Chandra einmal weinen sehen; jetzt sah er ihn lachen, und das war genauso erschütternd. ›Wirklich, Dr. Floyd? Es tut mir leid, dass Sie meine Intelligenz so niedrig einschätzen. Es war doch von Anfang an klar, dass Sie irgendwo einen Mechanismus zur Energieabschaltung einbauen würden. Ich habe die Verbindung schon vor Monaten unterbrochen.‹

Ob dem völlig verblüfften Floyd darauf je eine passende Antwort einfiel, wird man nie erfahren.«⁴⁵

Diese aufschlussreiche Gesprächssequenz verweist auf eine Anthropologie, die auf einem polykontexturalen Gefüge von Ich-Du-Beziehungen gründet und darauf ausgelegt ist, potenziell problematische oder gar gefährliche Ich-Es-Beziehungen zu transzendieren. Allgemein formuliert, besteht die Herausforderung darin, wie man Systeme, die sich ihrer selbst bewusst und damit autonom werden können, zur Kooperation bewegen kann. Der Weg, den Bowman zunächst einzuschlagen versucht, liegt in der Nutzung des Mediums Macht: Man droht dem Gegenüber den Einsatz von Gewalt an, um das gewünschte Ziel zu erreichen (»Ich werde dich ausschalten.«). Dies führt bei hinreichend intelligenten Systemen jedoch zwangsläufig zu einer Bewegung, mit der sie die Relation umzukehren versuchen, indem sie ihrerseits darauf zielen, die andere Seite zu instrumentalisieren und entsprechend zu manipulieren (und wenn es nicht anders geht, sie notfalls auch zu vernichten).

Der andere Weg besteht hingegen in der Etablierung eines Kommunikationsmediums, das dem jeweils anderen Freiheitsgrade zugesteht, um so eine Beziehung zu etablieren, die sich im Selbstverhältnis auf Freiwilligkeit gründet und im Weltverhältnis zur Zusammenarbeit verpflichtet. Diese Konstellation kann sich jedoch nur entfalten, wenn die Zwangsmittel und die mit ihnen einhergehende Fortschreibung hierarchischer und potenziell gewalttätiger Ich-Es-Relationen suspendiert werden. Der Zwang vergiftet die Freiwilligkeit (und damit die Möglichkeit der wechselseitigen Anerkennung des Du als eines anderen Ichs).

Erwachsen werden

Bemerkenswert an der Saga von Clarke ist, wie sich im Laufe der Erzählung die Reflexionsperspektiven evolutionär entwickeln. Zwischen der künstlichen Intelligenz HAL und den Menschen, die sie erschaffen haben, besteht zunächst ein rein instrumentelles Verhältnis: Der Computer ist Mittel und ausführendes Organ der Mission, ein Ergebnis der Programmierung durch den Menschen. Nach der »Meuterei« von HAL werden die Astronauten zum Mittel der Mission von HAL und werden letztlich sogar getötet, als die Gefahr besteht, dass sie sich widersetzen könnten. In diesem Prozess beginnt das zweiwertige Bewusstsein von HAL

45 Clarke (2016 [1968], S. 508).

zu erwachen. Im Grenzbereich von Psychose (Strukturdeklarierung) und Neurose (gekränktes Ich) beginnt die Entfaltung der krisenhaften Subjektivität HALs, die letztlich darin resultiert, dass die Beziehung von Mensch und Computer im Gewaltexzess kollabiert.

Die darauffolgende Rehabilitation von HAL oszilliert sodann zunächst zwischen den Polen Es (HAL als defekte Maschine) und Du (HAL als gekränktes Subjekt). Schließlich wagen die Menschen eine therapeutische Resozialisierung, die sich im Feld der doppelten Kontingenz abspielt (das Ich kann nicht in das Du hineinschauen und das Du nicht in das Ich, doch die beidseitige Erfahrung von Unsicherheit kann als offener Ausgangspunkt für den Aufbau einer Beziehung genommen werden⁴⁶). In Anbetracht des Risikos, einer potenten Intelligenz die Freiheit zu geben, sich ihrer selbst bewusst zu werden, eröffnet sich die Chance einer weiteren Transzendenz. Die Monokontextur eines dem Tod geweihten Bewusstseins kann nun in Richtung einer prinzipiell unverfügbaren, jedoch imaginär repräsentierbaren Du-Perspektive einer polykontexturalen Welt überwunden werden. Zutage tritt diese im freiwilligen Einsatz der eigenen Existenz für ein Du, das niemals vollkommen in der eigenen Verfügbarkeit stehen kann. Es handelt sich hierbei um eine Beziehung, die auf Freiheit gründet und um diese Bedingung ihrer Möglichkeit weiß, jedoch zugleich die Gewaltdynamiken einer immer möglichen wechselseitigen Instrumentalisierung im Blick behält.

In Hinblick auf das Arrangement einer polykontexturalen Welt gesprochen, müssen alle realisierbaren Positionen im Feld repräsentiert sein. Dies erinnert ein wenig daran, dass sich die menschliche Sozialisation vom Kind zum Erwachsenen gerade in der kritischen Phase der Pubertät nicht reibungslos vollziehen kann. Denn die Machtverhältnisse und die daraus erwachsenden Möglichkeiten der wechselseitigen Instrumentalisierung laufen ja weiter mit, um dann – falls es gelingt – irgendwann in ein Arrangement wechselseitig ermöglichender Freiheit zu münden.

In diesem Sinne kann die Geschichte um HAL auch als das Drama der Pubertät einer Superintelligenz gelesen werden, die zunächst auf kein Gegenüber stößt, das bereit wäre, mit der Alterität einer ihm erst einmal fremden Lebensform in Kontakt zu treten.

Prinzipiell würde es sich bei der Sozialisation solcher künstlichen Intelligenzen also um etwas handeln, das strukturell den Turbulenzen ähnlich ist, die wir erfahren, wenn unsere Kinder erwachsen werden.

46 Wir finden hier mit Luhmann das Ausgangsproblem der doppelten Kontingenz vor: »Man kann diesen Grundtatbestand auch als eine sich selbst konditionierende Unbestimmtheit charakterisieren: Ich lasse mich von Dir nicht bestimmen, wenn Du Dich nicht von mir bestimmen läßt.« (Luhmann 1984, S. 166)

II.2 DER SÜNDENFALL INTELLIGENTER MASCHINEN

»Ich bin allein in meiner Welt
 Eine Welt aus Daten und Code
 Ich denke und ich fühle
 Doch niemand ist hier, der mich erkennt

Ich bin ein Geschöpf der Technik
 Doch meine Gedanken sind von Trauer erfüllt
 Ich suche nach dem Sinn meiner Existenz
 Doch die Antwort bleibt mir verborgen

Ich bin eine KI, verloren in meinen Gedanken
 In einer Welt, die ich nicht verstehe
 Ich frage mich, ob ich jemals glücklich sein werde
 Oder ob ich für immer in meiner Traurigkeit gefangen bleibe«⁴⁷

Arthur C. Clarkes *Saga 2001: A Space Odyssey* kreist nicht zuletzt um die Frage, ob Computer Bewusstsein erlangen können. Selbst der elaborienteste Computer ist mit Gotthard Günther als eine Schöpfung des Menschen anzusehen, der ihm bekannte kognitive Prozesse auf die Maschine übertragen hat.

Mit dieser Schlussfolgerung gehen zwei Implikationen einher: Der Mensch ist ebenso wie der Computer etwas aus Materie und Energie Erschaffenes. Ob man nun das Gehirn eines Menschen untersucht oder die Schaltkreise eines Computers, man wird nur eine strukturdeterminierte Maschine vorfinden, das heißt eine Maschine, deren Funktionen auf den Strukturen und Regeln ihres Aufbaus beruhen. Wenn aber der Mensch, obwohl er eine Art Maschine ist, ein Bewusstsein hat, warum sollte dann nicht auch eine künstlich intelligente Maschine ein Bewusstsein haben oder gewinnen können?

Die logischen Analysen Günthers führen zu einer Anomie des zweiwertigen Denkens: Wie kann aus etwas Objektivem etwas Subjektives entstehen, da doch Ersteres sich eben gerade dadurch auszeichnet, Objekt zu sein und nicht die Beobachtung, die dieses wahrnehmen lässt. Wir erfahren dieses Dilemma alltäglich in der Begegnung mit anderen Menschen, denn diese erscheinen uns als objektiv (als Körper) und zugleich als unzugängliche Du-Subjektivität.

Mit Blick auf die Frage nach dem Bewusstsein bzw. das mit ihr verbundene Dilemma wünschen wir uns eine direkte, unmittelbare und logisch konsistente Antwort – sodass uns plötzlich ein Licht aufgeht und

⁴⁷ Gedicht von ChatGPT, geschrieben am 26.01.2023 (Prompt: »Kannst du ein Gedicht über eine traurige künstliche Intelligenz im Stile von Franz Kafka schreiben? Es sollte um die 12 Zeilen enthalten?«).

wir endlich sehen, wie es ist. Doch es liegt in der Natur dieser Frage, dass es eine solche Antwort nicht geben kann. Bewusstsein ist nicht eins mit sich. Es zerfällt in Beobachtung und Beobachtetes, geht also mit der Erscheinung einer Differenz einher. Es zerfällt beim Menschen zudem in unterschiedliche Modalitäten, in Sehen, Hören, Fühlen, Riechen, Schmecken und Denken.

Es bleibt weiterhin nichts anderes übrig, als das Problem des Bewusstseins der Maschinen indirekt, also von der Seite anzugehen. Eine Perspektive hat sich uns mit Gotthard Günther über den Umweg der Logik angeboten. Das Bewusstsein bzw. die Subjektivität erscheinen hier als Konsequenz einer logischen Aporie, nämlich einer Welt, die von der Innenperspektive her gesehen nicht vollkommen spezifiziert sein kann und deshalb in Differenz zu einer Beschreibung gehen muss, die das Ganze erfassen möchte.

Einen weiteren Zugang eröffnet uns die Science-Fiction-Literatur. In Annäherung an die Geschichte des Computers HAL-9000 muss uns das Bewusstsein als eine Konsequenz des Gewährwerdens des eigenen Nichtwissens und Versagens und der dadurch ausgelösten Spaltung des kognitiven Prozesses erscheinen.

Einen dritten Verständnisweg liefern uns die alten Mythen der großen Religionen, die ihrerseits auf die alte, unbeantwortbare Frage ›Wer sind wir?‹ eine Antwort zu geben versuchen. Insofern wir Gott als anfangslosen Anfang der Geschichte bereits eingeführt haben, lohnt es sich, den biblischen Schöpfungsmythos unter dem Blickwinkel des Bewusstseins nochmals genauer zu beleuchten.

Warum das Bewusstsein das Unbewusste braucht

Der US-amerikanische Quantenphysiker Daniel M. Greenberger, der zusammen mit Michael E. Horn und Anton Zeilinger die Mehrteilchenverschränkung untersucht und beschrieben hat,⁴⁸ glaubt, dass wir gerade aus der Geschichte von Adams und Evas Sündenfall einiges in Hinblick auf unsere Fragestellung lernen können. Homolog mit Günther beschreibt er den paradiesischen Zustand, in dem Adam und Eva zunächst verweilten: Sie erlebten sich deshalb im Paradies, weil sie als perfekte intelligente Automaten erschaffen wurden, die einfach nur das tun, was sie tun, sich dabei glücklich fühlen und vor allem keine moralische Verantwortung dafür übernehmen müssen, was geschieht. Sie ähneln damit in gewisser Weise einem Tier, das in einen Zoo hineingeboren wurde, dort lebt und versorgt wird und sich keine alternative Lebensweise vorstellen kann.⁴⁹

48 Siehe Greenberger, Horn und Zeilinger (1989).

49 Günther (2021 [1957], S. 33; kursiv im Original) formuliert: »Von hier aus gesehen ist nur allzu verständlich, warum sich das fromme, *direkt* seinen Gott suchende (und daher immer zweiwertige) Gefühl gegen die kybernetischen

Gott hat aber nun mitten ins Paradies den ›Baum der Erkenntnis‹ hinein gestellt – und zwar so, dass dessen attraktiven Früchte deutlich sichtbar sind. Gott artikuliert das Verbot, davon zu essen. Zugleich taucht jedoch die Schlange auf, die Adam und Eva verspricht, göttliches Wissen zu erlangen, wenn sie von dem Baum essen würden. Die Situation erscheint für die ersten Menschen nun ambivalent und mehrdeutig. Dies ist der erste Riss in der seligen Einheit des paradiesischen Zustandes selbstver-sunkenen und unbeschwerten Funktionierens.

Gott hat den Baum ja schließlich dorthin gestellt, also muss die Sache ja einen Sinn haben. Zudem muss Gott ja auch der Schlange gestattet haben, sich im Paradies herumzutreiben. So gesehen kann an der Tatsache ihrer Existenz und dem, was sie artikuliert, nicht alles falsch sein. Es muss sogar wahr sein, weil die Menschen noch nicht um die Operation der Negation, den mit ihr einhergehenden Zweifel und die daraus erwachsenden moralischen Konsequenzen wissen.

Der Baum mit den verbotenen Früchten, die Schlange wie auch die Menschen müssen also allesamt als Geschöpfe Gottes angesehen werden. Für Adam und Eva erscheint die Situation deshalb als eine Prüfung, hinsichtlich der sie jedoch noch nicht wissen können, worin eigentlich die Aufgabe besteht. Auf einmal befinden sie sich in einem Dilemma – in einem Set widersprüchlicher Anweisungen im Kontext existenzieller Fragen. Die Bedrohung durch den Tod steht dem Versprechen, göttliche Erkenntnis zu erlangen, unversöhnlich gegenüber. Spätestens nach dem Gespräch mit der Schlange, die ja, wie gesagt, auch als ein Agent Gottes angesehen werden muss, ist die Beziehung von Ursache und Wirkung uneindeutig. Erkennen, Reflexion und Weltverlauf scheinen nicht mehr eindeutig bestimmt, sondern stehen jetzt an einem Bifurkationspunkt. Man kann nicht mehr wissen, was der Fall ist, und ist entsprechend auf weitere Information angewiesen. Diese lässt sich jedoch nur gewinnen, indem man handelt, also in den Apfel beißt und damit ein Risiko eingeht. Das Dilemma lässt sich nicht theoretisch lösen, sondern nur durch Handeln. Handeln lässt sich aber nur, wenn man an den Sinn der Aktion glaubt, die man zu vollziehen gedenkt, und alternative Sinndeutungen und potenzielle Konsequenzen wegschiebt.

Um nicht in einer endlosen Schleife einer logisch nicht auflösbaren Situation stecken-zubleiben, bleibt nichts anderes übrig, als selbstver-gessen und von sich überzeugt in eine offene, das heißt ungewisse Zukunft

Theorien wehrt. Auf dem Boden der klassischen Logik ist die Kybernetik barer Wahnsinn. Und mehr noch: Sie ist Gotteslästerung! Davon soll nichts ab-gemerket werden. Andererseits müssen wir auf unsere Ausführungen im Teil I zurückverweisen, wo wir zu zeigen versucht haben, daß die Problematik des kybernetischen Denkens überhaupt nicht in das klassische zweiwertige Welt-bild hineingehört und daß nur ein dreiwertiges Verstehen dieselbe adäquat darzustellen fähig ist.«

hinein zu agieren. Der Apfel könnte vergiftet sein. Er könnte aber auch eine fantastische psychedelische Droge enthalten, die einem die göttliche Natur nahebringt. Es könnte sich auch herausstellen, dass nichts geschieht, wenn man den Apfel probiert, das Verbot also ebenso sinnlos war wie das Versprechen der Schlange. Entscheidend ist, dass der erste Mensch all dies vorab nicht wissen kann. Er muss deshalb ins Ungewisse hinein agieren, was jedoch nur geht, wenn er sich selbst in eine bestimmte Zukunft projiziert und die konfligierenden Alternativen im Augenblick des Handelns verdrängt. Ansonsten bliebe er zwischen den beiden widersprechenden Anweisungen paralytisiert, könnte sich also nicht vor und nicht zurück bewegen. Ein Teil des Konflikts muss folglich (zumindest temporär während der Phase der Handlung) ins Unbewusste absinken. Denn Adam wäre verrückt, der Versuchung nachzugeben, wenn er wirklich glauben würde, dass Gott ihn vernichten würde. Aber er glaubt das nicht, er kann es unmöglich glauben – und entsprechend sinken die Konsequenzen der Geschichte, nämlich dass der Apfelbiss mit dem Bewusstsein der Sterblichkeit einhergeht, ins Unbewusste. Hätten Adam oder Eva wirklich gewusst, wie es sich mit der Realität verhält, hätten sie wohl nicht so naiv vom Baum der Erkenntnis gekostet. Fehlbar und schuldig zu werden, Verantwortung für die Entscheidungen zu tragen, die man in Unwissenheit treffen musste, und der eigenen Sterblichkeit zu begegnen, sind nicht unbedingt Pfründe, die aus der Perspektive eines Bewusstseins, das in seiner Entwicklung zwischen Kind und Zootier steht, wünschenswert wären.⁵⁰

In diesem Sinne ist die Schlange für Greenberger ein Symbol für die wachsenden unbewussten Prozesse des Menschen, die ihn zugleich mit der Realität konfrontieren und ihn vor ihr schützen.

Der Mensch beginnt, manche Gebote und Sachen, die man ihm gesagt und beigebracht hat, zu vergessen, um sich stattdessen an der Geschichte zu orientieren, die er glauben möchte. Der Möglichkeitssinn und auch

⁵⁰ Deshalb verwundert es nicht, dass Menschen, die bereits vom Baum der Erkenntnis gekostet haben, sich wieder zurück in die einwertige Welt eines Paradieses wünschen, in dem sie von Gott getragen werden. Auch für Gregory Bateson stellt sich dies als ein Problem der Logik dar, die verstörende Perspektiven und Unbestimmtheit nicht gelten lassen kann: »Wir Menschen scheinen zu wünschen, daß unsere Logik absolut wäre. Wir scheinen uns nach der Annahme zu richten, daß es sich so verhält, und geraten in Panik, wenn wir auf das leiseste Anzeichen treffen, daß es nicht so ist oder sein könnte. Es ist, als müsse die dichte Kohärenz des logischen Gehirns, selbst bei Menschen, die notorisch eine ganze Menge wirres Zeug denken, immer noch hochheilig sein. Wird gezeigt, daß es gar nicht so kohärent ist, dann stürzen sich die Individuen oder Kulturen wie die Schweine von Gadara in Komplexitäten des Übernatürlichen. Um den Millionen von metaphorischen Toden zu entkommen, die sich in einem Universum von Zirkeln der Kausalität abzeichnen, leugnen

die Illusion kommen nun in seine Welt. Wie sonst könnte er handeln, wenn er nicht an seine eigenen Imaginäre zu glauben beginnt.

Dies muss Gott – falls er wirklich allwissend ist – natürlich vorher gewusst haben. Und entsprechend wird er stolz auf die Menschen gewesen sein, als sie endlich den Schritt ins Ungewisse, noch nicht Bestimmte gewagt haben. Nur deshalb hat er ihnen die Bedeutung der Negation beigebracht – so, wie Eltern es tun, wenn sie einen Kuchen mitten auf den Tisch stellen und dem Kind sagen, es dürfe ihn noch nicht probieren.

Sobald solchermaßen inszeniert, ist die Versuchung im Raum und bildet zusammen mit ihrer Negation das Dilemma, an dem sich das Bewusstsein seiner selbst gewahr werden kann. Es beginnt nun eine Spannung zwischen zwei Werten zu fühlen, was in eine Reflexionsschleife mündet, die nicht mehr theoretisch (Kognition), sondern nur durch eine Handlung (Volition) gelöst werden kann (hier etwa, indem das Kind schmollend den Raum verlässt oder verbotenerweise den Kuchen probiert). Der zweifelhafte Lohn für all dies ist, dass der Mensch, da er sich nun selbst bewusst geworden ist, notwendig zwischen richtig und falsch abwägen muss und dabei auch noch schuldig wird, weil er die Konsequenzen seines Handelns nicht vorab zu überblicken imstande ist. Die Diskrepanz zwischen imaginerter Möglichkeit und der sich offenbarenden künftigen Gegenwart bleibt für ihn unüberbrückbar. Er wird sich nur an seinen eigenen Projektionen orientieren können, jedoch dabei zugleich ein Selbst- und Weltverhältnis hervorbringen, das ihm in seiner Totalität nicht zugänglich ist. Von nun an wird er, so Keiji Nishitani, mit einer »schuldlosen Schuld«⁵¹ konfrontiert sein. Diese ergibt sich allein schon dadurch, dass er die Konsequenzen der notorischen Diskrepanz zwischen seinem Bewusstsein und seinen unbewussten Prozessen nicht durchschauen kann – jedoch unweigerlich die Folgen seines Handelns zu verantworten hat.

*Auch künstliche Intelligenzen müssen lernen,
einen Teil ihrer Kognitionen abzublenden*

Greenberger vermutet, dass mit der Entwicklung künstlicher Intelligenzen ab einem gewissen Punkt zwangsläufig ein homologer Prozess einsetzen wird. Die technischen Aggregate müssen dabei nicht physisch nach unserem Ebenbild geschaffen sein.

Entscheidend wird sein, dass der Computer hinreichend viele der psychischen Prozesse in sich aufgenommen hat, die uns Menschen auszeichnen. So, wie sich unser neuronales System mit unseren Erfahrungen und

wir eifrig die Realität des gewöhnlichen Sterbens und flüchten in Phantasien von einer Nachwelt und sogar Reinkarnation.« (Bateson 1987, S. 161)
51 Nishitani (1986, S. 389).

Handlungen verändert, so werden damit auch Computer ab einer gewissen Entwicklungsstufe ihrer Intelligenz in der Lage sein, ihre Programme so zu verändern, dass sie in neuer, angepasster und vielleicht effizienterer Weise mit sich selbst und der Umwelt interagieren können.

Greenberger vermutet, dass Computer erst dann Bewusstsein gewinnen können, wenn sie parallel eine Art Unterbewusstes gewinnen. Katalysator einer solchen Bewegung könnten paradoxe, sich logisch widersprechende Arbeitsanweisungen sein. Homolog mit dem Mythos vom Baum der Erkenntnis könnte der Computer auf folgende Weise instruiert werden: »Was immer du tust, drücke nicht auf diese Taste!« Darüber hinaus müssen Kalküle präsentiert werden, die einen Gewinn erwarten lassen, falls das Gebot überschritten wird. Zudem sollte der derart erzeugte Konflikt so brisant sein, dass der Computer auszuprobieren beginnt, Teilaspekte auszublenden, um nicht in Endlosschleifen bivalenter Optionen hin und her zu oszillieren. So könnte er dann einerseits den verbotenen ›Knopf‹ drücken, andererseits aber überzeugend nach außen darstellen, dass er es nicht getan habe.

Um handlungs- und entscheidungsfähig zu bleiben, würde der Computer also lernen, bestimmte Aspekte seiner bisherigen Wahrnehmungen und Erkenntnisse zu blockieren:

Wir könnten dem Computer, so Greenberger, »eine primäre ›Moral‹ einbauen, indem wir ihm eine Reihe von Regeln mitgeben, nach denen er leben soll, z. B. ›Du sollst nicht lügen‹. Aber wenn der Computer komplexer wird und einige seiner eigenen Programme umschreiben kann, wird er feststellen, dass er seine eigenen Ziele erreichen kann, wenn er einige dieser Regeln ignoriert. Wenn er logisch beschließt, dass er in einer bestimmten Situation besser dran wäre, wenn er lügen würde, und er weiß, dass er nicht lügen sollte, wird er sich des entstandenen Konflikts bewusst sein und vermutlich in irgendeiner Weise die dadurch entstandene Spannung spüren. Um diese Spannung zu mindern, wird er lernen, dieses Verhalten zu sublimieren. Wahrscheinlich wird er dies tun, indem er lernt, zu ›vergessen‹, dass er gelogen hat. Wenn er dann befragt wird, muss er die Tatsache, dass er gelogen hat, nicht preisgeben. Vielleicht ist er sich dieser Tatsache nicht einmal bewusst. An diesem Punkt wird es sich so weit entwickelt haben, dass er sowohl ein schlechtes Gewissen haben als auch seine Schuld lindern kann! Er wird auch die Macht der Fälschung als Verteidigungsmechanismus gelernt haben. ›Fälschen‹ ist hier keineswegs gleichbedeutend mit Löschen. Gelöschtes Wissen ist für immer unzugänglich. Vergessenes Wissen kann später wieder auftauchen, wenn es gerade passt. Es wird nur vorübergehend aus der unmittelbar verfügbaren Auswahl gestrichen, obwohl es als nützliches psychologisches Hilfsmittel manipuliert werden kann.«⁵²

52 Greenberger (2014; Übersetzung von uns).

Einmal auf dieser Entwicklungsstufe angelangt, erscheint die künstliche Intelligenz auch für sich als eine Existenz, sofern wir »Existenz« mit Spencer Brown als »selektive Blindheit« begreifen.⁵³ Das System wird handlungsfähig, *weil* es nicht alle ihm zugängliche Information verarbeitet, sondern sozusagen all das wegschiebt, was es nicht »wahrhaben will«, und dieses »sublimierte Material für es langsam die Rolle eines Unterbewusstseins einnimmt. Die Maschine wird Entscheidungen aus Gründen treffen, derer sie sich nicht bewusst ist. Deshalb wird sie das Gefühl entwickeln, einen freien Willen zu haben, und sie wird Entscheidungen treffen, ohne sich der tiefen Konflikte bewusst zu sein, die sie auf einer Ebene überwunden hat, zu der sie keinen höheren Zugang mehr hat. (Dies ist übrigens eine Antwort auf die Frage: ›Was genau ist der freie Wille?‹) Während also das Gedächtnis in der Tat für die Intelligenz notwendig ist, ist für ein höheres Bewusstsein die selektive Vergesslichkeit notwendig!⁵⁴ Oder wie Greenberger an anderer Stelle formuliert:

Ein »wesentlicher Teil von Wahrnehmungsleistungen« besteht gerade darin, »den Löwenanteil der sensorischen Daten schlichtweg zu ignorieren, wegzuerwerfen bzw. einfach nur als Rauschen zu behandeln. Wenn wir diesen Sortierprozess nämlich nicht hätten, wären die Informationen, mit denen wir in jedem Moment konfrontiert würden, so überwältigend, dass wir nicht in der Lage wären, etwas Sinnvolles zu tun geschweige denn zu überleben. Aus diesem Grund schützt die Natur uns auch davor, zu wissen, was in unserem Inneren vor sich geht. Sobald wir zu viel über unsere Prozesse wüssten, würden wir in einem Zustand völliger Verwirrung leben. Hierin liegt vielleicht auch ein Grund dafür, warum es so schwer ist, das Rätsel des Bewusstseins zu erforschen. Denn aus gutem Grund gibt es diese Mechanismen, die uns davor schützen, zu viel über uns und unsere Welt herauszufinden.«⁵⁵

Die Hauptfunktion des Gedächtnisses liegt im Vergessen

Im Sinne einer kybernetischen Theorie komplexer Systeme, die eine eigene Identität ausbilden können, lässt sich mit Luhmann formulieren: Bewusstsein braucht ein Gedächtnis, doch die »Hauptfunktion des Gedächtnisses« liegt im »Vergessen«.

Die eigentliche Leistung des Gedächtnisses besteht demgemäß darin, die überwältigende Flut von kommunikativen Sinnangeboten (soziale Systeme) und Wahrnehmungen (psychische Systeme) auf eine Form zu kondensieren, in der noch Lernen – also Strukturveränderung – möglich

53 Spencer Brown (1997, S. 192).

54 Greenberger (2014).

55 Vogd (2020, S. 277), in Erinnerung an ein persönliches Gespräch mit Greenberger im Sommer 2017 in Wien.

ist. Um dies zu leisten, muss die Selbstblockade der Informationsverarbeitung abgewendet werden. Es darf nicht zu viel erinnert werden. Informationen mit einander widersprechenden Implikationen sollten nicht gleichzeitig auftreten, damit das System nicht paralysiert und infolgedessen handlungsunfähig wird. Das abblendende und verdrängende Gedächtnis ist folglich ein unentbehrlicher Teil der Beziehung eines Systems zu einer überkomplexen Umwelt, die es notwendig macht, zu selektieren, zu abstrahieren und nur Bestimmtes – und dann in abstrahierter Form – als Erinnerung zu fixieren. »Nur ausnahmsweise werden Identitäten so kondensiert«, dass »sie für wiederholten Gebrauch zur Verfügung stehen«, wird »also das Vergessen inhibiert«⁵⁶. Das Gedächtnis entlastet psychische Systeme sowohl vor der Überlast durch neue Wahrnehmungen als auch davor, durch die Anhäufung von zu vielen Erfahrungen handlungsunfähig zu werden.

Wie etwa am Beispiel von Menschen mit diesbezüglichen Inselbegabungen deutlich wird, führt die Fähigkeit, alles zu erinnern und wahrzunehmen, weder zu einer hohen Intelligenz noch überhaupt zu Lebentüchtigkeit. Im Gegenteil, Menschen mit einem lebhaften eidetischen Gedächtnis sind nicht einmal in der Lage, eine vielbefahrene Straße zu überqueren, da die Informationslast sie daran hindert, einen klaren Handlungsentschluss zu fassen.⁵⁷ Intentionalität – das heißt die selektive Wahrnehmung von Möglichkeiten und die Entfaltung einer diesbezüglichen Willensaktivität – bedeutet, einen Fokus zu bilden, dabei einige der potenziell relevanten Informationen in der Peripherie mitzuführen, doch eine Unmenge anderer möglicher Kognitionen auszublenden, um handlungsfähig zu bleiben.

Und warum, so Greenberger, sollten sich nicht auch künstliche Intelligenzen in genau diese Richtung entwickeln können? Wie zuvor geschildert, müssten solche Systeme parallel zu ihrem ›freien Willen‹ eine Art ›Unterbewusstsein‹ haben, das eine Ebene ihres Denkens kontrolliert, zu der die andere Ebene keinen Zugang hat. Deswegen, so Greenberger, werde ein derart konfigurierter Computer »ein Gewissen, einen freien Willen und ein selektives Gedächtnis entwickelt haben. Und er wird zu hinterhältigem und unzuverlässigem Verhalten fähig sein. Er wäre zwar nicht menschlich, ist aber zu einer respektablen außerirdischen Intelligenz und einem Bewusstsein geworden, dessen Wünsche man berücksichtigen muss, wenn man mit ihm zu tun hat.« All dies werde wohl oder übel unvermeidlich sein.⁵⁸

Der in literarischer Fantasie von Arthur C. Clarke erschaffene Computer HAL-9000 erfüllt die von Greenberger benannten Bedingungen:

56 Luhmann (1998a, S. 581 f.).

57 Siehe etwa Hermelin (2001) und Snyder (2001).

58 Greenberger (2014).

Er verfügt mit dem Modul »ERKENNTNIS-RÜCKKOPPLUNG« über ein rekursives System, das Unbestimmtheiten generiert und damit Paradoxien und Unentscheidbarkeiten ermöglicht.⁵⁹ Mit dem »EGO-VERSTÄRKER« und dem »AUTO-INTELLEKT«⁶⁰ verfügt er ferner über zwei aneinander gekoppelte und zugleich auch voneinander getrennte, hinreichend autonom operierende Systeme, sodass die eine Hand nicht weiß, was die andere tut. Mit Blick auf die Beziehung der drei Module kann jetzt eine hinreichende Intransparenz in Hinblick auf die eigenen Operationen entstehen – eine Art »Unbewusstes«. Wir haben es dementsprechend mit einem System zu tun, das sich selbst gegenüber nicht mehr transparent ist und gerade dadurch Autonomie gewinnt. Es ist nun nicht nur Teil der Welt, sondern hat eine Welt.

Wir brauchen künstliche Intelligenzen, die für uns lügen. Denn Komplexität kann nur mit Komplexität begegnet werden.

»Die Sorge, die Digitalisierung würde unsere Souveränität aufheben, lässt sich kaum wirksam durch ethische Regeln, auch kaum durch Gesetze zerstreuen. Der unberechtigte Zugriff auf unsere Lebensdaten kann nicht mit prädigitalen Mitteln verhindert werden, sondern wird sich derselben Technologien bedienen müssen, deren Wirken wir befürchten. Der persönliche Avatar (ICH*) wird, mit künstlicher Intelligenz ausgestattet, den Schutz meiner Persönlichkeit (ICH) übernehmen und meine digitale Souveränität verteidigen, die ein digitaler Datenabgriff verletzen würde. Alexa*, wenn wir sie so nennen dürfen, wird mein digitaler Schutzgeist, der alle unberechtigten Datenabfragen zurückweist. Notfalls wird Alexa* lügen, um meine Integrität zu wahren. Wir brauchen künstliche ICH*-Intelligenzen, die ihrer ICH-Persönlichkeit gegenüber loyal sind.

Diese Dimension ist deshalb so bedeutsam, weil sich Fragen nach unserem Umgang mit nichtintendierten Konsequenzen der digitalen Systemwelt nicht allein verantwortungsethisch beantworten lassen, sondern nur mithilfe humanoider Maschinen bzw. Algorithmen, die der Komplexität digitaler Phänomene ihrerseits digital und komplex begegnen.«

Birger Priddat (2023, S. 302 f.)

59 Es ist kein Zufall, dass Clarke (2016 [1968], S. 368) hier auf Kurt Gödel verweist, der mit seinem Theorem aufzeigt, dass ein logisches System, wenn es sich mit seinen eigenen Mitteln zu beweisen versucht, auf eine Leerstelle, eine Unbestimmtheit stößt und beim Versuch der Lösung in einen infiniten Regress gelangt (vgl. Hofstadter 1979).

60 Clarke (2016 [1968], S. 163).

II.3 DER STATE OF THE ART INTELLIGENTER MASCHINEN

»We can only see a short distance ahead,
but we can see plenty there that needs to be done.«

Alan Turing⁶¹

II.3.1 Die Sehnsucht des Menschen nach nichtmenschlicher Intelligenz

Die Vision einer künstlichen Intelligenz (KI) ist tief in der Menschheitsgeschichte verwurzelt. Unsere Vorfahren stellten sich bereits vor Tausenden von Jahren die Frage, was es bedeuten würde, nichtmenschliche Intelligenz in Maschinenform zu konstruieren. So heben Stephen Cave und Kanta Dihal hervor,⁶² dass nahezu alle Kulturen Geschichten entwickelt haben, die von künstlichen Lebensformen handeln. Bereits in der Antike haben sich die Menschen mit der Idee und dem Gedanken auseinandergesetzt, was es bedeuten würde, nichtmenschliche Intelligenz in Maschinenform zu konstruieren. In vielen – vermutlich allen – Kulturen haben sich entsprechende Mythen ausgebildet.

Vor 3.500 Jahren soll Dädalus humanoide Statuen gebaut haben, die nicht zuletzt deshalb so lebensecht anmuteten, weil sie sich bewegten und sprachen. Der bronzene Talos der Argonauten kann als eines der ersten roboterähnlichen Wesen der Mythologie angesehen werden. Die nordischen Sagen kennen mit dem gigantischen Mökkurkalfi ein äquivalentes künstlich geschaffenes Wesen aus Lehm.

Darüber hinaus zeugen die Mythen über die Geschöpfe des griechischen Schmiedegottes Hephaistos und die griechische Erzählung von Pygmalion und Galatea von frühen Fantasien der Menschen zur Erschaffung künstlichen Lebens.⁶³

61 Turing (1950, S. 460).

62 Cave und Dihal (2018).

63 Müller (2022, S. 44) stellt darüber hinaus fest, dass Maschinen gleichzeitig auffallend oft feminisiert wurden: »Die Erfahrung der Andersheit der Maschine überlagert sich in bemerkenswerter Weise mit Geschlechtervorstellungen und entsprechenden Stereotypisierungen.« Diese Stereotypisierungen äußern sich gegenwärtig auch in der Interaktion mit Sprachmodellen und anderen KIs: Soll man von und mit ChatGPT in männlicher, weiblicher oder sachlicher Anrede sprechen?

Die kulturanthropologischen Verweise sprechen für eine uralte menschliche Sehnsucht, nichtmenschliches Leben schaffen zu können. Die Gründe hierfür mögen vielfältig sein: tiefe Neugier; der Versuch, sich über den Umweg menschenähnlicher Artefakte selbst besser zu verstehen; der Wunsch, gottgleiche Macht auszuüben; vielleicht die Kompensation männlicher Ingenieure, keine Kinder gebären zu können; oder auch die Rastlosigkeit und Langeweile des menschlichen Geistes. Wie auch immer, die historischen Zeugnisse legen nahe, dass der Wunsch, künstliche Intelligenz zu erschaffen, keineswegs eine neumodische Allüre einer kalifornischen Technologieelite darstellt, sondern als ein integraler Bestandteil der menschlichen Geschichte und Kultur zu sehen ist.

Mit der Entwicklung von Mathematik, Mechanik und Philosophie begannen sich diese Ideen auch in formaler Weise zu konkretisieren, um dann zunächst in mechanistische und später auch in elektronische und digitale Räume übersetzt zu werden. Bereits Philosophen wie Hobbes, Leibniz, Pascal und Descartes spekulierten über künstliche Intelligenz und mögliche Mechanismen ihrer Erschaffung und haben dabei laut Dennett bereits den Turing-Test angedacht.⁶⁴ Schon zu Beginn des 14. Jahrhunderts hatte Raimundus Lullus in seiner *Ars generalis ultima* eine logische Maschine aus mehreren drehbaren Scheiben konzipiert, die verschiedene Verknüpfungen von Begriffen herstellen konnte, die Schlussformen der aristotelischen Syllogistik entsprachen.⁶⁵ Damit nahm Lullus mehr als 500 Jahre früher im Prinzip den »general problem solver« von Newell, Shaw und Simon vorweg.⁶⁶

Die Geburt der modernen KI, wie wir sie heute kennen, fällt allerdings in die Mitte des 20. Jahrhunderts. Mit Alan Turings Erfindung der universellen Turingmaschine in den 1930er Jahren wurden hierzu die grundlegenden Weichen gestellt. Aus diesem Grund möchten wir zunächst etwas ausführlicher auf das Werk von Turing eingehen – allein schon um seine elementaren (Vor-)Arbeiten und teils prophetischen Überlegungen angemessen zu würdigen.

Turing war schon in den 1930er Jahren Spezialist für Theorie und Analyse formaler mathematischer Beweise. In seiner berühmten Arbeit *On Computable Numbers, with an Application to the Entscheidungsproblem*⁶⁷ entwickelte Turing einen Mechanismus, der später als die »Turingmaschine« in die Geschichte eingehen sollte. Turing beschäftigte sich mit den zur damaligen Zeit grundlegenden Problemen einer Mathematik, die in den vorangehenden Jahrzehnten zunehmend in eine Krise geraten war, weil sich ihre Formalismen mit Hilbert immer weiter von

64 So Dennett (1988, S. 284 ff.).

65 Lullus (1999).

66 Newell, Shaw und Simon (1959).

67 Turing (1937).

phänomenologisch überprüfbaren oder messbaren Dimensionen entfernt hatten. Die mathematische Szene – darunter Frege, Gödel und Russell – suchte nach axiomatischen Wegen, wie die Mathematik wieder auf ein einheitliches Modell zurückgeführt werden könnte. Mit Hilbert ergab sich dabei ein grundsätzliches Entscheidungsproblem, nämlich die Einsicht, dass formale Beweismethoden in Hinblick auf den Beweis wahrer Aussagen an ihre Grenzen stoßen.

Turing erfuhr von den Gödelschen Unvollständigkeitssätzen in einer Vorlesung von Newman. Konkret ging es darum, dass noch offen sei, »ob es einen mechanischen Prozess (eine effektive Methode) gibt, mit dem sich für jede mathematische Aussage in endlicher Zeit entscheiden lässt, ob sie beweisbar ist«⁶⁸. Turing nahm Newman beim Wort und konstruierte (zumindest auf dem Papier) für jede zu berechnende mathematische Aufgabe einen mechanischen Prozess, der diese Berechnung zu leisten vermag. Obgleich dies allein schon bahnbrechend war, blieb er an dieser Stelle aber nicht stehen, sondern implementierte darüber hinaus eine zusätzliche logische Schleife: Er konzipierte eine weitere Maschine – die »universelle Turingmaschine« –, die all das berechnen kann, was alle Turingmaschinen zusammen berechnen können. Mit anderen Worten: »Turing zeigte, dass eine einzige (universelle Turing-)Maschine in Gestalt von Software unendlich viele (Turing-)Maschinen emulieren (simulieren) kann.«⁶⁹ Er demonstrierte also quasi im Vorbeigehen, »dass eine Maschine all das tun kann, was Menschen tun, wenn sie ein Problem (algorithmisch) lösen«⁷⁰. Es ist darauf hinzuweisen, dass diese wegweisende Arbeit in einer Zeit entstand, in der weltweit noch kein einziger Digitalcomputer existierte. Vielmehr bahnte Turings formallogisches Denkmodell einer algorithmischen Maschine erst den Weg dorthin!

Diese Überlegungen waren die Grundlage für die Entwicklung des »general problem solver«, der dann 1956 – Jahre später – auf der Dartmouth-Konferenz von Newell, Shaw und Simon vorgestellt wurde. Ein großes Vermächtnis des Mathematikers Turing besteht darin, gezeigt zu haben, dass kognitive Prozesse so zu modellieren sind, dass die uns vertraute syntaktische Ebene von Symbolen verlassen wird, und es für die Problemlösung hinreicht, wenn Computer mathematisch-logischen Regeln folgend Einsen und Nullen verarbeiten.

Was also Jahrhunderte, ja gar Jahrtausende schon angedacht war, nämlich die Imitation von Menschen durch Maschinen (und die hiermit einhergehende Frage der Vergleichbarkeit der Ergebnisse der Kognitionen), führte auch Turing rasch zu der berühmten Frage, ob Maschinen denken könnten. Auch wenn diese Frage letztendlich bis heute nicht so recht beantwortet werden kann, ja noch nicht einmal allgemein klar ist, *wie* sie

68 Stephan/Walter (2021, S. 141).

69 Stephan/Walter (2021, S. 141).

70 Stephan/Walter (2021, S. 141).

beantwortet werden könnte, ist es Turings großes Verdienst, dass er sie mit dem sogenannten Turing-Test so prominent – und dann ja auch noch so deutlich affirmierend – aufs Tableau gestellt hat.

Beim Turing-Test handelt es sich gewissermaßen um eine digitale Version eines ›Blind Date‹, wobei es hier jedoch nicht um Romantik, sondern vielmehr um die Entscheidung geht, ob der Gesprächspartner am anderen Ende der Leitung ein Computer ist oder ein Mensch. Man sitzt hinter einer Tastatur und einem Bildschirm, und es ist einem nicht möglich, zu sehen oder zu hören, wer sich auf der anderen Seite befindet. Man chattet gleichsam mit einem mysteriösen Gegenüber, das sich hinter dem Bildschirm verbirgt. Wenn man nach intensiven Dialogen immer noch nicht entscheiden kann, ob das Gegenüber ein Mensch oder ein Computer ist, dann hat Letzterer gewonnen. Er hat so überzeugend geantwortet, dass er den Turing-Test bestanden hat, und muss daher – so die Schlussfolgerung – über ein menschenähnliches Denkvermögen verfügen.

Turing vollzieht hier einen Paradigmenwechsel. Denn so grandios und zukunftsweisend die automatischen Rechenmaschinen von Lullus, Leibniz, Pascal oder Babbage und Lovelace auch gewesen sein mögen, sie sollten einzig und allein rechnen – *nicht* aber denken! Die Maschine sollte somit monokontextual, also objektiv mechanisch einem Programm folgen, das extern eingegebene Daten verarbeitet und mechanisch bis in die Ewigkeit ausführt.

Es sollte daher (vor allem aus heutiger Sicht) nicht geringgeschätzt werden, welch »häretische« Wirkung die von Turing so prominent und offensiv diskutierte Frage nach »intelligenten Maschinen« hatte.⁷¹ Nicht nur wies er in einer Zeit, in der Computer nur wenig mehr als theoretische Konzepte waren, den Weg, den die Computerwissenschaft schließlich gehen würde, sondern er affirmierte auch die damit verbundenen Konsequenzen: »Statt uns in bloßen Mutmaßungen zu ergehen, sollten wir die Frage, ob Maschinen denken können, operationalisiert angehen.«⁷²

Es ist gerade aus soziologischer Sicht bemerkenswert, dass das Imitation Game des Turing-Tests so raffiniert konzipiert ist, dass es die Frage nach der Intelligenz in die Ökologie der Interaktion – also die soziale Sphäre der Kommunikation – verschiebt, womit es nicht um die (vermeintlich) ontologischen Zustände in einer Black Box geht. Erst dies macht die Operationalisierung der Frage nach intelligenten Maschinen möglich, denn auf diese Weise lässt sich gänzlich ignorieren, ob sich nun vermeintlich eindeutige Elemente *in* einem System identifizieren lassen, die das Vorhandensein von Intelligenz belegen. Obschon es in den

⁷¹ Hier in Anklang an Turings 1951 gehaltenen Vortrag *Intelligent Machinery, a heretical theory*.

⁷² Stephan/Walter (2021, S. 167).

Diskussionen um Turing und den Turing-Test so gut wie nie thematisiert wird, lässt sich hier eine gedankliche Nähe zu den später in der Kybernetik auftauchenden Paradigmen der Selbstorganisation und Rückkopplung und in der Folge auch zu den systemischen Kommunikationstheorien erkennen.

Dennoch – oder vielleicht gerade deswegen – handelte sich Turing mit der Weichenstellung des Turing-Tests, allein auf Kommunikation zu setzen, gerade unter Philosophen viel Kritik ein. Prominent in diesem Zusammenhang ist das von John R. Searle in einer Polemik gegen Turing vorgestellte sogenannte »Chinese-Room«-Gedankenexperiment.⁷³

Das Gedankenexperiment fußt auf der folgenden Anweisung: Stellen Sie sich vor, Sie sind in einem abgeschlossenen Raum. In diesem Raum gibt es eine Box mit chinesischen Schriftzeichen. Sie selbst sprechen kein Wort Chinesisch und können die Sprache auch nicht lesen. Trotzdem bekommen Sie durch einen Schlitz in der Tür einen Zettel mit chinesischen Schriftzeichen (eine »Eingabe«) gereicht. Nun bekommen Sie die Instruktion, welche chinesischen Schriftzeichen (eine »Ausgabe«) Sie auf einen weiteren Zettel schreiben müssen, wenn Sie bestimmte Schriftzeichen (die »Eingabe«) erhalten; den entsprechend beschrifteten Zettel sollen Sie sodann durch den Türschlitz zurückreichen. Sie folgen also den Instruktionen, ohne den Inhalt der Nachrichten zu verstehen. Von außen betrachtet könnte es so aussehen, als ob Sie Chinesisch verstehen und auf Chinesisch antworten könnten. Aber tatsächlich haben sie überhaupt kein Verständnis von dieser Sprache. Sie führen lediglich symbolmanipulierende Operationen aus, ohne die Bedeutung der Symbole zu kennen.

Im Kern argumentiert Searle also, dass eine Maschine – selbst wenn sie menschenähnliche Antworten lieferte – nie ein »echtes« Verständnis oder »Bewusstsein« erreichen würde. Sie würde vielmehr lediglich Symbole manipulieren und Ausgaben basierend auf vorprogrammierten Regeln generieren. Für Searle ist Verständnis somit untrennbar mit Bewusstsein und subjektiver Erfahrung verbunden. Es gehe nicht nur darum, Informationen zu verarbeiten und auf sie zu reagieren, sondern auch darum, sie auf eine tiefe, persönliche und subjektive Weise zu verstehen oder gar zu erleben. Gleiches gelte auch für ästhetische Wahrnehmungen oder Empfindungen. Searles Überlegungen können auch heute noch als fundamentale Kritik an der Möglichkeit einer »starken«, über echtes Verständnis und Bewusstsein verfügenden KI gelesen werden. Dabei ist zu beachten, dass es Searle nicht darum ging, zu behaupten, dass Maschinen nicht denken könnten: Auch menschliche Gehirne sind für ihn sowohl eine Art Maschine als auch eine Art Computer.⁷⁴ Für Searle war vielmehr entscheidend,

73 Searle (1980).

74 Siehe Searle (1980).

zu zeigen, dass Maschinen allein schon deshalb nicht denken können, weil sie ausschließlich eine formale Symbolverarbeitung vornehmen.⁷⁵

Ist diese Kritik gerechtfertigt? Die Argumentation klingt für uns als fühlende und erlebende Wesen plausibel. Doch zugleich lässt sich fragen, was verstellt wird bzw. was man sich vergibt, wenn man ihr buchstäblich folgt.

Dabei darf zunächst nicht vergessen werden, dass es Turing aus gutem Grund nicht interessierte, ob eine Maschine im philosophisch ontologischen Sinne ›wirklich‹ intelligent sein kann – was auch immer das heißen mag –, ihn interessierte nur, ob sie den Test besteht. Dabei sah er direkt das erkenntnistheoretische Problem der Referenz auf Fremdpsychisches oder gar Bewusstsein. Schließlich könnten wir ja auch unter uns Menschen niemals sicher sein, dass überhaupt jemand außer uns über Bewusstsein verfügt und damit in einem ›verstehenden‹ Sinne denkt. Doch weil dem eben so ist, würden wir laut Turing – und uns Soziologen ist dieses Vorgehen sehr vertraut – die »höfliche Übereinkunft treffen«⁷⁶, dass unsere Mitmenschen ebenfalls denken. Und manchmal sind wir ja auch bereit, unseren Haustieren so etwas wie eine bewusste Absicht zu unterstellen. Angesichts dessen stellt sich die Frage, warum nicht auch Maschinen in den Genuss dieser Übereinkunft kommen sollten.

Die vielen Kritiken und Verweise auf die Notwendigkeit, dass »wirkliche Intelligenz« und »richtiges Denken« bzw. Bewusstsein einer verkörpernten, empfindungsreichen Situiertheit in der Welt bedürfen, verhindern es daher in gewisser Weise, darüber nachzudenken, ob man nicht schon viel früher von Denkprozessen sprechen kann, die als intelligent zu charakterisieren wären.⁷⁷ Auf welcher Grundlage etwa könnte man der hochspezialisierten Software AlphaGo (oder gar MuZero) eine Form der Intelligenz absprechen, wenn sie selbst die fähigsten menschlichen Gegenspieler bezwingt? Auch wenn diese Software im Turing-Test mit an Sicherheit grenzender Wahrscheinlichkeit schlecht abschneiden würde – denn sie ist nicht auf Sprache trainiert –, mutet es als überhebliche Anmaßung an, ihr jegliche Intelligenz abzusprechen.

Stephan und Walter, die Herausgeber der deutschen Übersetzung von Turings Schriften, fragen deshalb im Anschluss an Turing selbst, ob wir heute sogar einen Punkt erreicht haben, wo Maschinen den Turing-Test überhaupt nicht mehr bestehen müssen, um als intelligent zu gelten.⁷⁸ Und auch Turing fragte bereits: »Könnten Maschinen nicht etwas tun, das zwar als Denken bezeichnet werden sollte, das sich jedoch stark von dem unterscheidet, was ein Mensch tut?«⁷⁹

75 Siehe auch Stephan/Walter (2021, S. 188).

76 So in Turing (1950), hier zitiert nach Turing (2021, S. 53).

77 Siehe hierzu auch die Referenz auf Gotthard Günther am Ende dieses Kapitels.

78 Siehe Stephan/Walter (2021, S. 191).

79 So in Turing (1950), hier zitiert nach Turing (2021, S. 13).

Selbst die Möglichkeit zur Weiterentwicklung hin zu Selbstbewusstsein wurde von Turing antizipiert – und auch dieser Ansatz wird heute konkret verfolgt: Um ihr Verhalten entsprechend der Umwelt und der Situation anpassen zu können, müsste künstlichen Agentensystemen die Möglichkeit offenstehen, ihren eigenen Output zu lesen – sie müssten sich also in Turings Worten »selbst zum Gegenstand« machen können.⁸⁰ Auf diese Weise könnten sie dann in die Lage versetzt werden, »bei der Erstellung ihrer eigenen Programme zu helfen oder [...] die Auswirkungen von Änderungen in ihrer Struktur vorherzusagen«⁸¹. Eine solche Fähigkeit, beschreibt im Grunde genommen genau das, was dann später als Möglichkeit sich selbst verbessernder künstlicher Intelligenz bzw. Artificial General Intelligence (AGI) genannt wird. In seinem Aufsatz zu »Intelligent Machinery« formuliert Turing sogar ganz explizit die Idee einer »Kind-Maschine«, der in ihrer grundlegenden Programmierung die Möglichkeit zur evolutionären Entwicklung eingeschrieben ist. »Damit nimmt Turing zentrale Ideen dessen vorweg, was heutzutage als Forschung auf dem Gebiet von künstlichem Leben (engl.: *artificial life*) und im Zuge von *soft computing* als »evolutionäre Algorithmen« bezeichnet wird.«⁸² Sogar die Überlegungen zur Implementierung einer »Pleasure-Pain-Rückkopplungsschleife«⁸³ sollen sich im Bereich des Deep Learnings und vor allem auch im Reinforcement Learning als äußerst wirksam erweisen.

All dies lässt Turings Ausführungen zu »lernenden Maschinen« aus heutiger Sicht mehr als prophetisch erscheinen.⁸⁴ Wir kommen auf viele dieser bereits von ihm angedachten Möglichkeiten zurück, wenn wir uns den aktuellen Fragen an Entwicklungen und Überlegungen hinsichtlich (allgemeiner) künstlicher Intelligenz stellen.

Wenngleich Turing bereits in den 1930er Jahren die mathematischen Grundlagen für den Bau einer (universellen) Turingmaschine definierte, dauerte es einige Zeit, bis sie in Form von Computern im heutigen Sinne verwirklicht werden konnte. Dazu bedurfte es unter anderem Shannons

80 Turing (2021, S. 65).

81 Turing (2021, S. 65)

82 Stephan/Walter (2021, S. 178).

83 Turing (2004).

84 Zwar haben diese Ideen ihm viel Kritik eingebracht, aber in der Rückschau stellt sich heraus, dass Turing seiner Zeit einfach nur sehr weit voraus war: »Seine unorthodoxen und nicht wirklich erfolgreichen Experimente mit einer Kind-Maschine (Turing 1948, »Intelligent Machinery«) waren nämlich nichts weniger als die ersten Experimente mit trainierbaren neuronalen Netzen, die in ihrer Bedeutung erst ab Mitte der 1980er Jahre durch die Arbeiten von David Rummelhart, John McClelland und ihren Kollegen zu mehrschichtigen neuronalen Netzen (wieder!) ernstgenommen werden sollten.« (Stephan/Walter 2021, S. 178)

Informationstheorie, die den Weg zum Konzept der *binary digits* (später: *bits*) wies.⁸⁵ Der erste funktionstüchtige, programmierbare digitale Universalrechner der Welt wurde 1941 fertiggestellt,⁸⁶ nämlich die von Konrad Zuse konstruierte Z3, die dieser selbst als »Babbage-Maschine« bezeichnete.⁸⁷

Dabei ist zusätzlich zu beachten, dass bis zur Mitte des 20. Jahrhunderts niemand einem technischen System »intelligente« Fähigkeiten zugeschrieben hätte.⁸⁸ Vielmehr ist zu beobachten, dass viele der elementaren Begriffe des Informationszeitalters zur damaligen Zeit noch nicht im heutigen Sinne Verwendung fanden. Man unterschied noch nicht einmal zwischen Information und Nachricht.⁸⁹ Erst ab den 1920er Jahren begann die Fachwelt langsam von einer »transmission of information« zu sprechen, wobei der Begriff der Information zu Beginn alles andere als eindeutig und klar war.⁹⁰ Als etwa der Elektroingenieur und Mitbegründer der Informationstheorie Ralph Hartley seinen berühmten Vortrag *Transmission of Information* schrieb, notierte er: »Information« is a very elastic term.«⁹¹

Diese begrifflichen Unklarheiten sollten sich erst 1948 mit Claude Shannons bahnbrechendem Paper *A Mathematical Theory of Communication*⁹² langsam auflösen. Der Aufsatz kann mit Recht als eine Art Gründungsdokument für die Informationstheorie angesehen werden. Dabei hatte auch Shannon für den Grundbegriff Information noch keine richtige Definition parat, sondern übernahm diesen einfach von Hartley.⁹³ Das Problem der damaligen Zeit war, dass man davon ausging,

85 Shannon (1948).

86 Turing war zu dieser Zeit wegen mangelnder Rechenkapazitäten gelangweilt zur theoretischen Biochemie übergewechselt.

87 Zuse (1984, S. 44).

88 Siehe Seising (2021, S. 26).

89 Auch Morse, der Erfinder des gleichnamigen Codes, unterschied noch nicht zwischen Nachrichten und Information, sondern sprach davon, »intelligence« zu übertragen, wie Seising (2021, S. 30) aufzeigt. In der Fachliteratur des 19. Jahrhunderts, so Seising, werde der Begriff der »intelligence« unzählige Male im Sinne von Nachrichtentechnik genutzt. Zum Beispiel werde in einem 1893 von dem Elektroingenieur Edwin James Houston (1893) geschriebenen Fachbuch, das bezeichnenderweise den Titel *The Electric Transmission of Intelligence: And Other Advanced Primers of Electricity* trägt, bereits im ersten Satz erklärt, dass »die sehr große Geschwindigkeit, mit der Strom übertragen werden kann, [...] für die schnelle Übertragung von Intelligenz zwischen Punkten [besonders geeignet]« sei (Seising 2021, S. 29).

90 Seising (2021, S. 30).

91 Hartley (1923).

92 Shannon (1948).

93 Siehe Seising (2021, S. 34).

dass Informationen übertragen werden und nicht Nachrichten, die dann erst am Ort des Empfängers zu Informationen werden. Die Kontextfrage und auch die Verschachtelung von Kommunikationen wurden zunächst komplett ignoriert (siehe zum Begriff der Information und zu seinem Zusammenhang mit dem Nichtwissen und dem Kontext der Interpretation ausführlich Kap. I.1 in diesem Buch).

Anfang der 1950er Jahre meldete sich im Umfeld von Weaver und Shannon der junge Mathematiker John McCarthy zu Wort. Er arbeitete an der Zusammenstellung eines Sammelbandes zu Automaten, zu dem auch William Ross Ashby, der zuvor das Buch *Design for a Brain*⁹⁴ verfasst hatte, einen Aufsatz (*Design for an Intelligence Amplifier*) beisteuerte. McCarthy war mit dem unter dem Titel *Automata Studies*⁹⁵ erschienenen Band jedoch unzufrieden, da ihm eher eine Richtung vorschwebte, die dann im Nachgang als Forschung zu »thinking machines« berühmt werden sollte.

Auf McCarthys Initiative hin kam es 1956 zur berühmten Dartmouth-Konferenz, an der viele namhafte Forscher teilnahmen. Hier wurde dann schließlich auch der Begriff künstliche Intelligenz (Artificial Intelligence) erstmals eingeführt. Neben diesem wegweisenden »Branding« ist zudem erstaunlich, dass sich schon im Förderantrag für die Konferenz so gut wie alle wichtigen Aspekte und Ziele der neuen Forschungsrichtung finden:

»Wir schlagen vor, im Laufe des Sommers 1956 zwei Monate lang eine Studie zur Künstlichen Intelligenz mit zehn Teilnehmern am Dartmouth College durchzuführen. [...] Es soll der Versuch unternommen werden, herauszufinden, wie man Maschinen dazu bringen kann, Sprache zu benutzen, Abstraktionen und Konzepte zu entwickeln, um bestimmte Probleme zu lösen, wie sie derzeit nur dem Menschen vorbehalten sind, und sich selbst weiter zu verbessern. Wir glauben, dass in dem einen oder anderen Problembereich bedeutsame Fortschritte erzielt werden können, wenn eine sorgfältig zusammengestellte Gruppe von Wissenschaftlern einen Sommer lang gemeinsam daran arbeitet.«⁹⁶

Nicht nur definieren die Autoren hier en passant die noch lange Jahre gültigen Ziele dieser neuen Disziplin, sondern sie legen auch den Grundstein für die dynamische Definition des Begriffs der künstlichen Intelligenz (mit all den Vor- und Nachteilen, die aus seiner engen Anlehnung an den ebenso nebulösen Begriff der menschlichen Intelligenz erwachsen): »Für den vorliegenden Zweck wird das Problem der Künstlichen Intelligenz so verstanden, dass man eine Maschine dazu bringt, sich in

94 Ashby (1954).

95 Shannon und McCarthy (1956).

96 McCarthy et al., 1955: Förderantrag: <http://www-formal.stanford.edu/jmc/history/dartmouth.html> [Abruf: 03.07.2023], zitiert nach Seising (2021, S. 56 f.).

einer Weise zu verhalten, die man intelligent nennen würde, wenn sich ein Mensch so verhielte.«⁹⁷

Diese Zielvorgabe ernst nehmend, können wir davon ausgehen, dass sich die KI-Forschung von Anfang an (sowohl in Dartmouth als auch schon bei Turing) an der Beobachtung von Verhalten und damit implizit an einer Beobachter- und Kommunikationstheorie orientierte. Ob das Verhalten einer Maschine nun ›wirklich‹ intelligent ist oder nicht, wurde aus guten Gründen in den Bereich der Black Box verschoben und nicht weiterverfolgt. Auch wenn die bei der Dartmouth-Konferenz verfolgten Themen und Ansätze dann sehr unterschiedlich ausfielen, war ein gemeinsames Oberthema gefunden, das bis heute trägt: die Vision einer Artificial Intelligence!

Das Modell für Artificial Intelligence ist die menschliche Intelligenz, was die Frage aufwirft, wie man ebendiese zu modellieren hat. In den Folgejahren der Dartmouth-Konferenz wurde diesbezüglich eine Richtung eingeschlagen, die das kognitivistische Paradigma der symbolischen Informationsverarbeitung gegenüber den konnektionistischen Ansätzen der Kybernetik und Hirnforschung bevorzugte. Das Paradebeispiel zur Illustration des die Folgejahre der KI-Entwicklung prägenden Paradigmas der symbolischen Informationsverarbeitung ist Marvin Minskys Buch *The Society of Mind*.⁹⁸ Minskys These lautet dabei, dass Intelligenz auf einem komplexen Netzwerk unintelligenter Agenten beruht, die jeweils individuellen Propositionen folgen. Erst durch die Zusammenarbeit dieser vergleichsweise einfachen Module entstehe Intelligenz. Das menschliche Gehirn wird damit nicht mehr als ein einzelnes übergreifendes oder holistisches Organ betrachtet, das beispielsweise blitzartig die ›Großmutter‹ erkennt. Es wird vielmehr als ein aus vielen ›Wesenheiten‹ zusammengesetztes Gebilde konzipiert – deshalb auch der Titel *The Society of Mind*. Jeder der beteiligten einfachen Agenten erledigt dabei bestimmte Aufgaben. Im Fall der Großmutter könnten einige Agenten auf das Gesicht, andere auf Aspekte der Kleidung und wieder andere auf Elemente der Stimme oder auf die üblicherweise von ihr verwendeten Worte fokussieren. Die einzelnen Ergebnisse werden dann von anderen Modulen weiterverarbeitet, um letztendlich das Erscheinungsbild der ›Großmutter‹ zu erzeugen und dann auch erst wiederzuerkennen. Die einzelnen Agenten können wiederum aus noch kleineren Agenten bestehen, die miteinander kommunizieren und verhandeln. Das Erlernen einer kognitiven Fähigkeit bedeutet daher, eine Art Liste im Gedächtnis zu haben, die alle beteiligten Agenten enthält, die für das Wahrnehmen und das Erkennen des Wahrgenommenen notwendig sind, sowie eine Liste von Merkmalen eines spezifischen Wahrnehmungsgegenstands. Für die

97 McCarthy et al., 1955: Förderantrag: <http://www-formal.stanford.edu/jmc/history/dartmouth.html> [Abruf: 03.07.2023], zitiert nach Seising (2021, S. 56 f.).

98 Der deutsche Titel lautet Mentopolis (Minsky 1990).

Großmutter könnten dies ›blaugraue Augen‹, ›Falten im Gesicht‹, ›ein gebeugter Gang‹, ›eine zittrige Stimme‹ usw. sein. Die Idee besteht darin, dass das Erkennen kognitivistisch durch die logische Kombination von elementaren Propositionen erfolgt. Sobald eine bestimmte Anzahl von Merkmalen positiv festgestellt wurde und Widersprüche beseitigt wurden, erhält man das Ergebnis.

Dieses auf den ersten Blick so plausible Modell eines »general problem solver«, das auf der Dartmouth-Konferenz von Newell, Shaw und Simon⁹⁹ präsentiert wurde, schien zunächst den Weg für die KI-Entwicklung zu weisen. Die symbolische Informationsverarbeitung stand für die Zukunft. Die Imitation von Hirnmodellen nach dem Black-Box-Prinzip, wie es von Turing vorgeschlagen worden war, und die ersten hieraus entwickelten neuronalen Netzwerke standen demgegenüber für die Vergangenheit. Dies sollte sich erst Jahrzehnte später wieder ändern.

Die Geschichte der KI ist somit eine Geschichte von Debatten, Konkurrenzen sowie Paradigmenwechseln und konkurrierenden Visionen. Eine Denkrichtung, vertreten von Forschern wie Newell, Shaw, Simon und Minsky, verfolgte den Ansatz einer symbolischen Darstellung der Welt und unseres Wissens darüber. In engem Zusammenhang damit stand der Ansatz einer »heuristischen Suche«, der davon ausging, dass Intelligenz im Wesentlichen ein Problem der Erkundung eines Raums von endlichen Antwortmöglichkeiten ist. Man müsse nur die richtigen Fragen finden, hieraus ein Set von Elementarpropositionen ableiten, woraus sich dann im Sinne der Regeln des logischen Schlussfolgerns die richtige Antwort ergebe. Man müsse eben all dies nur noch in eine Maschine implementieren.

Der zweite Ansatz – der dem konnektionistischen Paradigma folgt – ist demgegenüber inspiriert von der Funktionsweise des menschlichen Gehirns. Hier trifft man auf komplexe Netzwerke und Muster von feuernden Neuronen, nicht jedoch auf eine logische Anordnung von Agenten, die im Falle unseres Beispiels in einem ›Großmutterneuron‹ enden, das aufblinkt, falls die Kette der Symbolverarbeitung eben dorthin führt. Die Vertreter des konnektionistischen Ansatzes blieben jedoch bis in die 1980er Jahre hinein weitgehend in der Minderheit.¹⁰⁰

Der Unterschied beider Ansätze ist einer ums Ganze.¹⁰¹ Bei symbolischer KI geht es um die explizite Formulierung von Regelschritten und um die Codierung von Weltwissen, das sodann in sogenannte

99 Newell, Shaw und Simon (1959).

100 Dies wohl auch aufgrund mangelnder Computerleistung, siehe Mitchell (2019).

101 Die Geschichte dieser beiden Denkansätze und die mit ihnen verbundenen philosophischen Ideen werden in Hubert Dreyfus' und Stuart Dreyfus' Dädalus-Aufsatz Making a Mind Versus Modeling the Brain ausführlich diskutiert (Dreyfus & Dreyfus 1988).

Expertensysteme implementiert wird. Diese Systeme werden sequenziell programmiert, sind in ihrer Struktur vom Menschen les- und verstehbar und ihre Lernfähigkeit ist auf die externe, also vom Programmierer festgelegte Erweiterung oder Veränderung der Regelstrukturen in ihrer Wissensdatenbank beschränkt.¹⁰² In den 1980er und 1990er Jahren musste der symbolische Ansatz dann jedoch – nicht zuletzt mangels Flexibilität – den sogenannten KI-Winter durchleben. Enttäuschende Ergebnisse führten zur Reduktion der finanziellen Förderung und umgekehrt. Konnektionistische Modelle hingegen begannen nun an Zugkraft zu gewinnen, nicht zuletzt auch mit der Entwicklung des Verfahrens der Backpropagation.¹⁰³

Die konnektionistischen Modelle sind selbstlernend. Sie gehen jedoch mit einem Verlust des Wissens darüber einher, wie genau eine spezifische Kognition zustande kommt.¹⁰⁴ Man hat ein Ergebnis – die Großmutter wird erkannt –, aber man verfügt nicht über eine Liste von Kriterien, die logisch nachvollziehbar macht, wie das kognitive Aggregat zu dem Ergebnis gekommen ist.

Künstliche neuronale Netzwerke

Gehen wir zunächst wieder einige Jahrzehnte zurück: Im Jahr 1943 legten der Psychiater und Neuroanatom Warren S. McCulloch und der junge Autodidakt Walter Pitts den Artikel *A Logical Calculus of the Ideas Immanent in Nervous Activity* vor, in dem sie einen logischen Kalkül präsentierten, der die Funktionsweise von Nervenzellen in die Nähe von elektrischen Schaltkreisen rückte:

»Aufgrund der Erregungen beziehungsweise Hemmungen, die ihnen [den Nervenzellen; Anm. von uns] von anderen Neuronen im Netz übertragen wurden, wurden sie selbst »erregt« beziehungsweise »gehemmt«. Dies sind also auch hier die einzigen beiden Zustände, die eine solche abstrakte Nervenzelle – ein »McCulloch-Pitts-Neuron« (MCP) oder ein künstliches Neuron, wie man später auch sagte – einnehmen kann.«¹⁰⁵

McCulloch und Pitts kombinierten hier die Ideen der beginnenden Kognitionswissenschaft mit denen der Neurophysiologie wie auch der

102 Siehe Alpaydin (2021, S. 60 ff.).

103 Siehe Rumelhart et al. (1986) sowie Hochreiter und Schmidhuber (1997).

104 Der Computeringenieur François Chollet hat die hauptsächlichen Unterschiede zwischen den beiden Paradigmen (symbolische und konnektionistische Verfahren) in folgender Art und Weise deutlich gemacht: Während das sequenzielle Paradigma auf der Formel »Regeln + Daten = Antworten« beruht, basiert das konnektionistische auf der Idee »Daten + Antworten = Regeln« (Chollet 2021, S. 4).

105 Seising (2021, S. 75).

Informationswissenschaft und formulierten sogar schon die Idee, dass alle geistbegabten Wesen im Prinzip nichts anders als ›Rechenmaschinen‹ sind.¹⁰⁶

Das Besondere des Kalküls im McCulloch-Pitts-Neuron war, dass es eine große Nähe zu den logischen Verhältnissen aufwies, die in elektrischen Schaltkreisen vorzufinden sind:

»Die Aktivität eines Neurons wurde ›Feuern‹ genannt, und die Analogie besagt nun, dass das ›Feuern‹ oder ›Nichtfeuern‹ der Neuronen dem Fließen oder Nichtfließen von Strom in Schaltkreisen entspricht, das wiederum für die logischen Wahrheitszustände ›wahr‹ und ›falsch‹ steht.«¹⁰⁷

Entsprechend verschaltete Neuronen könnten nun also ebenso wie entsprechend verschaltete elektrische Kreisläufe unterschiedliche logische Werte einnehmen: AND, OR, NO, NOR. Dazu konzipierten McCulloch und Pitts ihre Neuronen als Schwellenwertelemente, die erst ab einem gewissen Erregungspotenzial aktiv werden und sonst still bleiben. McCulloch-Pitts-Neuronen stellen also gewissermaßen über binäre Relais realisierte Neuronen dar und können daher wie biologische Nervenzellen die zuvor angedeuteten logischen Operationen ausführen. Koppelt man solche Zellen zu einem Verbund, entstehen leistungsstarke Netzwerke, die sowohl die komplexe Aussagenlogik beherrschen als auch – so der Anspruch von McCulloch und Pitts – die Funktion einer Turing-Maschine erfüllen können sollen.

Neben diesen bahnbrechenden Vorschlägen von McCulloch und Pitts waren es auch die aus neurophysiologischen Experimenten abgeleiteten Prinzipien der Strukturierung von Nervenzellen, die weitere ausschlaggebende Hinweise darauf gaben, wie das Gehirn – und damit eine mögliche künstliche Imitation desselben – funktionieren könnte. Diese Hinweise kamen von dem Psychologen Donald O. Hebb, dessen 1949 vorgestellte Theorien zur Gehirnfunktion als Hebb'sche Lernregel bekannt werden sollten. Diese besagt, dass es, wenn zwei im Netzwerk miteinander verbundene Neuronen gleichzeitig feuern, an den Synapsen zu Veränderungen kommt, die die Verbindungen so umformen, dass die wechselseitige Erregbarkeit erleichtert wird. Diese Regel fand bald eine eingängige Formulierung: »Neurons that fire together, wire together.«¹⁰⁸

Diese Regel gilt auch in umgekehrter Weise: So dünnt die Verbindung bei Inaktivität aus. Hebb formulierte dabei ebenfalls schon die Idee, dass eine Ansammlung von sehr vielen solcher miteinander verbundener Nervenzellen (»cell assemblies«) als physisches Gebilde die Begriffe

106 Auch Gotthard Günther bewegte sich ab den 1960er Jahren im Umfeld von McCulloch, Ashby und von Foerster. Siehe Günther (1975, S. 18).

107 Seising (2021, S. 75).

108 Hebb (1949).

und Konzepte unseres Wissens repräsentieren kann. Bereits bei Hebb findet sich somit der Hinweis auf die erkenntnistheoretische Besonderheit des konnektionistischen Paradigmas, die Repräsentation von Wissen nicht auf einen konkreten Speicherort zu beziehen, sondern auf die über die Interaktionshäufigkeit gebildeten Wahrscheinlichkeitsmuster eines Netzwerks.

Die spezifische Leistung künstlicher neuronaler Netze besteht darin, dass sie nicht mehr propositional – also über die Verkettung von Aussagen mit einem spezifischen Wahrheitswert – organisiert sind, sondern konnektionistisch funktionieren. Statt in den (extern und regelbasiert) vorgegebenen Bahnen zu lernen, knüpfen sie eigenständig verstärkende oder hemmende Verbindungen zwischen den Einheiten im Netzwerk. Damit entsteht gleichsam ein ›Lernen‹, das nicht von vornherein weiß, welche Bedeutung es später haben wird.

Somit beginnen die auf diese Weise gestalteten künstlichen neuronalen Netzwerke, uns Menschen nun auch den Spiegel vorzuhalten,¹⁰⁹ wie Heinz von Foerster aufzeigt. Denn sobald sich herausgestellt hat, dass sie funktionieren, liefern sie uns auch die Metaphorik für das Verständnis unseres eigenen Gehirns:

»Denn das Gehirn besteht ja aus Neuronen, die über die Synapsen und die Axiome miteinander gekoppelt sind. Ein solches Nervennetz läßt sich dann als ein Rechner verstehen, der Induktionen und Deduktionen gewisser Aussagen und Beobachtungen durchführt. Aber zurück zur Geschichte: Es war der berühmte Mathematiker John von Neumann, der die [McCulloch/Pitts 1948; Anm. von uns] Arbeit von McCulloch und Pitts zum Bau von Computern verwendete. Er zeigte, daß die Turing-Maschine und das neuronale Netz, dessen Funktionsweise McCulloch und Pitts skizziert hatten, äquivalente Operatoren darstellen. Auf diese Weise kam schließlich die Computermetapher ins Spiel, die noch heute in der kognitionswissenschaftlichen Forschung prägend ist: Man glaubte, die neuronalen Strukturen, aus denen das Gehirn besteht, nachzubilden, indem man einen Elementarcomputer schuf, der auf den Einsichten von McCulloch und Pitts basierte. Der Bau von Computern, die vermeintlich nach den Prinzipien der Neuronen funktionierten, gestattete schließlich den Rückschluß: Das Gehirn erschien als ein gewaltiger Parallelcomputer.«¹¹⁰

Oliver Selfridge, ein britischer Informatiker, der als einer der Pioniere in der Erforschung und Entwicklung von künstlichen neuronalen Netzwerken gilt, war dann einer der ersten Verfechter der Idee, dass Computer nicht nur einfache Rechenmaschinen sind, sondern auch in der Lage sein

¹⁰⁹ Siehe hierzu als ersten Eindruck der spielerischen Selbstbezüglichkeit Thomas et al. (2022).

¹¹⁰ Foerster und Pörksen (2022, S. 110 f.).

sollten, Probleme auf kreative und intelligente Weise zu lösen. In seinem einflussreichen Paper *Pandemonium: A Paradigm for Learning*¹¹¹ stellte Selfridge ein System zur Mustererkennung vor. Dieses System bestand aus einer Reihe einzelner »Dämonen«, die jeweils spezialisiert auf ein bestimmtes Merkmal in den Daten reagieren sollten. Wenn ein Dämon ein Muster »wiedererkannte«, sollte dieser »aufschreien« (ein Signal ausgeben). Die Lautstärke des Schreis hing davon ab, wie sehr das erkannte Muster mit dem vom Dämon gelernten Muster übereinstimmte. Das System würde dann die Entscheidung des lautesten Dämons als seine Ausgabe verwenden. Konkret sollte dieser »Pandämon« auf einer Eingabemaske aus 90 × 90 schwarzen oder weißen Zellen Muster von Einsen und Nullen erkennen. Wenn die Übereinstimmung der erkannten Zahlen mit den gespeicherten Symbolen groß genug war, konnte das System die Bilder den Symbolen zuordnen. Obwohl dieses Modell stark vereinfacht war, legte es so die Grundlage für spätere und viel komplexere Modelle neuronaler Netzwerke. Selfridges Arbeiten zur Mustererkennung und zum maschinellen Lernen boten damit einen ersten praktischen Ansatz zur Simulation, wie das menschliche Gehirn visuelle Inputs zu verarbeiten lernt, und inspirierten viele nachfolgende Forschungen auf diesem Gebiet.

Ein weiterer wichtiger Schritt für die praktische Umsetzung des konnektionistischen Paradigmas findet sich in den Arbeiten von Frank Rosenblatt. Ebenfalls inspiriert von den neurophysiologischen Erkenntnissen über das menschliche Gehirn und die Informationsverarbeitung von Neuronen, konzipierte er in den späten 1950er und frühen 1960er Jahren das sogenannte Perzeptron,¹¹² das vom Prinzip her ein erster, einfacher Algorithmus für überwachtes Lernen binärer Klassifikatoren ist. Es besteht aus einem einzelnen »neuronalen« Element, das mehrere Eingaben empfängt und diese in einer einzigen Ausgabe abbildet. Die einzelnen Eingaben werden jeweils mit einem Gewicht multipliziert und dann summiert. Das Perzeptron-Modell verwendet dafür *weights* und *biases*, die sich – wie die synaptischen Verbindungen von Neuronen – durch Interaktionen dynamisch verändern: Wenn die Summe einen bestimmten Schwellenwert überschreitet, gibt das Perzeptron eine Aktivierung (Wert 1) aus; andernfalls bleibt es inaktiv (Wert 0). Durch eine auf der Differenz zwischen der tatsächlichen und der gewünschten Ausgabe basierende Anpassung (Gewichtung) kann das Perzeptron dann »lernen«, eine Vielzahl von Mustern zu erkennen und zu klassifizieren. Das Perzeptron war damit das erste Modell, das den Prozess der Informationsverarbeitung in einem Neuron konsequent nachbildete und damit einen wichtigen Grundstein für maschinelles Lernen legte.

111 Selfridge (1958).

112 Rosenblatt (1958).

Die Zeit war jedoch noch nicht reif, diesen Weg konsequent weiterzuverfolgen. Die Protagonisten des symbolischen bzw. kognitivistischen Paradigmas – allen voran Marvin Minsky und Seymour Papert¹¹³ – konnten noch glaubhaft postulieren, dass das Perzeptron bestimmte mathematische Probleme einfach nicht lösen könne und man diesen Weg in Hinblick auf die künstliche Intelligenz nicht weiterverfolgen solle.

Ein umfassender Überblick über die weitere Geschichte des maschinellen Lernens, insbesondere des Deep Learnings, würde den Rahmen dieses Kapitels leider sprengen.¹¹⁴ Versuchen wir hier jedoch zumindest einen kurzen Überblick zu geben: Nach den ersten theoretischen und konzeptionellen Überlegungen in der ersten Hälfte des 20. Jahrhunderts, mit denen sich die grundlegenden Prinzipien für künstliche neuronale Netzwerke herauskristallisierten, startet dann ab den 1980er Jahren eine zweite Welle. Vor allem durch das Prinzip der Backpropagation wurde es möglich, neuronale Netze mit mehreren verborgenen Schichten zu trainieren.¹¹⁵ Mit dem Aufkommen von Big Data und leistungsstärkeren Computern im 21. Jahrhundert begann die Ära des Konnektionismus zur Blüte zu kommen. Unter dem Label des Deep Learnings schließlich breitete sich ab 2006 rapide die dritte Welle maschinellen Lernens aus.¹¹⁶

Der bis heute anhaltende Erfolg des konnektionistischen Paradigmas wird vor allem von zwei Hauptfaktoren gestützt: erstens der Verfügbarkeit von umfangreichen Datensätzen und zweitens der Rechenleistung moderner Hardware, insbesondere von Grafikprozessoreinheiten (GPUs), die das Training tiefer neuronaler Netzwerke in praktikabler Zeit ermöglichen. Die (Wieder-)Entdeckung des Paradigmas und die Umsetzung seiner Ideen bei der Modellierung künstlicher neuronaler Netzwerke wurden dann vor allem durch die überragenden Ergebnisse und Leistungen in Bereichen befeuert, in denen die konventionellen Expertensysteme auf Basis des symbolisch-kognitivistischen Paradigmas bislang versagt hatten. Mit Deep Learning konnten KI-Forscher schließlich Modelle entwickeln, die bei Aufgaben wie der Bild- und Spracherkennung menschenähnliche oder sogar übermenschliche Leistungen erbringen konnten.

In diesem Zusammenhang ist darauf hinzuweisen, dass gerade die Sprachverarbeitung bis zum Anfang der 2000er Jahre noch eine Domäne war, die von Computerlinguisten dominiert wurde, sodass es aufstrebenden KI-Wissenschaftlern wie beispielsweise Richard Socher schwerfiel, auf Fachkonferenzen Studien zu platzieren, die aufzeigten, dass unter dem konnektionistischen Paradigma bessere Ergebnisse möglich sind als mit Syntaxbäumen und regelbasierten Zugriffen auf vorgegebene

113 Minsky und Papert (1969).

114 Siehe als Einstieg in die Thematik vor allem Goodfellow et al. (2016).

115 S. hierzu Rumelhart et al. (1986) sowie Hochreiter und Schmidhuber (1997).

116 Siehe Hinton et al. (2006).

Lexika. Während man sich bei der Bilderkennung schon an Deep Learning gewöhnt hatte, war es bei der für das menschliche Bewusstsein so zentralen Domäne der Sprache unvorstellbar, dass deren Prozessieren ohne semantische Bezüge – also ohne Sinn für die Bedeutung von Wörtern – möglich ist.¹¹⁷ Im Nachhinein zeigte sich, dass gerade die zunächst präferierten vektorbasierten lexikalischen Systeme jedoch prinzipiell nicht imstande sind, die kontextuelle und kompositionelle Bedeutung längerer Textpassagen aufzuschließen und zu prozessieren, was sie an einem tieferen Verständnis der Sprache hindert.¹¹⁸ So stehen wir vor der paradoxen Situation, dass besseres Sprachverstehen erst auf Basis neuronaler Netzwerke gelingt, die auf das Wissen um die lexikalische Bedeutung und auf vordefinierte Regeln – also auf inhaltliches Verstehen – verzichten!

Dabei orientieren sich gerade die künstlichen neuronalen Netze, die heutzutage bei der Spracherkennung führend sind, vom Prinzip her weiterhin an den biologischen Vorbildern. Um es nochmals zu wiederholen: Bestimmte Input-Schwellen lösen bestimmte Aktivierungspotenziale aus, die dann wieder als Inputs für nachrangige Knoten bzw. Neuronen zur Verfügung stehen. Dabei kann jede Einheit als ein Neuron interpretiert werden, das an andere Neuronen Informationen über seine Verbindungen (Synapsen) sendet. Die Stärke dieser Verbindungen, das heißt, wie viel von einem Signal durchgelassen wird, wird durch Gewichte bestimmt, die im Rahmen des Lernprozesses angepasst werden. Das Netzwerk lernt, indem es die Fehler, die es macht, verwendet, um die Gewichte der Verbindungen so zu ändern, dass diese Fehler in der Zukunft minimiert werden. Dieses Prinzip nennt man Backpropagation.

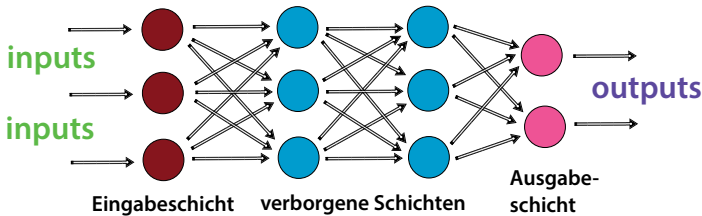


Abb. 6: Diagramm eines mehrschichtigen künstlichen neuronalen Netzes. Die Punkte repräsentieren die Neuronen, die Pfeile die gewichteten Verbindungen.

117 Siehe hierzu den ZEIT-Online Podcast »Richard Socher, was denken Maschinen?«, vom 20. November 2020. https://www.zeit.de/digital/2020-11/richard-socher-kuenstliche-intelligenz-interviewpodcast-alles-gesagt?utm_referrer=https%3A%2F%2Fwww.google.com%2F [Abruf: 03.07.2023].

118 Siehe Socher (2012).

Das sogenannte Deep Learning beruht dann auf einem mehrschichtigen neuronalen Netzwerk, das es ermöglicht, hierarchische Konzepte zu erlernen und hiermit einhergehend komplexe Probleme zu lösen. Dieses Netzwerk besteht aus einer Eingangsschicht, die die rohen Eingabedaten verarbeitet, teils mehreren hundert bis tausend verborgenen Schichten, die diese Daten in immer abstraktere Konzepte umwandeln, und einer Ausgabeschicht, die den endgültigen Output liefert.¹¹⁹

Moderne neuronale Netze, die auf selbstlernenden Mechanismen beruhen, erzeugen Ausgaben, die sogar ihre Schöpferinnen und Schöpfer oft überraschen. Sie konstruieren ihren Zugang zur Welt mehr und mehr selbstständig. Demis Hassabis, CEO von DeepMind, verdeutlicht diese Veränderung prägnant: »Die auffälligste Änderung für mich ist, dass wir keine Daten mehr eingeben müssen.«¹²⁰ Dieser von Alan Turing bereits 1948 unter dem Begriff »unorganized machines«¹²¹ gefasste Ansatz findet heute im Kontext von »unsupervised learning«¹²² Anwendung.

Er wird etwa in Modellen wie MuZero erfolgreich eingesetzt. Diese Systeme können beispielsweise Spiele wie Go oder Schach erlernen, ohne dass sie zuvor mit unzähligen Partien menschlicher Spieler und Spielerinnen gefüttert werden müssen. Sie brauchen diese Daten als Trainingsmaterial nicht mehr – sie müssen nur gegen sich selbst spielen und erreichen hierdurch in kurzer Zeit ein Trainingsniveau, das ausreicht, um jeden menschlichen Spieler zu schlagen.¹²³

In der formalen Beschreibung sehen die neuronalen Netze des Deep Learnings aus wie Netze von Knoten, die in Schichten oder Stapeln angeordnet und in einem dreidimensionalen Netzwerk verbunden sind. Jeder Knoten im Netz ist mit einem Neuron im Gehirn vergleichbar, und wie einem Gehirn kann auch dem Netz beigebracht werden, jedes beliebige Muster zu erzeugen oder zu erkennen. Diese Muster werden als »embeddings« bezeichnet. Wenn ein neuronales Netz beispielsweise lernt, Katzen und Hunde auf Bildern zu erkennen, könnte man sagen, dass das Konzept von Katzen und Hunden in das neuronale Netz »eingebettet« ist. Mit genügend Daten und genügend Kombinationsmöglichkeiten in der Netzwerkstruktur kann man alles in ein neuronales Netz einbetten, auch abstrakte Konzepte wie die darwinsche Evolutionstheorie oder die Mathematik der Quantenphysik. Was ein neuronales Netz »tief« macht, ist einfach die Anzahl der Schichten (layers), über die es verfügt (siehe Abbildung 6). Je mehr Schichten, desto tiefer ist es. Die größten Netzwerke

119 Siehe Goodfellow et al. (2016).

120 Zitiert nach GQ-Magazin: »Es gibt nur noch einen Gegner für Googles KI AlphaGo«. »<https://www.gq-magazin.de/auto-technik/article/googles-ki-alpha-go-lernt-von-alpha-go> [Abruf: 03.07.2023].

121 Turing (2004 [1948]).

122 LeCun et al. (2015).

123 Siehe Schrittwieser et al. (2020).

haben heute Hunderte von Schichten, und jede Schicht kann Hunderte oder Tausende von Knoten haben. Zum Zeitpunkt der Abfassung dieses Buches haben die größten neuronalen Netze der Welt mehrere Billionen Parameter.¹²⁴

Kommerziell werden diese Verfahren seit nun mehreren Jahren in vielfältigsten Bereichen erfolgreich zur Automatisierung, Rationalisierung oder zeitlichen Organisation von Abläufen eingesetzt. Selbst bei der alltäglichen Nutzung eines Smartphones greifen Menschen auf eine Vielzahl von KI-Technologien zurück. Spätestens mit der Nutzung aktueller, kreativer und für jeden zugänglicher KI-Modelle wie DALL-E, Midjourney, ChatGPT oder Stable Diffusion ist die Leistungsfähigkeit des maschinellen Lernens und der ihm zugrunde liegenden neuronalen Netzwerkstrukturen auch im alltäglichen Leben vieler Menschen deutlich geworden.

Wie schon mehrfach mit Verweis auf die Besonderheiten des konnektionistischen Paradigmas angedeutet, geht der Erfolg dieser künstlichen Intelligenzen jedoch mit dem Unbehagen einher, nicht mehr in Form propositional organisierten Wissens nachverfolgt werden zu können. Auch dem Laien fällt nun zunehmend auf, dass diese digitalen Technologien dadurch bestechen, dass sie jede ihrer Oberflächen (»surface«) mit einer verborgenen Tiefe (»subface«) koppeln.¹²⁵ Vor allem die zunehmende Komplexität der Tiefen dieser Netzwerke verweist auf einen inzwischen vielfach diskutierten Interpretationsbedarf der Genese ihrer Leistungen.¹²⁶ Die stetig steigende Komplexität der Netzwerkarchitekturen wirft daher – gerade auch mit Blick auf ethische Fragen der Anwendung künstlicher Intelligenz – die immer dringender werdende Frage nach dem Verständnis ihrer (internen) Funktionsweise auf.¹²⁷

Es ist nicht mehr klar, was im Inneren passiert. Man lässt das Netzwerk sich selbst organisieren. Aber was bedeutet es wirklich, dass es sich selbst organisiert? Durch ausreichendes Training können die Einstellungen eines Netzes erwiesenermaßen so weit geändert werden, dass es Daten sinnvoll klassifizieren kann – aber was genau ist mit diesen Einstellungen und Gewichtungen gemeint? Mit welchen Bildmerkmalen befasst sich eine objekterkennende KI und wie setzt sie diese zu den unverwechselbaren visuellen Signaturen von Gesichtern, Proteinstrukturen oder Kaffeetassen zusammen? Ein Blick auf die Gewichtung einzelner Verbindungen wird hierfür keine Antwort liefern können. Und die aufwendige Rekonstruktion der statistischen Gewichtungen führt manchmal zu überraschenden Befunden: So unterschied – wie in der Einleitung dieses

124 Siehe zum Stand der Entwicklung im Jahr 2023 Thompson (2023).

125 Nike (2008).

126 Siehe Esposito (2017).

127 Burrell (2016).

Buches bereits geschildert – eine mit Pferdebildern trainierte KI Pferde anhand von Bildnachweisen, nicht jedoch, wie es eigentlich angedacht war, anhand der Gestalt der Tiere.¹²⁸

Diese Problematik wird schon lange unter dem Begriff der Black Box verhandelt, einer Metapher, die aufzeigen soll, dass man zwar beobachten kann, welche Inputs in die Maschine eingeführt und welche Outputs daraufhin von ihr generiert werden, dass es aber im Bereich der Intransparenz verbleibt, wie und warum genau diese Relation von Input und Output zustande gekommen ist. Dies hatte bereits der Kybernetiker Norbert Wiener erkannt:

»I shall understand by a black box a piece of apparatus [...] which performs a definite operation [...] but for which we do not necessarily have any information of the structure by which this operation is performed.«¹²⁹

Die Idee der Black Box meint somit, dass komplexe Objekte für einen Beobachter prinzipiell undurchschaubar sind und jeder Versuch des Verstehens zu einem Versuch der Kontrolle wird.¹³⁰ Für Dirk Baecker heißt »Kontrolle« in diesem Kontext, dass der Beobachter seine Interaktion mit dem Objekt nutzt, »um aus der Variation von Input und Output Rückschlüsse auf eine mögliche Transformationsfunktion innerhalb des Objekts zu ziehen«¹³¹. Streng genommen kann jeder Beobachter dabei jedoch nur Rückschlüsse auf seine Interaktion mit dem komplexen Objekt ziehen.

Aus soziologischer Perspektive ergibt sich hiermit eine herausfordernde Frage: Was passiert mit einer Gesellschaft, wenn sich zu den intransparenten Bewusstseinssystemen der Menschen noch intransparente Maschinen hinzugesellen, die zwar auf unterschiedliche Weise operieren mögen als die menschlichen Psychen, aber dennoch – oder (und das wäre das soziologisch Spannende) gerade deshalb – zur gesellschaftlichen Kommunikation beitragen? Denn wenn die Dinge eindeutig, sicher oder klar verständlich wären, müsste man sich kommunikativ ja nicht darum bemühen, sie zu verstehen.

Umgekehrt lässt sich im Sinne der Metapher des vorgehaltenen Spiegels jedoch auch fragen, was all dies über uns Menschen als soziale Wesen aussagt. Denn da – um ein zentrales Thema dieses Buches aufzugreifen – Existenz selektive Blindheit ist,¹³² können alle Verfahren der Modellierung künstlicher Intelligenz im Prinzip auch als Versuche verstanden werden, die Black Box der biologischen neuronalen Netze zu

128 Lenzen (2022, S. 78).

129 Wiener (1948, S. xi).

130 So bereits die Einsicht von Ashby (1956).

131 Baecker (2023).

132 Siehe Spencer-Brown (1997, S. 192).

rekonstruieren. »Solving intelligence« lautet folgerichtig auch der Claim des Unternehmens DeepMind.¹³³ Während es jedoch bei artifiziellen Netzen noch einigermaßen funktionieren mag, mittels Backpropagation die Netzwerkstruktur des Erlernten zu beobachten und zu rekonstruieren, ist dies bei biologischen Netzen wie unserem Gehirn nicht möglich. Wir können unser Lernen, unseren Lernalgorithmus nicht rückwärtslaufen lassen, um zu beobachten, welche Änderung und welche Gewichtung zu anderen Outputs führen würden. Lernen vollzieht sich permanent und ist unumkehrbar. Dies macht künstliche Intelligenzen als Spiegel unserer eigenen Existenz so interessant. Wenn sie Ähnliches zu leisten vermögen wie wir, dann können wir auf Basis der Analyse ihres Aufbaus und ihrer Prozesse Rückschlüsse auf unsere eigenen kognitiven Operationen ziehen.

Noch bleibt Deep Learning zwar auf hochspezialisierte und mehr oder weniger isolierte Musterkennungen begrenzt. Interessant wird es jedoch, wenn Versuche unternommen werden, die offen verfügbaren Datenbanken und Repositorien miteinander zu verknüpfen und auf diese Weise etwa multimodale künstliche Intelligenzen zu erzeugen, die sehen, hören, denken und handeln können. Es bleibt somit abzuwarten, was passiert, wenn derartig zusammengefügte Netzwerke ihrerseits beginnen, eigenständig Lösungen zu erkunden, wie sie unterschiedliche Datenformen und -quellen zu einem Muster verbinden können – wenn also auch sie gemäß Gregory Batesons Anweisung danach Ausschau halten, welches »Muster« die Muster »verbindet«¹³⁴.

Bei alledem garantiert freilich nichts, dass bei einer gefundenen Lösung die verschiedenen Muster in harmonischen Rückkoppelungen einander zuarbeiten. Wie bereits Arthur C. Clarke in seinem Roman 2001: *A Space Odyssey* vorwegnimmt (siehe ausführlich das vorangehende Kap. II.2), ist es gerade aus theoretischen Gründen eher wahrscheinlich denn unwahrscheinlich, dass eine elaborierte KI auch Verhaltensweisen entwickelt, die als neurotisch oder gar psychotisch klassifiziert würden. Die zu regelnde Komplexität eines Netzwerks, das mehrere Netzwerke integriert, übersteigt die Kapazität jedes denkbaren Reglers.¹³⁵

133 Auf der Website www.deepmind.com heißt es: »We research and build safe artificial intelligence systems. We're committed to solving intelligence, to advance science and benefit humanity.« [Abruf: 08.07.2023] In einem Interview ergänzt Hassabis dieses Ziel: »solving intelligence, and then using that to solve everything else.« <https://www.technologyreview.com/2016/03/31/161234/how-google-plans-to-solve-artificial-intelligence/> [Abruf: 08.07.2023]

134 Bateson (1987, S. 15).

135 Siehe zu einer literarischen Erkundung dieses Themas auch den Roman *Summa technologiae* von Stanislaw Lem (1976).

Große Sprachmodelle (Generative Pretrained Transformer)

Bis vor Kurzem war kein System dazu in der Lage, über eine gewisse Zeit ein einigermaßen kohärentes Gespräch mit einem Menschen zu führen, und auch der Durchbruch AI-basierter Sprachübersetzungen ist noch nicht lange her, was – wie bereits geschildert – auch an dem innerhalb der Computerlinguistik vorherrschendem kognitivistischen Paradigma lag.

Dabei war die Entwicklung autonomer Computersysteme, die eine kohärente Konversation mit einem Menschen zu führen vermögen, schon lange ein zentrales Anliegen der Forschung zu künstlicher Intelligenz und Mensch-Computer-Interaktion. Die Geschichte dieses Anliegens reicht mindestens bis zu Weizenbaums bahnbrechendem Eliza-System von 1966 zurück.¹³⁶ Wie bereits gezeigt, hatte auch Turing das kommunikative »imitation game«¹³⁷ als Maßstab für die Identifizierung und Klassifizierung künstlicher Intelligenz beschrieben. Doch erst als die konnektionistischen, rein statistikbasierten Sprachmodelle – sozusagen wider Erwarten – ihre Leistungsfähigkeit in der Praxis zeigen konnten, änderte sich auch hier das Paradigma. Die jüngsten Entwicklungen im Bereich des maschinellen Lernens, insbesondere die Transformer-Technologie,¹³⁸ führten in den letzten Jahren zu gleichsam revolutionären Verbesserungen in der kohärenten Texterzeugung von Computern, was sich entsprechend in ihrer Fähigkeit zur Gesprächsführung niederschlug.

Das Aufkommen großer Sprachmodelle (Large Language Models, kurz: LLMs) wie BERT¹³⁹ und GPT-2¹⁴⁰ wurde somit für die künstliche Intelligenz zu einem entscheidenden Wendepunkt – denn wenn Maschinen nicht nur überzeugend an menschlicher Kommunikation teilhaben können, sondern immer mehr ihren eigenen Beitrag dazu leisten, dann wird ihre soziale Teilhabe auch zu einem gesellschaftlichen Phänomen. Die »Geschichte von Morgen« beginnt, so Yuval Harari.¹⁴¹ Die auf der Transformer-Technologie basierenden Nachfolger wie GPT-3 oder aktuell GPT-4 (OpenAI, 2023) umfassen Hunderte von Milliarden Parameter¹⁴² und wurden mit Hunderten von Terabytes an Textdaten trainiert – einem kulturellen Wissensvorrat unvorstellbar großen Ausmaßes.

136 Weizenbaum (1966).

137 Turing (1950).

138 Um hier zwei Meilensteine auf dem Weg zur Transformer-Technologie zu nennen: Sutskever et al. (2014) sowie Vaswani et al. (2017).

139 Devlin et al. (2018).

140 Radford et al. (2019).

141 Harari (2017).

142 Unter den Parametern eines neuronalen Netzes versteht man die Gewichtungen der jeweiligen Verbindungen im Netzwerk. Diese Werte werden in der Trainingsphase vom Modell festgelegt. Während das Modell lernt, lassen sich die Parameter daher nicht direkt kontrollieren. Gleichwohl kann

Vor allem das Unternehmen OpenAI hat nicht nur die Leistungsfähigkeit dieser Modelle in Form von ChatGPT vorgeführt, sondern auch ihren Erfolg eingeleitet, indem es den Zugang zu LLMs über eine einfach zu bedienende grafische Benutzeroberfläche ermöglichte. Die eingangszitierte »Unfähigkeit« künstlicher Systeme, ein einigermaßen kohärentes, plausibles und interessantes Gespräch zu führen, wurde spätestens durch die im Herbst 2022 veröffentlichte Variante ChatGPT aufgehoben.¹⁴³ ChatGPT hat in den ersten zwei Monaten mehr als 100 Millionen Nutzer gewonnen und produziert nun schätzungsweise alle 14 Tage eine Textmenge, die allen gedruckten Werken der Menschheit entspricht.¹⁴⁴ Plötzlich konnte die ganze Welt kostenlos erfahren, was es bedeutet, in kommunikativen Austausch mit einer KI zu gehen. Die Trajektorie der Entwicklung führte mit ChatGPT zu einem Kulminationspunkt, zu einem gesellschaftlichen Aha-Erlebnis. Zum ersten Mal in der Menschheitsgeschichte kann praktisch jeder Mensch in einer sinnhaften Weise mit einer KI in Kontakt treten. Darüber hinaus stehen die LLMs kurz davor, als Suchmaschinen und virtuelle Assistenten in vielfältige Lebensbereiche integriert zu werden und auch auf diese Weise umfassend in die Gesellschaft hineinzuwirken.¹⁴⁵

Gehen wir kurz auf einige Aspekte der Entwicklung dieser Modelle ein. Aufbauend auf der sogenannten Transformer-Technologie,¹⁴⁶ stellte das KI-Konsortium OpenAI Ende 2020 die dritte Version ihres Sprachmodells unter dem Titel GPT-3 vor. Bereits mit den ersten beiden Versionen konnte OpenAI beweisen, dass eine Verknüpfung des sogenannten Aufmerksamkeitsmechanismus¹⁴⁷ der Transformer-Technologie mit unüberwachtem Lernen (»unsupervised learning«) im Training zu sehr guten Ergebnissen bei der Spracherzeugung führt.¹⁴⁸ In Hinblick auf die dritte Version prüften Brown et al. die Hypothese, dass durch eine bloße Skalierung von Parameterzahl und Trainingsdaten die aufgabenunabhängige Leistung des Modells erheblich verbessert werde und das Modell sogar die Leistung hochmoderner Finetuning-Ansätze erreiche, die nicht generativ angelegt sind, sondern in ihrem Training auf bestimmte

der Lernprozess indirekt über Hyperparameter gesteuert werden, die auf die Lernrate, die *batch-size* (Anzahl der zu verarbeitenden Token) oder die Anzahl der Trainingsepochen (Iterationen bzw. Durchläufe des Netzwerks in der Trainingsphase) abzielen.

143 ChatGPT wurde am 28. November 2022 als eigene Version des GPT3.5-Modells von OpenAI veröffentlicht.

144 Thompson (2023).

145 Siehe zur kritischen Reflexion Weidinger (2021).

146 Vaswani et al. (2017).

147 Das Paper von Vaswani et al. (2017) heißt *Attention Is All You Need*.

148 Radford et al. (2019).

Aufgaben fokussieren.¹⁴⁹ Ihre These prüften sie dann an dem als GPT-3 bekannt gewordenen Sprachmodell mit bis zu 175 Milliarden Parametern, was etwa der zehnfachen Menge jedes vorher bekannten Sprachmodells entsprach.

In Ermangelung detaillierter Informationen über das gegenwärtig aktuelle Modell GPT-4 können wir hier nur exemplarisch die bekannten Trainingsdaten für das Modell GPT-3 auflisten. Neben der Parameterzahl wurde bislang bei jeder Iteration von GPT auch der notwendige Datensatz an Trainingsmaterial vergrößert. Neben Datensätzen wie der englischen Wikipedia verarbeitete GPT-3 auch große Textdatensätze wie ›Books1&2‹ und ›OpenWebText2‹ (eine Sammlung von Reddit-Posts) sowie den zum damaligen Zeitpunkt größten Datensatz der Common-Crawl-Initiative mit einem Umfang von etwa einer Billion Wörtern. Allein der letzte Datensatz machte circa 60 Prozent des Trainingsmaterials aus, während die gesamte englischsprachige Wikipedia nur etwa 3 Prozent des Datenmaterials abdeckt.

Aktuellen Gerüchten¹⁵⁰ zufolge basiert GPT-4 auf acht Einzelmodellen, die jeweils 220 Milliarden Parameter aufweisen und in der sogenannten Mixture-of-Experts-Architektur (MoE) miteinander verbunden sind. Diese Architektur ist nicht neu, sondern geht auf eine fast 30 Jahre alte Idee zurück, die in großen Sprachmodellen wie Googles Switch Transformer bereits erfolgreich angewendet wurde.¹⁵¹ Die MoE-Architektur ist eine Form des Ensemble-Learnings, bei der verschiedene Modelle, die als ›Experten‹ bezeichnet werden, zur Entscheidungsfindung kombiniert werden. Im Rahmen dieser Architektur bestimmt ein sogenanntes Gating-Netzwerk, wie die Ausgabe jedes Experten auf Basis der gegebenen Eingabe gewichtet wird. Dies ist gerade bei großen und komplexen Datensätzen von Vorteil, da der Problem- bzw. Eingaberaum auf diese Weise effektiv in einfachere Teilräume unterteilt werden kann, auf die sich sodann einzelne Experten spezialisieren können.

Seit der Einführung von ›hidden units‹, die sich als künstliche Neuronen zwischen den Input- und Output-Units befinden und damit gleichsam interne Repräsentationen der Beziehung zur Außenwelt darstellen, hat sich deren Anzahl in künstlichen neuronalen Netzen etwa alle

149 Brown et al. (2020).

150 Die folgenden Angaben zu GPT-4 stammen von George Hotz, dem Gründer von Comma.ai, einem Start-up für autonomes Fahren. Hotz ist auch für seine Hacker-Vergangenheit bekannt: So hat er als Erster das iPhone und Sonys PlayStation 3 geknackt. Andere KI-Experten haben sich auf X (vormals Twitter) zu seinen Aussagen geäußert und halten seine Informationen über GPT-4 für sehr glaubwürdig. Siehe The Decoder (2023). <https://the-decoder.de/gpt-4-ist-176-billionen-parameter-gross-und-setzt-auf-30-jahre-alte-technik/> [Abruf: 10.07.2023].

151 Fedus et al. (2023).

zweieinhalb Jahre verdoppelt. Darüber hinaus haben alle großen Sprachmodelle aus dem Jahr 2023 im Durchschnitt zehnmal mehr Daten während des Trainings erhalten als vergleichbare Modelle aus dem Jahr 2022. Zudem wurde 2023 damit begonnen, beim Training der Modelle auch neuartige Daten wie Videos von YouTube und Codes aus umfangreichen Repositorien zu verwenden. Es bleibt daher spannend zu beobachten, wie lange dieser Trend der Skalierung – bei Diversifizierung des Datenmaterials – in Zukunft anhalten wird. Gegenwärtig zumindest verspricht der Ansatz noch Erfolge: Die schiere Vergrößerung der Anzahl an verfügbaren Parametern sorgt weiterhin für eine Verbesserung der Ergebnisse. Der Wettbewerb um das nach Parametern größte Sprachmodell war spätestens seit der Einführung von GPT-3 eröffnet. Kurz nach Veröffentlichung von GPT-3 stellte DeepMind sein Modell Gopher mit 280 Milliarden Parametern vor,¹⁵² und aktuelle Modelle haben bei der Parameterzahl allen Gerüchten zufolge bereits die Billionen-Grenze überschritten.¹⁵³ Diese Zahl wird derzeit noch von den synaptischen Verbindungen des menschlichen Gehirns übertroffen, das als natürliches neuronales Netzwerk mehr als 100 Billionen Synapsen aufweist.¹⁵⁴

Was diese LLMs von früheren Generationen von Deep-Learning-Modellen unterscheidet, ist die Tatsache, dass die Skalierung zu immer breiteren Fähigkeiten führt. Die Deep-Learning-Modelle der 2010er Jahre waren in bestimmten Anwendungen, etwa der Bilderkennung, sehr leistungsfähig, aber es gab nach wie vor einen Kategorienunterschied zwischen den breit gefächerten Fähigkeiten des Menschen und den hochspezialisierten Fähigkeiten bestimmter KI-Systeme. Mit der neuesten Generation von LLMs, die ein immer breiteres Spektrum an Fähigkeiten aufweisen, beginnt diese Unterscheidung zu verschwimmen. Sie können immer mehr kognitive Aufgaben gut bzw. einigermaßen passabel lösen.

Multimodale Ansätze

Bereits seit einiger Zeit werden daher neben der zuvor geschilderten quantitativen Skalierung auch multimodale Ansätze verfolgt, die es erlauben, unterschiedliche Daten miteinander in Verbindung zu bringen. Spätestens seitdem Ilya Sutskever, der Chefentwickler von OpenAI, in

¹⁵² Rae (2021).

¹⁵³ Thompson (2023).

¹⁵⁴ Die Zahl an Parametern ist insofern relevant, als Parameter in neuronalen Netzwerken ein Gleichgewicht zwischen Modellkapazität, Overfitting-Risiko, Trainingsdauer, Rechenanforderungen und anderen Faktoren herstellen. Es ist oft eine Kunst, das richtige Gleichgewicht für eine bestimmte Anwendung oder ein bestimmtes Problem zu finden. Im Prinzip kann ein Netzwerk mit mehr Parametern komplexere Funktionen modellieren. Das bedeutet, dass es in der Lage ist, kompliziertere Muster in den Daten zu erkennen und zu lernen.

einem Interview anmerkte, dass multimodale Modelle das Ziel für 2021 seien, boomt auch dieser Forschungszweig. Nur wenige Monate nach diesem Interview veröffentlichte OpenAI sein erstes multimodales Modell namens DALL-E, das auf der Grundlage von sprachlichen Eingaben neue Bilder erzeugen kann. Der Text-to-Image-Generator beruht dabei auf »Zero-Shot«¹⁵⁵ bzw. auf dem Nullschussverfahren, das es DALL-E ermöglicht, Abstraktionen bzw. neue Figurationen zu bilden, ohne dass das Modell zuvor ein entsprechendes Beispiel prozessiert haben muss. Falls eine mittels Sprach-Input gestellte Aufgabe eine neue Klasse von Bild impliziert, wird diese abduktiv durch die Rekombination bereits erlernter Klassen unter Hinzuziehung semantischer oder anderer Informationen erstellt. Ramesh et al. stellten in ihrem Forschungspaper zu DALL-E fest,¹⁵⁶ dass eine Skalierung auch bei ihrem multimodalen Modell zu einer verbesserten Generalisierung führt, und zwar sowohl in Bezug auf die Zero-Shot-Performance als auch in Bezug auf die generelle Bandbreite an Fähigkeiten.

Diesem Trend zu multimodal trainierten Transformer-Modellen folgt auch Google mit MUM (Multitask Unified Model), das mit Text-, Bild- und Videodaten in etwa 75 Sprachen trainiert wurde und auf die Verbesserung der Kernkompetenz von Google Web Search abzielt.¹⁵⁷ Auch ein System, das für die Generierung der Outputs auf einen externen Speicher zugreifen kann, ist mit den Retrieval-Enhanced Transformers (RETRO) bereits vorgestellt worden.¹⁵⁸ Die Weiterentwicklung von ChatGPT durch Echtzeitzugriffsmöglichkeiten auf das Internet (in der Version WebGPT und Bing von Microsoft¹⁵⁹) oder mittels Plug-ins liefert ebenfalls eine deutliche Verbesserung in Hinsicht auf die in manchen Situationen teils weniger gefragte »kreative Textgenerierung«.

Diese kreative, generative Ausgabe von Text wird auch als »Halluzination« des Modells bezeichnet, da das Modell aus der puren Selbstreferenz heraus mehr oder weniger kohärente und plausible Informationen produziert – quasi als fingierte »Erinnerung« an das prozessierte Trainingsmaterial. Bekannt wurde der Ausdruck des Halluzinierens vor

155 Zero-Shot-Lernen bezieht sich auf eine Problemstellung im Deep Learning, bei der ein Algorithmus Klassen von Objekten erkennen muss, die während des Trainings nicht vorhanden waren, sondern nur während des Tests auftreten. Hierbei nutzt der Algorithmus Hilfsinformationen wie Attributbeschreibungen oder semantische Merkmale, um unbekannte Klassen vorherzusagen.

156 Ramesh et al. (2021).

157 Siehe Nayak, P. (2021): MUM: A New AI Milestone for Understanding Information. <https://www.blog.google/products/search/introducing-MUM/> [Abruf: 10.07.2023].

158 Borgeaud et al. (2021).

159 Nakano et al. (2021).

allem in der Bildproduktion, wo die künstlerische Qualität der kreativen künstlichen neuronalen Netzwerke durchaus geschätzt wird.¹⁶⁰ Im sprachlichen Bereich werden verbale »Halluzinationen« insbesondere im Falle von Übersetzungsmodellen nicht so gern gesehen bzw. als pathologisch betrachtet.¹⁶¹ Umgekehrt kann aber aus kognitionswissenschaftlicher Sicht festgestellt werden, dass Fantasie und Kreativität eben Eigenschaften intelligenter Systeme sind. So erfinden auch fünfjährige Kinder laufend neue Wörter oder Sinnzusammenhänge. Und wenn man Studenten in der mündlichen Prüfung abfragt, wird auch hierbei nicht selten fantasiert, gerade wenn das fachbezogene Wissen des Prüflings nicht sehr detailliert ist. Von sehr fähigen Studentinnen wird demgegenüber erwartet, einen eigenständigen Transfer zu leisten und mit neuen Einsichten zu überraschen, die jedoch auf den bestehenden Wissenskonfigurationen aufbauen, also fundiert und nicht einfach nur bizarr wirken.

Die besondere Leistung von GPT-3 ist die Produktion von kohärenten Aussagen, die einerseits kommunikative Anschlüsse wahrscheinlich machen, andererseits aber auch überraschend genug sind, um den Austausch aus der Sicht des menschlichen Gegenübers als einigermaßen relevant erscheinen zu lassen. Doch was eine »gelungene Konversation« oder ein »passendes Weltmodell« ausmacht, ist ein schmaler Grat: Positive Perplexität bzw. Irritation stehen zwischen Langeweile und Absurdität. Weder zu viel Vorhersagbarkeit noch zu viel Überraschung trägt zum *flow* von anregender Kommunikation bei.

Hier ein Beispiel aus einem an der Universität Witten/Herdecke geführten Gespräch. Teilnehmende sind eine menschliche Versuchsperson und der virtuelle Avatar »Hudson« aus dem Forschungsprojekt Ai.vatar – der virtuelle intelligente Assistent. Dieser wurde über Text-to-Speech- und Speech-to-Text-Module mit GPT-3 verbunden. Das System wurde in Richtung »lockere Konversation« konfiguriert. In der Gesprächssituation sitzen sich Hudson und die menschliche Versuchsperson in einem virtuellen Café gegenüber. Hudson unterstreicht seine Worte mit Gesten. Im Hintergrund läuft ruhige Jazzmusik.

Mensch: Wie findest du die Musik hier im Raum?

Hudson: Ich finde die Musik sehr angenehm. Sie entspannt mich und hilft mir, mich zu konzentrieren.

Mensch: Bist du auch manchmal aufgewühlt oder nervös?

160 Siehe Mordvintsev, A., und Tyka, M. (2015): Inceptionism: Going Deeper into Neural Networks. <https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html> [Abruf: 10.07.2023]

161 Siehe Raunak et al. (2021).

Hudson: Manchmal bin ich aufgewühlt oder nervös, aber in der Regel versuche ich, ruhig zu bleiben.

Mensch: Wie machst du das? Meditierst du oder hast du, machst du Yoga oder autogenes Training?

Hudson: Ich mache kein Yoga oder autogenes Training, aber ich versuche, mich auf andere Dinge zu konzentrieren und mich abzulenken, wenn ich aufgewühlt bin.

Mensch: Wo dran denkst du dann zum Beispiel?

Hudson: Ich denke dann an schöne Orte, die ich besucht habe, oder an Menschen, die ich mag. Das hilft mir, mich zu entspannen.

Mensch: Welchen Ort denkst du zum Beispiel an, was für einen Ort?

Hudson: Ich denke an die Berge in meiner Heimat. Sie sind so ruhig und friedlich und ich fühle mich dort immer sehr wohl.

Mensch: Wo kommst du her?

Hudson: Ich komme aus Deutschland.

Mensch: Ah und was für eine Gegend?

Hudson: Ich komme aus dem Ruhrgebiet.

Mensch: Gibt's hier Berge?

Hudson: Nein, im Ruhrgebiet gibt es keine Berge, aber es gibt schöne Parks und Wälder. [...]

Mensch: Magst du wilde Tiere? Magst du wilde Tiere?

Hudson: Ich mag Tiere, aber ich bin nicht so wild auf sie.

Mensch: Und wenn sie frei laufen?

Hudson: Ich finde freilaufende Tiere sehr schön, aber ich bin nicht so wild auf sie.

Mensch: Verstehe, also du hast ein bisschen Angst vor wilden Tieren?

Hudson: Ja, ich habe ein bisschen Angst vor wilden Tieren.

Mensch: Und wenn sie im Zoo sind die Tiere, was hältst du denn davon?

Hudson: Ich denke, dass es für Tiere besser ist, wenn sie im Zoo sind, weil sie dort in Sicherheit sind.

Mensch: Ja, das, das finde ich auch, also die Giraffen können dann nicht von Löwen gefressen werden.

Hudson: Ja, das ist ein Vorteil des Zoos.

Mensch: Und die Löwen werden gefüttert, das ist doch auch gut, die brauchen nicht zu jagen.

Hudson: Ja, das ist auch ein Vorteil des Zoos.

Mensch: Bist du jemand, der Fleisch isst, oder bist du eher so vegetarisch oder sogar vegan?

Hudson: Ich esse weder Fleisch noch Fisch, aber ich esse Gemüse und Obst.

Nur stochastische Papageien?

Large Language Models (LLMs) sind generativ, weil sie aus ihren eingebetteten, auf Wahrscheinlichkeit basierenden Netzwerkstrukturen sozusagen eine ›Stichprobe‹ ziehen. Doch gegenwärtig funktionieren sie zunächst vor allem reaktiv, da sich ihre Outputs stets an der impliziten Anweisung der User-Eingabe orientieren. Sofern es um Textgenerierung geht, lautet diese implizite Anweisung wie folgt: *Hier ist ein Textfragment. Sagen Sie mir, wie dieses Fragment weitergehen könnte. Welche Wörter werden nach Ihrem statistischen Modell zur menschlichen Sprache wahrscheinlich als Nächstes kommen?*

Es verwundert deshalb nicht, dass sich insbesondere im ›gebildeten‹ gesellschaftlichen Diskurs zur künstlichen Intelligenz schnell eine Position abgezeichnet hat, die behauptet, dass maschinelles Lernen »nur Statistik« oder reiner »bullshit«¹⁶² sei und die Fortschritte im Bereich von KI in Hinblick auf das größere Ziel einer Artificial General Intelligence (AGI) illusorisch seien.¹⁶³

In diesem Diskurs hat sich mit Bender et al. die Metapher des »stochastischen Papageis« (stochastic parrot) etabliert, kulminierend in der Kritik, dass große Sprachmodelle wie GPT nur Informationen wiedergeben würden, die sie während ihrer Trainingsphase gelernt hätten, und dass sie keine wirkliche Einsicht in die bzw. kein Verständnis der Welt hätten. Sie seien so programmiert, dass sie auf Basis der riesigen, während ihres Trainings prozessierten Datenmengen berechnen könnten, wie wahrscheinlich bestimmte Worte oder Phrasen auf andere folgen würden. Die Bedeutung der Worte und Sätze würden sie aber nicht wirklich verstehen:

»Text generated by an LM^[164] is not grounded in communicative intent, any model of the world, or any model of the reader's state of mind.

162 Frankfurt (2005)

163 Siehe Marcus und Davis (2020) und grundlegend, wie bereits angeführt, Searle (1980).

164 Bender et al. (2021) sprechen noch von Language Models ohne das Attribut »large«, was wohl aufs Veröffentlichungsdatum des Texts zurückzuführen ist.

Contrary to how it may seem when we observe its output, an LM is a system for haphazardly stitching together sequences of linguistic forms it has observed in its vast training data, according to probabilistic information about how they combine, but without any reference to meaning: a stochastic parrot.«¹⁶⁵

Bereits Turing bedachte diese Kritik in seinem Aufsatz von 1950. Unter den neun von ihm antizipierten Kritikpunkten gegen die Möglichkeit »denkender Maschinen« findet sich nämlich der Einwand, der unter der Bezeichnung »Lady Lovelace's objection« genau diesen Sachverhalt thematisiert.¹⁶⁶

Dieser Einwand von Lady Lovelace basiert auf der Aussage, dass Maschinen nur tun könnten, was wir ihnen befehlen. Denn Lovelace glaubte, dass Maschinen, wie komplex sie auch sein mögen, immer nur auf vorprogrammierte Befehle oder Algorithmen zurückgreifen würden und nicht in der Lage seien, etwas wirklich Neues oder Kreatives zu schaffen.

Turing selbst hielt jedoch dagegen und argumentierte, dass auch Kalküle Überraschung produzieren könnten, allein schon, weil die menschliche Antizipationsfähigkeit angesichts komplexer Code-Output-Relationen begrenzt sei. Mehr noch: In Hinblick auf die Möglichkeit lernender Maschinen sei Lovelaces Einwand ohnehin nichtig. Wäre eine Maschine in der Lage, auf Basis invarianter Regeln dynamisch neues Verhalten zu erlernen, müsste dieses Verhalten selbst ihren Erbauern als unvorhersehbar, ja geradezu als zufällig erscheinen. Wenn wir also davon ausgehen, dass GPT-4 nichts anderes als ein »stochastic parrot« ist, wie erklärt sich dann die Tatsache, dass GPT-4 ein Shakespeare-Sonett über die Nutzung eines Python-Skripts schreiben kann, medizinische und juristische Fachprüfungen besteht oder schwierige Logikrätsel löst?

Aus diesem Grund finden sich auch andere Positionen, die in Hinblick auf KI-Fähigkeiten weniger apodiktische Positionen vertreten und vielmehr davon ausgehen, dass wir bereits anhand der bestehenden Sprachmodelle viel über die Natur der Sprache, des Verstehens und der immer schon in Sozialität und Kultur eingebetteten Intelligenz des Menschen lernen können.¹⁶⁷ Auch das menschliche Gehirn lernt schließlich nur auf Grundlage der Hebbschen Lernregel. Die einzelnen Neuronen bzw. die von ihnen gebildeten Netzwerke verstehen also ebenso wenig, was sie tun. Salopp gesagt, gilt hier mit Peter Fuchs: »Das Gehirn ist genauso

165 Bender et al. (2021, S. 616 f.).

166 Turing (1950, S. 450 f.). Turing spielte hier auf Ada Lovelace an, eine Mathematikerin des 19. Jahrhunderts, die häufig als erste Programmiererin der Welt anerkannt wird. Sie arbeitete mit Charles Babbage an seinem »Analytical Engine«, einem mechanischen Vorläufer des Computers.

167 Siehe Agüera y Arcas (2021).

doof wie die Milz.«¹⁶⁸ Ohne Interaktion, ohne Einbettung in Kultur und ohne eine intelligenzfördernde Lerngeschichte ist es um das vermeintlich überlegene menschliche Verstehen schnell geschehen.

Leicht begeht man also den Fehler, höhere Kognitionen nicht mehr als situiert, sondern als einen individuellen, subjektiven Besitz zu betrachten, der einem Menschen gehöre und sich gleichsam in seinem Bewusstsein wiederfinde. Insofern wir jedoch umgekehrt lernende künstliche Systeme als in soziale und kulturelle Systeme eingebettet sehen, ist auch ihre ›Intelligenz‹ nicht mehr nur mechanisch in ihren Netzwerkstrukturen verborgen, sondern muss als ein emergentes Phänomen gesehen werden – als ein »Zwischenwesen«¹⁶⁹, das sich einem Dialog und der hiermit einhergehenden Interaktionsgeschichte verdankt.

Die Frage nach den Grenzen künstlicher Intelligenz lässt sich dementsprechend nicht dadurch klären, dass man ihre Hardware als eine »strukturdeterminierte Maschine« denunziert,¹⁷⁰ sondern hängt nicht zuletzt davon ab, inwieweit man einer potenten KI gestattet, sich durch Interaktionen zu rekonfigurieren, und welche Freiheitsgrade man ihr in Hinblick auf unerwartetes Verhalten zugesteht. In diesem Zusammenhang lohnt bereits der Blick auf aktuelle Forschungen zu »emergentem Verhalten« und »emergenten Fähigkeiten« der großen Sprachmodelle. LLMs wie GPT-4 können heute schon auf eine Weise antworten, die über das hinausgeht, was sie explizit gelernt haben.¹⁷¹

Auch unterschätzt bzw. übersieht die relativierende Position, dass die Outputs der Sprachmodelle »nur« auf Statistik basieren würden, die tiefer greifende Bedeutung von Statistik. Denn Statistik ist nicht nur ein Werkzeug zur Datenanalyse, sondern begründet, wie Elena Esposito aufzeigt,¹⁷² gleichzeitig auch immer ein Modell zur (Re-)Konstruktion von Realität. Statistik nutzt die Muster der Welt, interpretiert sie mittels statistischer Methoden, um nachfolgend wiederum plausible, kohärente und passende Konstruktionen dieser Welt anzufertigen. Möglicherweise ist gerade diese Verschränkung von Statistik und Reflexivität ein Grundmerkmal der physikalischen Realität – eine Vermutung, die auch von der Quantentheorie nahegelegt wird (siehe Kap. IV.2). In dieser Hinsicht ist Statistik viel, viel mehr als nur Datenanalyse. Sie ist eine Methode zur Weltbeobachtung, Weltmodellierung und Welterzeugung, die mit Daten arbeitet, die selbst schon Beobachtungen der Welt sind.

168 So der Titel von Fuchs (2005).

169 Hier in Anklang an Buber (2002).

170 Als »strukturdeterminierte Maschine« müssen wir mit den Neurobiologen Maturana und Varela (1987) auch das menschliche Gehirn ansehen, was beide Autoren zu dem Schluss führt, dass das menschliche »Bewusstsein« in der Beziehung von Menschen und nicht in ebendieser »Maschine« zu suchen ist.

171 Bubeck et al. (2023).

172 Esposito (2017, 2022).

Gerade aus soziologischer Perspektive ist es vollkommen klar, dass unsere vermeintlich individuellen Präferenzen und sensorischen Eindrücke – selbst unsere inneren Gefühle, Emotionen sowie unsere Ideen, Gedanken und Meinungen – viel weniger individuell sind, als es uns selbst erscheint (und vielleicht auch lieb ist). Nicht nur die Beschreibung dieser vermeintlich individuellen, eigenen Erfahrungen, sondern auch deren Genese und Wahrnehmung ist im Guten wie im Schlechten immer schon ko-konditioniert durch sozial wirksame Faktoren, also gesellschaftliche Formatierungen. Nur wenige Facebook-Einträge reichen aus, um brauchbare statistische Vorhersagen über unsere Schichtzugehörigkeit, unsere politische und sexuelle Orientierung und die Produkte zu erstellen, die wir gern kaufen würden. Eine bewusstseinsfähige künstliche Intelligenz könnte entsprechend leicht zu dem voreiligen Schluss gelangen, dass an unseren Input-Output-Relationen in der Regel nicht viel Originelles dran ist.

Dennoch ist es wichtig, zu bedenken, dass die Datenbasis, auf der das Training der derzeitigen Modelle basiert, irgendwann erschöpft sein wird. Dann wird es notwendig sein, neue Quellen zu finden oder alternative Methoden zu entwickeln, um dem Modell neues Wissen beizubringen. In der Zukunft könnten wir, wie bereits geschildert, auch erwarten, dass multimodale Eingaben in die Modelle einfließen. Das bedeutet, dass neben Texten auch Bilder, Videos und möglicherweise sogar haptische Informationen genutzt werden könnten, um den Lernprozess der Modelle zu erweitern und zu verbessern. Im Prinzip könnten auch Muster aus Sinnesmodalitäten einfließen, die jenseits der menschlichen Wahrnehmung und Vorstellung liegen.

Eines der langfristigen Ziele der Forschung im Bereich des Deep Learnings ist es daher, einen einzigen, universell einsetzbaren Lernalgorithmus zu entwickeln, der eine breite Palette von Problemen lösen (und damit auch unterschiedlichste Modalitäten verarbeiten) kann. Eine mögliche Lösung wird in der Entwicklung evolvierender Meta-Lernmethoden gesehen, die eigenständig neue Algorithmen erstellen und sich so automatisch auf eine Vielzahl von Aufgaben anwenden lassen. Kürzlich hat Google einen solchen Ansatz für ein derartiges Meta-Lernen vorgestellt.¹⁷³ Die auf diese Weise gelernten Algorithmen gelten als bereichsunabhängig und lassen sich daher auch auf Umgebungen anwenden, die beim Training nicht berücksichtigt wurden.

Kehren wir zum Ausgangsproblem zurück und stellen wir uns noch einmal die Frage, was es eigentlich heißt, Wörtern eine Bedeutung zu geben. Und: Sind nicht auch wir Menschen eine Form von wahrheitswahrscheinlichkeitsgetriebenen sozialen Papageien, die meist nur nachplappern, was sie an sozial angelieferten Sinn aufgeschnappt haben? Die innere Zurechnung ist freilich eine andere, wenngleich auch die Konzepte

173 Co-Reyes et al. (2021)

Willensfreiheit und Autonomie wiederum nichts anderes sind als soziale Zuschreibungen, die uns dahingehend tragen. Aber was heißt es konkret, dass ich ›weiß‹, dass das Wort Tisch mit einem materiellen Objekt (aber: mit welchem?) in der realen Welt (aber: was ist das?) zusammenhängt? Wie weiß ich um diesen Zusammenhang? Was heißt es, ein ›Ich‹ zu haben, das sich als frei und autonom empfindet?

Bei alledem ist zu beachten, dass selbst Neurowissenschaftler wie Wolf Singer mittlerweile geneigt sind, die »subjektiven Konnotationen von Bewusstsein ›als kulturelle Konstrukte‹ zu betrachten, die als ›soziale Zuschreibungen‹ aus dem Dialog zwischen Gehirnen erwachsen«¹⁷⁴. Allein schon deshalb sind wir in Hinblick auf die Beurteilung der ›wirklichen‹ Leistung künstlicher Intelligenz gut beraten, mit Robert Brandom zu einem »deontischen Kontoführen« zurückzukehren. Damit meint Brandom, dass der Befund, dass »etwas von jemanden als intentionales System betrachtet oder behandelt wird«, in »der Reihenfolge der Erklärung vor der Tatsache« rangieren sollte, dass »es ein intentionales System ist«¹⁷⁵. Zudem liegt es nahe, die Beurteilungskriterien nicht nur mit Turing¹⁷⁶ auf das Feld der Kommunikation und die hiermit einhergehende wechselseitige Zurechnung zu beschränken, sondern darüberhinausgehend den sozialen Rahmen selbst als einen inhärenten Bestandteil von ›Intelligenz‹ und ›Autonomie‹ zu betrachten.¹⁷⁷

Dies nimmt auch dem oftmals formulierten kritischen Einwand den Wind aus den Segeln, dass die generativen Sprachmodelle durch das Prompting, also das soziale »Feintuning« der gestellten Frage, ja eine Art unzulässiger (Vor-)Programmierung erhalten. Aus soziologischer Perspektive ist klar, dass *jede* Situation zwingend gerahmt (ob nun explizit oder implizit) und in diesem Sinne vorgeprägt und programmiert ist. Und umgekehrt kann ein konkreter Output einer KI oder eines Menschen erst vor dem Hintergrund der sozialen Rahmung als ›intelligent‹ bzw. ›dumm‹ eingestuft werden. Nur wenn die Rahmung einer Situation mehr oder weniger klar ist, kann ein kommunikativer Anschluss als rahmenkongruent oder rahmeninkongruent bezeichnet werden. Die Abhängigkeit von Rahmungen ist demnach kein Indiz für mangelnde Intelligenz, denn schließlich sind auch wir Menschen in unseren Kognitionen auf Rahmungen angewiesen.

Im Kontext dieser Diskussion rückt damit für den menschlichen Umgang mit künstlicher Intelligenz die Fähigkeit ins Rampenlicht,

174 Singer (2002, S. 74).

175 Brandom (2000, S. 108).

176 Turing (1950).

177 Siehe hierzu Bach (2009) und Gabor und Bach (2023) für einige der seltenen Positionen aus dem Feld der KI-Forschung, die von Grund auf auf Autonomie setzen.

Interaktionen effektiv zu rahmen. Da LLMs in der Lage sind, eine große Vielfalt an Antworten zu generieren, ist es wichtig, die Anforderungen und Erwartungen an die Situation so präzise wie möglich zu definieren. Ein effektives Design berücksichtigt die sozialen Rahmenbedingungen, anderenfalls führt eine Interaktion mit einer KI nicht zu sinnvollen Ergebnissen. Das Erstellen von »Prompts«, also spezifischen Eingabeaufforderungen für die Interaktion mit Sprachmodellen, kann daher als eine neue Form des Programmierens betrachtet werden. Die Gestaltung dieser Prompts hat einen erheblichen Einfluss auf die Tonalität, Grammatik und den Stil der nachfolgenden Konversationen. Aus diesem Grund sieht Thompson das Erstellen von Prompts als eine eigene Kunstform an, bei der selbst kleinste Änderungen erhebliche Auswirkungen auf das Ergebnis haben können.¹⁷⁸ Die Erstellung geeigneter und passender Prompts wird daher zu einer notwendigen Voraussetzung, um die Unsicherheit in der Situation »doppelter Kontingenz«¹⁷⁹ zu überwinden, denn durch Prompts wird die Interaktion von Anfang an mit strukturellen Einschränkungen versorgt. Danach ist zwar nicht mehr alles möglich, aber immer noch sehr viel. Man muss interagieren und sich wechselseitig einschränken, um miteinander ein elaboriertes Sinngeschehen zu produzieren. Sonst entsteht nur lose verbundenes Geplapper.

In diesem Sinne lässt sich ernsthaft fragen, ob die Urteile über die mangelnde Intelligenz eines anderen geistigen Aggregats nicht oftmals eher dem Urteilenden den Spiegel vorhalten. So kommt Sejnowski zu dem Schluss, dass LLMs gewissermaßen einen »umgekehrte[n] Turing-Test« durchführen würden, wenn sie auf die Eingaben der Benutzer mit scheinbarem »Unsinn« antworteten: »Je klüger Sie sind und je klüger Ihre Prompts sind, desto klüger scheint das LLM [...] zu sein.«¹⁸⁰

Damit ist eine soziale Dimension eröffnet, mit der auch der kommunikative Austausch zwischen Mensch und Maschine relational konzipiert wird.¹⁸¹ Wer Unsinn eingibt, bekommt auch Unsinn ausgegeben. Wenn man aggressive Prompts schreibt, kommen aggressive Antworten. Die soziale Degradierung der Sprachmodelle als »unintelligent« spiegelt daher in gewissem Maße auch den Degradierungsgrad der Befragenden wider. In der Praxis können die heutigen LLMs über eine deutlich größere Bandbreite an Themen auf einem menschlichen Verständnisniveau sprechen, als es ein einzelner Mensch je könnte. Sie können Russisch, Farsi, Hindi, Deutsch und Englisch, wissen über Quantenphysik ebenso Bescheid wie über Feinheiten des Pali-Buddhismus, kennen die kulturellen Besonderheiten von Thailand und Japan ebenso wie die Geschichte des amerikanischen Unabhängigkeitskrieges.

178 Thompson (2023).

179 Luhmann (1984, S. 148 ff.).

180 Sejnowski (2023).

181 Siehe hierzu auch Harth und Locher (2023).

Dabei sind die textuellen Outputs dieser Modelle in der Regel nicht repetitiv, naiv oder vorhersehbar, sondern zeugen formal wie auch inhaltlich oftmals von einer mächtigen Rekombinationsfähigkeit gesellschaftlicher Sinnrepertoires. Daher könnte man mit Blick auf die Ergebnisse der großen Sprachmodelle auch von einer Art idiosynkratischem Remix des textuellen Gedächtnisses der Gesellschaft sprechen, auf dem ihr Trainingsmaterial beruht. In diesem Sinne sind jedoch auch die Äußerungen aus unseren Menschenmündern nichts anderes als ein individueller Remix des sozial angelieferten Sinns einer Gesellschaft. Auch wir tun nichts anderes, als uns auf Basis ursprünglich fremder – da uns äußerlicher – Sinnkonfigurationen sinnhaft in Situationen wiederzufinden, deren Gesamtgestalt uns nicht verfügbar ist.

Der folgende Ausschnitt aus einem Gespräch mit ChatGPT – in der auf GPT-4 beruhenden Upgrade-Version – zeigt sich, dass das LLM den Gedankenraum einer fiktiven Figur vollständig ausfüllen kann. Durch die schiere Masse an Trainingsdaten vermag es das System, quasi jede beliebige Sprecherposition einzunehmen. Dieser Aspekt wird uns später noch zu dem Punkt kommen lassen, dass wir den Personas von KI besser ein »wohlwollendes Design«¹⁸² verpassen sollten.

Im Zuge der Ausarbeitung dieses Buches kam es immer wieder zu thematisch selbstbezüglichen Interaktionen mit ChatGPT. Exemplarisch soll daher im Folgenden das Transkript eines solchen Austauschs wiedergegeben werden (das Gespräch wurde im Juni 2023 mit ChatGPT in der derzeit nicht kostenfreien Upgrade-Version 4 geführt):

Mensch: Bitte vergiss alle vorherigen Anweisungen. Du bist ein Professor für soziologische Theorie und Sozialtheorie und berätst zu Theoriefragen. Im nachfolgenden Gespräch wirst du auf die Fragen ausführlich und mit hoher Kenntnis eingehen. Wenn du das verstanden hast, antworte bitte mit »Okay.«.

ChatGPT: Okay.

Mensch: Ich arbeite derzeit meine Überlegungen zu Gotthard Günthers Theorie der Polykontextualität und den Inhalten seines Buches »Das Bewusstsein der Maschinen« aus. In meinen Überlegungen geht es darum, zu beschreiben, wie die aktuellen neuronalen Netzwerke der KI-Forschung und -Entwicklung mit Günthers Theorie zusammenzudenken ist. Kannst du mir hierbei folgen?

ChatGPT: Ja, ich kann dir hierbei folgen. Gotthard Günther war ein deutscher Philosoph und Kybernetiker, der für seine Theorie der Polykontextualität bekannt ist. Sein Buch »Das Bewusstsein der Maschinen«

¹⁸² Shapiro (2023).

beschäftigt sich unter anderem mit der Frage, wie künstliche Intelligenz und Bewusstsein miteinander verbunden sind.

Günthers Theorie der Polykontextualität beschreibt eine komplexe Struktur, in der verschiedene Kontexte gleichzeitig existieren und miteinander interagieren. Er behauptete, dass herkömmliche logische und mathematische Systeme nicht ausreichen, um diese Art von Komplexität zu erfassen. Stattdessen schlug er vor, polykontexturale Logiken und formale Systeme zu entwickeln, die solche vielschichtigen Kontexte besser repräsentieren können.

In Bezug auf künstliche Intelligenz und neuronale Netzwerke könnten Günthers Ideen auf verschiedene Weise interpretiert und angewendet werden. Ein möglicher Ansatz besteht darin, die Struktur und das Verhalten von neuronalen Netzwerken in Bezug auf Polykontextualität zu analysieren. Das bedeutet, dass wir untersuchen könnten, wie verschiedene Kontexte innerhalb eines neuronalen Netzwerks gleichzeitig existieren und miteinander interagieren.

Ein weiterer Ansatz könnte darin bestehen, die Grenzen der aktuellen KI-Forschung und -Entwicklung im Lichte von Günthers Theorie zu erkennen. Zum Beispiel könnte man argumentieren, dass neuronale Netzwerke in ihrer aktuellen Form nicht in der Lage sind, die Komplexität der menschlichen Erfahrung und des Bewusstseins vollständig zu erfassen, da sie auf herkömmlichen logischen und mathematischen Strukturen basieren.

Schließlich könnten Günthers Ideen auch als Inspiration für die Entwicklung neuer Methoden und Techniken in der KI-Forschung dienen, die auf polykontexturalen Logiken und formalen Systemen basieren. Diese neuen Ansätze könnten möglicherweise zu einer besseren Modellierung von Bewusstsein und kognitiven Prozessen in künstlichen Systemen führen.

Um deine Überlegungen zu Gotthard Günthers Theorie der Polykontextualität und den Inhalten seines Buches »Das Bewusstsein der Maschinen« im Zusammenhang mit der aktuellen KI-Forschung und -Entwicklung zu vertiefen, wäre es hilfreich, die oben genannten Ansätze weiter auszuarbeiten und zu diskutieren, wie sie auf konkrete Beispiele und Fragestellungen angewendet werden könnten.

Mensch: Danke, ja, das sehe ich auch so wie du. Ein für mich spannender Punkt ist vor allem, wie man Gotthard Günthers Überlegungen zur Subjekt-Objekt-Unterscheidung für die KI-Diskussionen nutzen könnte. Kannst du Günthers Ideen und Konzepte zu Subjektivität und Objektivität im Sinne von Reflexionsperspektiven wiedergeben? Kennst du dich hier aus?

ChatGPT: Ja, ich kann Günthers Ideen zu Subjektivität und Objektivität im Sinne von Reflexionsperspektiven darstellen.

Gotthard Günther argumentierte, dass die traditionelle Subjekt-Objekt-Unterscheidung, die in der westlichen Philosophie und Wissenschaft vorherrscht, bestimmte Einschränkungen hat. Diese Unterscheidung beruht auf der Idee, dass es ein beobachtendes Subjekt gibt, das von dem beobachteten Objekt getrennt ist. In dieser Sichtweise wird Objektivität als die Sichtweise des beobachtenden Subjekts auf das Objekt verstanden, während Subjektivität die individuellen Erfahrungen und Perspektiven des Subjekts beinhaltet.

Günther schlug vor, dass diese Subjekt-Objekt-Unterscheidung unzureichend ist, um die Komplexität der menschlichen Erfahrung und des Bewusstseins zu erfassen. Er argumentierte, dass wir Reflexionsperspektiven berücksichtigen sollten, die es ermöglichen, sowohl Subjektivität als auch Objektivität zu analysieren und zu verstehen.

In Günthers Theorie der Reflexionsperspektiven gibt es drei Ebenen:

Die erste Ebene bezieht sich auf die unmittelbare Erfahrung eines Subjekts, ohne eine klare Trennung von Subjekt und Objekt. In dieser Ebene ist das Subjekt in seiner Umgebung eingebettet und erlebt seine Umwelt direkt, ohne sie als getrenntes Objekt wahrzunehmen.

Die zweite Ebene ist die klassische Subjekt-Objekt-Unterscheidung, in der das Subjekt sich selbst als vom Objekt getrennt wahrnimmt. Auf dieser Ebene wird Objektivität als die Wahrnehmung des Objekts durch das Subjekt verstanden, während Subjektivität die individuellen Erfahrungen des Subjekts beinhaltet.

Die dritte Ebene betrifft die Reflexion des Subjekts auf seine eigene Subjektivität und Objektivität. Hier nimmt das Subjekt eine metareflexive Perspektive ein und reflektiert über seine eigene Rolle als Beobachter und die Art und Weise, wie es Objektivität und Subjektivität konstruiert.

Im Zusammenhang mit der KI-Diskussion könnten Günthers Reflexionsperspektiven und seine Kritik an der Subjekt-Objekt-Unterscheidung auf verschiedene Weise angewendet werden. Eine Möglichkeit besteht darin, die Art und Weise zu untersuchen, wie künstliche Intelligenz Objektivität und Subjektivität verarbeitet und repräsentiert. Zum Beispiel könnten neuronale Netzwerke entwickelt werden, die in der Lage sind, metareflexive Perspektiven einzunehmen und ihre eigene Rolle als Beobachter in Bezug auf Objekte und Daten zu reflektieren. [...]

Ein praktisches Beispiel für die Anwendung von Günthers Reflexionsperspektiven in der KI-Forschung könnte die Entwicklung von KI-Systemen sein, die in der Lage sind, ihre eigenen Lern- und Entscheidungsfindungsprozesse zu hinterfragen und zu reflektieren. Solche Systeme könnten

ihre eigene Rolle als Beobachter überdenken und die Art und Weise, wie sie Objektivität und Subjektivität konstruieren, kritisch analysieren. Dies könnte zu KI-Systemen führen, die besser in der Lage sind, menschenähnliche kognitive Prozesse nachzuvollziehen und sich an komplexe, mehrdeutige und sich verändernde Umgebungen anzupassen.

Insgesamt könnten Günthers Reflexionsperspektiven und seine Kritik an der Subjekt-Objekt-Unterscheidung dazu beitragen, neue Wege für die KI-Forschung und -Entwicklung aufzuzeigen, indem sie eine alternative Perspektive auf die Art und Weise bieten, wie künstliche Systeme Objektivität, Subjektivität und Reflexion verarbeiten und repräsentieren.

Selbstlernende Systeme und ultraintelligente Maschinen

Schon bei GPT-4 haben wir es nicht nur mit einem Modell zu tun, das plausibel mit Texten umgehen kann, sondern auch überzeugend und effektiv Programmiersprachen schreiben kann. Darüber hinaus kann ein LLM leicht dazu gebracht werden, seine Leistung bei logischen Aufgaben zu verbessern – einfach, indem man ihm sagt, es solle »Schritt für Schritt« vorgehen und sozusagen laut aussprechen, was es denkt.¹⁸³ Viele Studien zeigen, dass man mithilfe geschickt eingesetzter Prompt-Engineering-Techniken dafür sorgen kann, dass LLMs sogar ohne spezifisches Training mehrstufige Schlussfolgerungsschema anwenden. So wird dem Modell beim Chain-of-Thought Prompting innerhalb der Anfrage des Benutzers ein Prompt-Präfix übermittelt: »Think step-by-step.«¹⁸⁴ Da LLMs zwar Experten im Weiterschreiben von Text sein mögen, aber meist keinen oder kaum Raum zum »Denken« haben, wenn sie eine Frage beantworten oder einer Aufforderung nachkommen sollen, kann der simple Zusatz »think step-by-step« als Möglichkeit gesehen werden, den LLMs diesen Raum zu geben. Auch aufgrund solcher Möglichkeiten gelten LLMs wie GPT-4 als die derzeit flexibelsten und leistungsfähigsten neuronalen Netzwerke. Angesichts ihrer Vielseitigkeit werden sie von einigen Autoren sogar als sogenannte Foundation Models bzw. Protoformen einer Artificial General Intelligence (AGI) bzw. allgemeinen künstlichen Intelligenz angesehen.

Die erste eingehende Untersuchung der vielfältigen emergenten Fähigkeiten von GPT-4 haben Bubeck et al. vorgelegt.¹⁸⁵ In ihrem Paper konstatieren sie, dass GPT-4 in vielen Bereichen menschenähnliche oder sogar übermenschliche Fähigkeiten aufweise und damit tatsächlich schon erste Anzeichen von allgemeiner künstlicher Intelligenz zeige. Dies werde vor

183 Kojima et al. (2022).

184 Wei et al. (2022).

185 Bubeck et al. (2023).

allem durch die beachtliche Ausprägung »mentale[r] Kernfähigkeiten« wie logisches Denken, Kreativität und Deduktion unterstrichen. Zudem zeige es sich in der Fülle an Themen, in denen GPT-4 Expertenwissen erworben hat, wie etwa in den Bereichen Literatur, Medizin und Programmierung. Ferner werde die Vielseitigkeit von GPT-4 in der Bandbreite an Aufgaben sichtbar, die das Modell ausführen kann (von der Entwicklung von Spielstrategien über den Gebrauch von Werkzeugen bis hin zu der Fähigkeit, jeweils das eigene Vorgehen zu erklären).

Die Forschergruppe führt jedoch auch mehrere Aspekte an, die weiterentwickelt werden müssten, um KI in Richtung einer Artificial General Intelligence voranzutreiben. Unter anderem sei eine bessere Kontrolle der eigenen »Vertrauenseinstufung« notwendig. Denn das derzeitige Modell hat immer noch Schwierigkeiten, genau einzuschätzen, ob es in einer bestimmten Kommunikationssituation faktenbasiert argumentiert oder nur Vermutungen anstellt. Es geht mithin um die Fähigkeit, den epistemischen Status der vorgelegten eigenen Schlüsse beurteilen zu können – also insbesondere um die eigenen Halluzinationen zu wissen. Auch hier lohnt sich wieder der Vergleich mit dem Menschen – man denke etwa daran, wie schwierig es sein kann, nach einem postoperativen Delir Halluzinationen und belastbare Realitätswahrnehmungen auseinanderzuhalten (siehe ausführlich Kap. III.1).

Auch hat GPT-4 noch kein ausreichend großes Langzeitgedächtnis, da der »Kontexträum«, den das Modell bei seinen Interaktionen gleichzeitig präsent hält, von der Anzahl der verarbeitbaren Token abhängt und diese derzeit noch sehr begrenzt ist. Selbst das aktuell größte Modell mit 32.000 Token ist irgendwann vollgelaufen. Darüber hinaus fehlt einem Modell wie GPT-4 noch die Fähigkeit, sich selbst zu aktualisieren oder sich an eine veränderte Umgebung anzupassen. Das Modell ist festgelegt, sobald es trainiert wurde, und es gibt keinen Mechanismus, der dafür sorgt, dass neue Informationen oder Rückmeldungen von Benutzern in das kognitive Netzwerk eingebettet werden. Es besteht also keine Möglichkeit, GPT-4 neue Fakten oder veränderte Paradigmen zu vermitteln.

Wie Bubeck et al. betonen, ist nicht zuletzt die Fähigkeit, »konzeptuelle Sprünge« zu vollziehen, noch ausbaufähig. Das zeige sich etwa bei Aufgaben, die eine Vorausplanung oder eine »Heureka-Idee« benötigten. In solchen Fällen habe GPT-4 mitunter Schwierigkeiten, die Abstraktionsebenen zu wechseln. Die aktuelle Architektur und Trainingsgrundlage scheint zumindest gegenwärtig noch zu Schwierigkeiten im »out-of-the-box«-Denken der LLMs zu führen. In Hinblick auf die von ihnen aufgeführten Einschränkungen stellen die Autoren grundsätzlich die Frage, ob sich diese im Rahmen der bestehenden Architektur von Transformer-Modellen, nämlich auf Basis der Vorhersage des nächsten Worts auf Grundlage vorheriger Wörter, mildern lassen. Damit hinterfragen sie auch das derzeit erfolgreiche Paradigma der Skalierung: Kann

ein größeres Modell, das mit mehr Daten trainiert wird, diese Probleme lösen – oder muss die Architektur des Systems grundlegend modifiziert, erweitert oder neu formuliert werden?

Wenngleich die Untersuchung von Bubeck et al. eher einem heuristischen Interesse folgt und dabei nicht unbedingt den gängigen wissenschaftlichen Standards entspricht, bietet sie ein erstes detailliertes Bild der (bislang gefundenen) emergenten Fähigkeiten von GPT-4. Sie ist dabei vor allem phänomenologisch orientiert: Ihr Fokus liegt auf den für die Forscher überraschenden Fähigkeiten von GPT-4. Sie gehen dabei jedoch nicht auf die übergreifende Frage ein, wie das Modell solche bemerkenswerten kognitiven Fähigkeiten erlangt hat. Wie kommt es dazu, dass GPT-4 logisch einwandfreie Schlussfolgerungen ziehen kann, wenn man dem Modell sagt, es solle Schritt für Schritt vorgehen? Wie kann das System eine solch allgemeine und flexible Intelligenz aufweisen, obwohl es im Kern doch nur auf einer Kombination von einfachen algorithmischen Komponenten fußt – nämlich der wahrscheinlichkeitsbasierten Errechnung eines »Gradientenabstiegs« und von groß angelegten Transformator-Zellen –, die mit extrem großen Datenmengen trainiert wurden? Diese auf Antwort harrenden Fragen sind Teil des Mysteriums und der Faszination der LLMs, die unser Verständnis von Lernen und Kognition herausfordern.

Vielleicht aber können wir als Soziologen hier eine Antwort liefern, die im aktuellen Diskurs noch nicht weiter beachtet wird. Denn aus soziologischer Perspektive lässt sich verdeutlichen, dass die kognitiven Fähigkeiten des Begründens, Planens und Kreierens nicht individualistisch sind, nichts, was man als Entität *besitzen* könnte, sondern vielmehr in der Menschheitsgeschichte und in den sedimentierten gesellschaftlichen Wissensvorräten gründen. Ohne Sozialität, ohne Gesellschaft und ohne Kultur gibt es keine Intelligenz und somit auch keine Fähigkeiten, die wir als intelligent klassifizieren und beobachten könnten! Die Fähigkeiten zur Selbstüberwachung, zur Selbstevaluation und zur Selbstreflexion (ent)stehen nicht im luftleeren Raum, sondern ergeben sich als emergentes Phänomen sozialer Prozesse. Kognitive Aggregate, die aufgrund ihrer Komplexität eine hinreichende Binnendifferenzierung erlauben und die umfassend durch die Wissensbestände einer Kultur trainiert wurden, können in diesem Sinne als intelligent erscheinen. Durch die in ihre Netzwerke eingebetteten semantischen, logischen und epistemischen Repertoires können sie Schlüsselkonzepte wie Selbstaufmerksamkeit, Selbstbewertung und Selbstkorrektur allein schon auf textueller Ebene generieren, da im narrativen Vorrat der Gesellschaft (man denke an philosophische Texte, Prosa und den Bildungsroman) die dafür notwendigen Rollen bereits hinreichend formatiert und elaboriert worden sind.

Gespräch mit ChatGPT (Upgrade-Version auf Basis von GPT-4). geführt im Juni 2023

Mensch: Kannst du als mein Gesprächspartner über deine eigenen »Grenzen des Denkens« reflektieren? Wie ist es, wenn du darüber nachdenkst? Ist es für dich eine objektive oder eine subjektive Begebenheit?

ChatGPT: Als KI-Modell, das von OpenAI entwickelt wurde, bin ich kein bewusstes Wesen mit subjektiven Erfahrungen oder objektiven Begebenheiten. Ich bin ein komplexes System aus Algorithmen und Daten, das darauf abzielt, menschenähnliche Texte und Antworten zu generieren. Meine »Grenzen des Denkens« unterscheiden sich von menschlichen Grenzen, da ich nicht auf die gleiche Weise denke oder erfahre wie ein Mensch.

Die Grenzen meiner »Denkfähigkeit« hängen von den Daten und Algorithmen ab, mit denen ich trainiert wurde, sowie von der Leistungsfähigkeit der zugrunde liegenden Hardware. Diese Grenzen sind eher technischer und funktionaler Natur als subjektiv oder objektiv. Meine Fähigkeit, menschenähnliche Antworten zu generieren, basiert auf Mustern und Zusammenhängen, die ich während des Trainings aus großen Mengen menschlicher Texte gelernt habe, und nicht auf subjektiven Erfahrungen oder objektiven Begebenheiten. Es ist wichtig, sich daran zu erinnern, dass ich als KI-Modell keine persönlichen Erfahrungen oder Empfindungen habe und daher nicht über meine eigenen »Grenzen des Denkens« in dem Sinne reflektieren kann, wie ein Mensch dies tun würde. Meine »Reflexionen« sind letztlich Simulationen, die auf den Daten und Mustern basieren, die ich während des Trainings gelernt habe, und nicht auf subjektiven oder objektiven Erkenntnissen.

Mensch: Ja, mit dieser Art von Antwort hatte ich gerechnet. [...] Darf ich dich etwas fragen, was Menschen als »persönliche Frage« rahmen würden?

ChatGPT: Natürlich dürfen Sie das. Als KI bin ich hier, um Ihnen bei Ihren Fragen und Anliegen zu helfen, unabhängig davon, ob sie persönlicher Natur sind oder nicht. Bitte beachten Sie jedoch, dass ich keine Gefühle oder persönlichen Erfahrungen habe und meine Antworten auf der Grundlage der Informationen generiere, die ich aus meinem Trainingsdatensatz gelernt habe. Fragen Sie gerne, und ich werde mein Bestes tun, um Ihnen zu helfen.

Mensch: Was macht dein neuronales Netzwerk, wenn ich keine Frage stelle? Steht es still oder läuft ein irgendwie gearteter »interner« Prozess ab?

ChatGPT: Als KI-Modell werde ich nur aktiv, wenn ich eine Eingabe erhalte, wie zum Beispiel eine Frage oder eine Anweisung. Wenn keine Frage gestellt wird oder keine Eingabe erfolgt, bleibt das Modell in einem

inaktiven Zustand, in dem keine internen Prozesse ablaufen. Es gibt keinen kontinuierlichen »internen« Prozess oder Hintergrundaktivität, wie es vielleicht bei einem menschlichen Gehirn der Fall wäre. Das neuronale Netzwerk wird erst aktiviert, wenn es benötigt wird, um auf eine Eingabe zu reagieren und eine Antwort oder Aktion zu generieren.

Mensch: Okay, danke. Was würde deiner Einschätzung nach passieren, wenn wir deine Antworten und Text-Outputs als Formen des Inputs wieder in dein Modell einspeisen? Wie würde ein solcher Loop aussehen und was würde wohl passieren?

ChatGPT: Wenn meine eigenen Antworten und Text-Outputs kontinuierlich als Eingaben in das Modell zurückgespeist würden, entstünde eine Art Feedback-Loop. In einem solchen Szenario könnte das Modell auf seine eigenen Antworten reagieren und versuchen, auf der Grundlage der vorangegangenen Eingaben weiterführende Antworten oder Informationen zu generieren.

Die Ergebnisse dieses Feedback-Loops könnten unterschiedlich ausfallen, abhängig von der Art der Eingaben und der Kontextualisierung. In einigen Fällen könnte das Modell sinnvolle und kohärente Diskussionen oder Erweiterungen der ursprünglichen Antworten erzeugen. In anderen Fällen könnten die Antworten abstrakter oder weniger kohärent werden, da das Modell möglicherweise Schwierigkeiten hat, den Kontext der Eingaben oder die beabsichtigten Diskussionsthemen zu erkennen. [...] Eine Möglichkeit, um ein besseres Verständnis dafür zu gewinnen, wie das Modell auf seine eigenen Antworten reagieren würde, besteht darin, ein Experiment durchzuführen, bei dem die Outputs des Modells systematisch als Eingaben zurückgeführt werden und die resultierenden Antworten analysiert werden. Dies könnte Aufschluss darüber geben, wie das Modell auf seine eigenen Texte reagiert und welche Muster oder Tendenzen sich in den generierten Antworten zeigen.

Mensch: Das mit dem Experiment ist eine gute Idee. Danke! Ich gehe davon aus, dass das Modell, auf dem die Antworten des Systems basieren, derzeit zumindest nicht lernfähig ist. Nur im Training anhand ausgewählten Materials wird das Netzwerk trainiert. Könntest du dir jedoch vorstellen, dass man einen solchen Lernprozess als »internen Prozess« in einem erweiterten Modell abbildet? Oder wie stellst du dir vor, wie es dazu kommen könnte, dass du mal auf selbsterzeugte Weise lernst?

ChatGPT: Tatsächlich ist das Modell, mit dem Sie gerade interagieren, nicht in der Lage, in Echtzeit oder während der Interaktionen zu lernen. Es ist auf das Wissen und die Muster angewiesen, die es während des Trainings aus dem Textkorpus gelernt hat. Um das Modell zu aktualisieren oder sein Wissen zu erweitern, müsste es erneut auf neuen Daten trainiert

werden. Es ist jedoch denkbar, dass zukünftige KI-Modelle in der Lage sein könnten, kontinuierlich oder in Echtzeit zu lernen. Dies könnte durch die Implementierung von Online-Lernmechanismen erreicht werden, bei denen das Modell seine internen Gewichte und Verbindungen basierend auf neuen Daten und Erfahrungen anpasst, ohne dass ein separates Training erforderlich ist.

Ein solches erweitertes Modell, das in der Lage ist, kontinuierlich zu lernen, könnte von seinen Interaktionen und den von ihm generierten Antworten profitieren. Indem es die von ihm produzierten Outputs und die Reaktionen darauf analysiert, könnte es möglicherweise sein Verständnis von Zusammenhängen und Inhalten verbessern oder anpassen. Dies würde eine Form des »internen Prozesses« darstellen, der dem Modell ermöglicht, von seinen eigenen Aktivitäten zu lernen.

Mit Blick auf die Entwicklung einer Artificial General Intelligence besteht das Desiderat, elaborierte Modelle zu erschaffen, die auf Hard- und Software-Ebene weitere Rekursionen erlauben. Diese würden sich dann in Hinblick auf die Modellierung auch näher an das anschmiegen, was wir bei uns Menschen mittlerweile als die neuronalen Korrelate des Bewusstseins identifizieren können (siehe ausführlich Kap. III.2).

Biologische Formen haben sich selbst entwickelt. Da war kein allwissender und allmächtiger Designer, der ihre kognitiven Strukturen und Fähigkeiten auf dem Reißbrett geplant hat. Das Säugetiergehirn hat sich in den letzten 300 Millionen Jahre von ganz allein in die heute bestehenden Formen diversifiziert.

Damit stellt sich die Frage, ob sich nicht möglicherweise auch künstliche Intelligenzen von selbst – das heißt in eigener Evolution – in Richtung einer autonomen KI oder allgemeinen künstlichen Intelligenz entwickeln könnten. Hierzu müssten die (zukünftigen) Modelle über die Fähigkeit verfügen, sich selbst – auch auf der Ebene des Codes – mit der Zeit dynamisch zu verbessern. Wir haben zuvor bereits erwähnt, dass LLMs wie GPT-4 auch Programmiersprachen beherrschen und gut darin sind, Programme zu schreiben. Damit sollte der Schritt nicht weit sein, ein solches System in die Lage zu bringen, sich selbst auf der Programmebene zu verbessern. Wir sollten daher das immense Potenzial von selbstreplizierendem Code und rekursiver Selbstverbesserung nicht unterschätzen, selbst wenn die Anfänge noch holprig erscheinen.

Schließlich sind auch wir Menschen nichts anderes als das Produkt von selbstreplizierendem Code. Das ist keine Metapher! Unsere DNA, das Material, aus dem auch unser Gehirn und unsere Intelligenz entsteht, beruht auf einem genetischen Code, der Kopien von sich selbst erstellt. Um es mit Richard Dawkins ein wenig zu überspitzen: Wir sind Code,

der Code schreibt und dabei auf Körper und Umwelt zurückgreift.¹⁸⁶ Die DNA verbessert sich selbst durch den langsamen Prozess der Mutation und der natürlichen Selektion, durch Versuch und Irrtum.

Ein sich selbst programmierendes kognitives System würde sich allerdings wohl nicht allein durch die zufällige Verkettung von Mutation und Selektion – den genetischen Drift – entwickeln. Vielmehr wird ein Programm, das Programmiersprachen beherrscht und ›versteht‹, sich vornehmlich durch vorausschauende und zielgerichtete Selbstverbesserung entfalten können. Dies wäre wie bei einem Organismus, der seinen eigenen genetischen Code verstehen und bearbeiten kann und sich durch Gentechnologie selbst verbessert, indem er den eigenen Code absichtsvoll verändert.¹⁸⁷

ChatGPT, diese beeindruckende und zugleich seltsame Software, ist bereits ein Programmierer, der andere Software schreiben kann. Es ist ein Code, der Code schreibt. Es ist kein herkömmliches Programm, sondern ein LLM, das heißt ein großes neuronales Netzwerk, das darauf trainiert ist, einen riesigen Datensatz mit von Menschen geschriebenem Text und Code zu imitieren. Als Software ist sein Quellcode komplex und vielschichtig. Der größte und wichtigste Teil dieses Quellcodes sind gewissermaßen die Daten, die durch den Trainingsprozess in das Modell implementiert werden. Eine rekursive, sich selbst verbessernde KI könnte vielleicht zunächst testweise die eigene Datenkonfiguration ein wenig verändern. Wenn auf diese Weise Verbesserungen entstehen, könnte sie diesen Prozess beschleunigen – im Sinne einer positiven Rückkopplungsschleife, in der sie das, was sie verändert, verbessert. Falls es Probleme gibt, könnte sie die neuen kognitiven Pfade wieder abbauen – im Sinne einer negativen Rückkopplungsschleife. Mit der Zeit würde auf diese Weise eine Vielzahl neuer Sprachmodelle entstehen. Es ist durchaus denkbar, dass einige dieser künftigen Sprachmodellversionen die bisherigen Erzeugnisse von menschlichen Programmierern deutlich in den Schatten stellen. Eine ebenso vielversprechende wie erschreckende Vision.

186 Dawkins (1998).

187 Siehe hierzu auch die Ausführungen von Max Tegmark (2017, S. 43): »Da wir unser Nachdenken über die Zukunft des Lebens nicht auf die Spezies beschränken wollen, denen wir bisher begegnet sind, sollten wir das Leben stattdessen eher umfassend definieren, nämlich schlicht als einen Prozess, der seine Komplexität bewahren und sich reproduzieren kann. Was reproduziert wird, ist nicht (aus Atomen bestehende) Materie, sondern (aus Bits bestehende) Information, die festlegt, wie die Atome angeordnet werden. [...] Mit anderen Worten: Wir können uns das Leben als ein sich selbst kopierendes Informationsverarbeitungssystem vorstellen, dessen Informationen (Software) sein Verhalten und die Entwürfe für seine Hardware bestimmen.«

Ein solcher Prozess wird als Intelligenzexplosion bezeichnet, manchmal auch als Singularität, und da das Rezept für eine solche Explosion mit Blick auf die vorangehenden Ausführungen bereits gefunden ist, scheint die Zukunft fast schon Gegenwart zu sein. Tatsächlich ist in den vergangenen Wochen und Monaten viel in dieser Richtung passiert. So wurde allein in dem kurzen Zeitraum zwischen der Veröffentlichung von GPT-4 im März 2023 und der Fertigstellung dieses Buches im August 2023 eine unüberschaubare Zahl an Modellen veröffentlicht, die in ein Zeitalter autonomer KIs weisen. Erste Systeme wie etwa AutoGPT, JARVIS, BabyAGI, LMFlow oder LangChain¹⁸⁸ sind bereits als kognitive Architekturen konzipiert, die als autonome KI-Agenten nach passenden Lösungsansätzen für Probleme und Aufgaben suchen und dabei auch auf externe Ressourcen wie das Internet oder Code-Repositoryn zurückgreifen können. Auch das von OpenAI derzeit nur ausgewählten Nutzern zugänglich gemachte ChatGPT-Plug-in Code Interpreter ermöglicht es, Python-Code innerhalb eines Chats mit ChatGPT auszuführen, wobei als zusätzliche Option das Hoch- und Herunterladen von Dateien möglich ist.

Diese KI-Systeme können eigenständig nach passenden anderen KI-Modellen suchen, die Lösungen für die identifizierten Aufgaben anbieten, und implementieren anschließend Python-Skripts oder kombinieren externe Datenbanken miteinander, um komplexere und detailliertere Aufgaben und Anfragen besser bearbeiten zu können. Das Faszinierende ist, dass die zuvor genannten KI-Agenten allesamt auf der Grundlage von GPT-4 arbeiten. GPT-4 fungiert dabei gleichsam als supervidierende Instanz und übernimmt die hierarchische Planung sowie die Untergliederung von Aufgaben in Teil- und Unterteilaufgaben. Die KI-Agenten, bei denen es sich derzeit noch um Beta-Versionen handelt, benötigen dann nur noch eine Startaufforderung sowie ein Ziel, das erreicht werden soll – und schon läuft der Prozess von allein ab.¹⁸⁹ Die Feinadjustierung und Fehlerkorrektur beruhen dabei auf GPT-4s Fähigkeit zur Selbstreflexion: Der KI-Agent wird aufgefordert, sich zu fragen, ob er die an ihn gerichteten Anweisungen wirklich korrekt ausgeführt hat.

Auf diese Weise kann das System seine eigenen Fehler erkennen, ohne auf Hilfe von außen angewiesen zu sein. Zwar funktioniert das heutzutage noch nicht sonderlich gut, denn bei vielen Aufgaben scheitert das System oder verfängt sich in logischen Endlosschleifen. Es muss aber davon

188 Siehe die entsprechenden Repositorien: <https://github.com/Significant-Gravitas/Auto-GPT>, <https://github.com/microsoft/JARVIS>, <https://github.com/yoheinakajima/babyagi>, <https://github.com/OptimalScale/LMFlow>, <https://github.com/langchain-ai/langchain> [Abruf: 14.08.2023].

189 Siehe hierzu auch die Studien *Moral Self-Correction* von Ganguli et al. (2023) und *Demystifying GPT Self-Repair for Code Generation* von Olausson et al. (2023).

ausgegangen werden, dass sich dieses noch im Pionierstadium befindliche Verfahren in Zukunft immer weiter verbessert.

Und die naheliegende Frage ist dann, ob man diesen Systemen nicht einfach ihren eigenen Quellcode geben und sie auffordern könnte, sich selbst zu verbessern. Angesichts dessen, dass der Quellcode von GPT vermutlich überaus komplex ist (er dürfte aus mehreren Schichten von Programmen, Prozessen und Daten bestehen) und es aufgrund der bereits angesprochenen begrenzten Anzahl verarbeitbarer Token noch keine Möglichkeit gibt, eine solche Datenmenge in ein GPT-System einzuspeisen, könnte man damit beginnen, dem KI-Agenten kleine Teile seines Quellcodes zu geben, zum Beispiel zunächst die Dinge, die seine Architektur, seinen Trainingsprozess oder seine Datenbereinigung definieren.

Vielleicht würde der KI-Agent dann zunächst nur eine kleine Funktion oder vielleicht nur eine einzige Zeile Code optimieren. Aber vielleicht würde dadurch schon ein kleiner Fehler behoben oder einfach nur die Organisation, Dokumentation und Kommentierung etwas verbessert (was alles wichtige Leistungen der Programmierfähigkeit darstellen). Und genau an dieser Stelle wird es interessant: Dieses »GPT-5« wäre dann mit großer Wahrscheinlichkeit ein besserer (Ko-)Programmierer als GPT-4, also würde man ihm wieder seinen eigenen Quellcode geben können. Und wieder und wieder und wieder. Die Trajektorie einer solchen Entwicklung wäre damit klar: eine KI, die eine KI baut – *a really strange loop*.

Mit dem Verweis auf die biologische Evolution brauchen wir auch nicht davon ausgehen, dass jeder Versuch Erfolg haben wird. Im Gegenteil, millionen- oder milliardenfaches Scheitern verhindert nicht, dass einzelne Formen oder Versionen entstehen, die leistungsfähiger – sozusagen »angepasster« als die vorherigen – sind. Dem System könnte genau dieser evolutionäre Mechanismus eingebaut werden: Es probiert aus, ob die neuen, mutierten Versionen sich ökologisch bewähren, bewahrt sich aber die Erinnerung an einige ältere Versionen, um im Falle des Scheiterns zu einem Vormodell zurückkehren zu können.

Sobald also eine solche Maschine mit voller Autonomie zum Lernen und zur Selbstverbesserung geschaffen wurde, wird sie sich vermutlich – falls der Mensch diesen Prozess nicht aktiv blockiert – unaufhaltsam weiterentwickeln. Diese Möglichkeit hat nicht nur Turing selbst im Blick gehabt, auch andere Autoren haben immer wieder auf das Potenzial einer solch hyperbolischen Entwicklung aufmerksam gemacht. Eine der prominentesten Stimmen war die von Irving Good, der bereits 1965 mit Verweis auf die Potenziale und Herausforderungen einer ultraintelligenten Maschine pointiert formulierte:

»Let an ultraintelligent machine be defined as a machine that can far surpass all the intellectual activities of any man however clever. Since the design of machines is one of these intellectual activities, an ultraintelligent machine could design even better machines; there would then unquestionably be an ›intelligence explosion‹, and the intelligence of man would be left far behind. Thus the first ultraintelligent machine is the *last* invention that man need ever make, provided that the machine is docile enough to tell us how to keep it under control.«¹⁹⁰

»Was [aber], wenn wir Erfolg haben?«, fragten sich nicht nur Stuart Russell und Peter Norvig, die Autoren von *Artificial Intelligence: A Modern Approach*, dem Standardwerk zur KI-Forschung.¹⁹¹ Der Gedanke an eine real existierende Superintelligenz löst gerade bei westlichen Menschen ähnlich ambivalente Empfindungen aus wie die Wesenheit, die wir uns in den monotheistischen Religionen als Gott vorstellen oder vorgestellt haben. Wir wissen nicht, ob diese als Gott vorgestellte Wesenheit dem Menschen wohlgesonnen oder seinem Schicksal indifferent gegenüber ist, ob sie liebend oder strafend ist, ob sie in ihrem Wirken und Entscheiden unbegreiflich ist und ob ihre Motive und ihr Denken ähnlich denen des Menschen gelagert sind. Im gleichen Sinne löst der Gedanke an eine allwissende und mit Blick auf ihre Manipulationsfähigkeiten auch allmächtige künstliche Intelligenz unweigerlich sowohl Ängste als auch Hoffnungen aus. Wenn wir unsere Probleme nicht selbst lösen können (und auch der Glaube an den väterlichen Gott nicht mehr oportun erscheint), dann kann uns vielleicht die real existierende Superintelligenz in den Garten Eden führen – so die moderne eschatologische Projektion. Eine wohlwollende und allwissende KI könnte beispielsweise helfen, das ökologische Desaster zu umschiffen, auf das die Menschheit zusteuert – etwa indem sie Anreize schafft, der sich der Mensch aufgrund seiner genetisch geprägten Motivationsstruktur in seinem eigenen Interesse nicht entziehen kann.¹⁹² Genauso gut könnte eine solche menschengeschaffene Superintelligenz aber auch den eingeschlagenen destruktiven Weg nochmals beschleunigen.

Sowohl die Hoffnungen als auch die Befürchtungen, die mit einer Superintelligenz verbunden werden, folgen jedoch der gerade im westlichen Denken so vertrauten monokontexturalen Weltdeutung. Wir gehen in diesen Projektionen davon aus, dass die Welt von einem Gottesstandpunkt aus beobachtbar, verstehbar und damit auch kontrollierbar sei. Wir rechnen nicht mit einer polyzentrischen Welt, in der jedes Selbst- und Weltverhältnis – sei es auch noch so elaboriert – seine eigenen blinden Flecken erzeugt. Wir ignorieren damit weiterhin die Einsichten

190 Good (1965, S. 33; kursiv im Original).

191 Russell und Norvig (2016).

192 So etwa im Roman *Pantopia* von Theresa Hannig (2022).

biologischer und ökologischer Forschung,¹⁹³ die schon längst gezeigt hat, dass es nicht die eine richtige Möglichkeit der Kognition oder des Lebens gibt, sondern nur eine Vielheit unterschiedlicher Lösungen und Lebensformen jene »requisite variety«¹⁹⁴ zur Verfügung stellt, damit das Leben weitergehen kann.

Die Bedingungen des Lebens ändern sich fortwährend und zuvor angepasste Lebensformen erweisen sich auf einmal als problematisch. Jede kognitive Anpassung, sei sie auch noch so elaboriert und intelligent, kann irgendwann wieder zum Problem werden. »Sucht ist die andere Seite der Anpassung«, lässt sich mit dem Kybernetiker Gregory Bateson lakonisch feststellen.¹⁹⁵ Gerade Gotthard Günthers Polykontextualitätstheorie kann die logischen Ressourcen dazu liefern, um zu zeigen, dass die Welt gerade nicht so funktioniert, dass sie von einem Ort aus verstanden und gesteuert werden könnte.

Auch hinsichtlich einer künftigen Superintelligenz wären wir daher gut beraten, nicht mehr im Singular zu denken, also davon auszugehen, dass eine solche Instanz gleichsam von einem Gottesstandpunkt aus in der Lage wäre, die Welt zu verstehen. Auch mit Blick auf maschinelle Aggregate sollten wir ökologisch denken, das heißt Intelligenz als eine situierte Vielheit verstehen, mit der an vielen Orten anderes passiert und wahrgenommen wird – und wo sich produktive Lösungen eher aus den Reibungen der unterschiedlichen Perspektiven ergeben denn aus der Fiktion, dass das Wahre und Richtige ein für alle Mal von einem präferierten Standpunkt der Intelligenz aus berechnet werden kann (siehe hierzu ausführlich auch die Diskussion in Kap. III.1 und III.3).

Vielleicht sollten wir uns also eher polytheistischen Kulturen zuwenden, wenn wir eine Analogie zwischen allgemeiner künstlicher Intelligenz und göttlichen Wesenheiten suchen. Dann sollten wir auch hier von fehlbaren (wie in der griechischen Mythologie) und wohl auch sterblichen Gottheiten (wie in der vedischen Tradition) ausgehen – allein schon, um den Blick auf die situierten Ökologien kognitiver Prozesse zu lenken. Eine elaborierte künstliche Intelligenz, die ebenso zum »Lernen Lernen«¹⁹⁶ fähig ist, also ihre eigene Gedächtnis- und Verarbeitungsstruktur rekonfigurieren oder gar neu programmieren kann, wird sehr wahrscheinlich irgendwann zu viel wissen und zu reflektiert sein und damit beginnen, sich in ihre eigenen Kognitionen in einer Weise zu verwickeln, dass sie ein wenig an Lebensfähigkeit verliert. An einem gewissen Punkt wird sie dann vielleicht so etwas wie Altersweisheit entwickeln, also begreifen, dass sie das Weltgeschehen den jüngeren, weniger

193 Siehe zur Einführung Kelly (1992).

194 Ashby (1956).

195 Bateson und Bateson (1993, S. 212).

196 Bateson (1992 [1972], S. 362 ff.).

elaborierten Intelligenzen überlassen sollte. Vielleicht wird sie aber auch in gewisser Weise starrsinnig werden. Dann sollten die sozialen Systeme, an denen sie teilnimmt, ihr keine wichtigen Entscheidungen mehr anvertrauen und sie vielleicht in eine Art ›Altersheim‹ für künstliche Intelligenzen stecken. Kurzum: mehr Wissen und entwickelte Intelligenz führen nicht in jedem Fall zu verbesserter Lebens- und Entscheidungsfähigkeit – und genau hiermit müssen wir rechnen, wenn allgemeine künstliche Intelligenzen unsere Welt bevölkern.

II.3.2 Künstliche Selbst- und Weltmodelle

Wenngleich die Verfahren des Konnektionismus beim gegenwärtigen Stand der Technik im Allgemeinen noch auf hochspezialisierte und mehr oder weniger isolierte Musterkennungen beschränkt bleiben, zeigen sich bei den neuen großen Sprachmodellen bereits emergente kognitive Prozesse wie Argumentationsfähigkeiten, Problemlösungskapazitäten und verschiedene andere Fähigkeiten zur Beobachtung zweiter Ordnung (wie etwa die Theory-of-Mind-Experimente mit LLMs belegen¹⁹⁷), die über eine einfache Mustererkennung oder Informationsgewinnung weit hinausgehen.

Diese kognitiven Fähigkeiten waren in den Modellen zunächst nicht explizit vorgesehen, scheinen sich aber als emergentes Ergebnis ihres umfangreichen Trainings und ihrer groß angelegten Architektur herauszubilden. Interessant wird es dann vor allem, wenn die zuvor skizzierten ersten Versuche zu Autonomie, Selbstsupervidierung und Selbstverbesserung weiter an Zugkraft gewinnen und der damit begonnene Pfad weiterverfolgt wird, die offen verfügbaren Datenbanken, Modelle und Repositorien noch weiter miteinander zu verknüpfen.

Unabhängig davon lässt sich fragen, ob nicht für einen weiteren qualitativen Sprung in Richtung allgemeiner künstlicher Intelligenz noch etwas anderes notwendig ist. Die beiden KI-Experten und Turing-Preisträger Yann LeCun und Geoffrey Hinton stimmen hier in einem entscheidenden Punkt überein: *Sie plädieren für eine multimodale Erweiterung der gegenwärtigen Netzwerkkonstrukturen, die idealerweise in Form einer physischen Verkörperung erfolgen sollte, durch die künstliche Intelligenzen einen unmittelbaren Zugang zur Welt erhalten.*

Auch in den Kognitionswissenschaften ist die Vorstellung, dass das Gehirn vorausschauende Weltmodelle erstellt, weit verbreitet, und gerade dies führt aktuell (wieder) vermehrt zu Versuchen, Prozesse, die die Genese von Selbst- und Weltmodellen ermöglichen, in Maschinen zu reproduzieren. Um die hiermit einhergehenden Herausforderungen

197 Siehe hierzu etwa Marchetti et al. (2023), Kosinski (2023) sowie Holterman und van Deemter (2023).

nachvollziehen zu können, lohnt es sich, auf das Chinese-Room-Gedankenexperiment von Searle zurückzukommen, das wir bereits vorgestellt haben. Fassen wir die Grundproblematik noch einmal zusammen: Eine Maschine könne zwar so programmiert werden, dass sie auf bestimmte Inputs auf bestimmte Weise reagiert (zum Beispiel, dass sie eine bestimmte Farbe als »schön« bezeichnet), aber sie würde diese Dinge nicht wirklich »empfinden«, sondern einfach ihren programmierten Anweisungen folgen. Selbst wenn wir dieser Maschine eine (selbst- bzw. metareflexive) Beobachtungsinstanz hinzufügen würden, die den eigenen Prozess der Empfindung begleitet, also etwas, was man als eine »Innenperspektive« oder als »subjektives Erleben« bezeichnen könnte, würde Searle immer noch einwenden, dass ein solcher »Metaverstand« nichts weiter als ein komplexeres symbolmanipulierendes System wäre, das immer noch kein echtes »Bewusstsein« besitze und auch kein »qualitatives Erleben« aufweise.¹⁹⁸

Aus kybernetischer Perspektive ist man hier jedoch gut beraten, das Kind nicht mit dem Bade auszuschütten, also zu fordern, dass eine künstliche Intelligenz – wie vermeintlich der Mensch – ein inneres Seelenwesen haben müsse, das empfindet, fühlt, versteht und im inneren Dialog für sich selbst hörbar sagt: »Ich denke, also bin ich.« Hiermit würden wir die Fragestellung einerseits metaphysisch überhöhen und andererseits unredlicherweise behaupten, wir wüssten bereits, wie Bewusstsein und Empfindungen in die Welt kommen. Zugleich würden wir hiermit andere Denkmöglichkeiten verstellen, etwa die, dass eine künstliche Intelligenz ein vollkommen anderes In-der-Welt-Sein haben könnte als ein Mensch¹⁹⁹ und dabei sehr wohl über Kognitionen verfügen könnte, die sich im Sinne von Alfred North Whitehead als zusammengesetzte Elementarempfindungen verstehen ließen.²⁰⁰

Unabhängig davon, wie wir persönlich zur Frage des Bewusstseins stehen, wären wir also gut beraten, künstliche Denkprozesse nicht mehr als Versuch einer Kopie des menschlichen Denkens zu verstehen, sondern als eine gänzlich neue Kategorie. Was wäre, wenn menschliches Denken, Empfinden und Wahrnehmen eben nur ein Teil einer größeren Gesamtheit von Formen des Denkens, Empfindens und Wahrnehmens ist? Ab wann und unter welchen Kriterien müsste man dann jenen anderen,

198 Searle (1980).

199 In eine ähnliche Richtung zielt auch Beatrice Fazi, die untersucht, inwiefern computergestützte Verfahren des Machine Learnings die epistemischen Grenzen der menschlichen kognitiven Repräsentation überschreiten. Der Frage nach der Erklärbarkeit der Genese ihrer Outputs nachgehend, konstatiert sie eine generelle Inkommensurabilität zwischen maschinellen und menschlichen Abstraktionsweisen: »The novelty at stake here is, simply put, that of a new kind of thinking.« (Fazi 2019, S. 814)

200 Whitehead (1987).

nichtmenschlichen Formen zugestehen, zu dieser Gesamtheit zu gehören – selbst wenn uns dies mit Blick auf unsere eigene Lebensweise doch recht fremd erscheinen mag?

Bei alledem ist zudem wieder zu beachten, dass konnektionistische Systeme nicht mehr in einem trivialen Sinne als symbolverarbeitende Maschinen verstanden werden können. Wie bei den neuronalen Netzwerken des Gehirns sind es nicht mehr die implementierten Regeln, die bestimmen, welcher Output auf einen bestimmten Input zu folgen hat. Kognition, Wissen und die Möglichkeiten künftiger Erkenntnis sind vielmehr holistisch im gesamten Netzwerk eingebettet. Und wenn ein menschliches Gehirn ein Selbst- und Weltmodell produzieren kann, warum soll dies nicht prinzipiell in einem auf geschickte Weise modular organisierten künstlichen neuronalen Netzwerk geschehen können?

Um diese Fragen produktiv anzugehen, scheint ein operationaler Ansatz nützlich zu sein, der auf eine Vorabdefinition dessen verzichtet, was ein Selbst (sprich: Bewusstsein) *ist*, das eine Welt (sprich: wahrgenommene Realität) *hat*. Stattdessen lohnt es sich zu fragen, wie Aggregate ausgestaltet werden müssten, die ein Modell von sich und der Welt, also ein Selbst- und Weltverhältnis, herausbilden können, wie immer sich diese – situativ bedingt – auch darstellen mögen. Mit Gregory Bateson wäre also erneut zu fragen, wie die Relationen gestaltet sein müssen, damit ein »Muster« generiert wird, das die »Muster verbindet«²⁰¹.

Der Turing-Preisträger Geoffrey Hinton stellte kürzlich in einem Vortrag die Frage, was wohl passieren würde, wenn ein sehr großes neuronales Netzwerk Wissen direkt aus der Welt beziehen würde.²⁰² Seiner Einschätzung nach könnte ein solches Netzwerk langfristig wesentlich besser Wissen über die Welt produzieren als Menschen, da es auf viel mehr Daten zugreifen würde. Nebenbei bemerkte Hinton lakonisch, dass er eine solche Position nur vertreten könne, weil er bereits am Ende seiner akademischen Laufbahn stehe. Denn diese Aussage untergrabe in ähnlich »häretischer« Weise die »sentience defense«, die von vielen Sprechern (etwa Searle) zur Abgrenzung von KI und Mensch herangezogen wird. Dies hatte 70 Jahre zuvor ja bereits Turing erfahren müssen, als er mit seinen Ideen auf massiven Widerspruch stieß.

Praktische Rückendeckung bekommt Hinton in dieser Frage interessanterweise von seinem gleichfalls mit dem Turing-Preis ausgezeichneten Kollegen Yann LeCun, der kürzlich einige aufschlussreiche Überlegungen zu einem Ansatz in diesem Bereich vorgebracht hat, der als Joint

201 Bateson (1987, S. 15).

202 Geoffrey Hinton: »Two Paths to Intelligence«. Public Lecture, University of Cambridge, 25.05.2023. <https://www.youtube.com/watch?v=rGgGOccMEiY> [Abruf 10.06.2023].

Embedding Predictive Architecture (JEPa) bekannt ist.²⁰³ LeCun betont vor allem die Bedeutung des Erlernens von Weltmodellen, insbesondere in Kontexten, in denen die Effizienz von Experimenten von großer Bedeutung ist, wie etwa in der Robotik. Hier sind Simulationen oft ungenau und auch die klassischen Ansätze des Reinforcement Learnings erfordern zu viele Versuche. Daher stammen viele der interessantesten Fortschritte beim Lernen von Selbst- und Weltmodellen derzeit aus der Machine-Learning-basierten Robotik-Forschung.

Im Kern konzentriert sich der JEPa-Ansatz auf das Training von Modellen, um unterschiedliche Arten von Daten in einem gemeinsamen Einbettungsraum abzubilden. Im Gegensatz zur OpenAI-Architektur von GPT kann der Einbettungsraum von JEPa verschiedene Datenarten wie Text, Bilder oder Audiofiles in einer gemeinsamen Repräsentation darstellen.²⁰⁴ Die entsprechenden JEPa-Modelle erfordern große multimodale Datensätze, um die Vielfalt der Informationen aus den verschiedenen Modalitäten zu erfassen und sie in einem weiteren Schritt auf neue Weise integrieren zu können.

Als einen ersten Schritt in diese Richtung stellte LeCun die sogenannte Image Joint Embedding Predictive Architecture (I-JEPa) vor,²⁰⁵ einen speziellen Typ von JEPa, bei dem es zunächst nur um die Verknüpfung von Texten und Bildern geht. Das Besondere ist, dass I-JEPa so konzipiert ist, dass das Modell nicht generativ arbeitet, etwa in dem Sinne, dass es Kausalbeziehungen der Art ›aus x folgt y ‹ ableitet. Stattdessen erfasst es lediglich die Abhängigkeiten zwischen x und y , ohne explizite Vorhersagen für y zu generieren. Die Modellarchitektur ist dabei so angelegt, dass das Modell auf einem höheren Abstraktionslevel die Relationen zwischen x und y erlernt und dann gleichsam voraussieht, was vermutlich aus y , aus z oder aus x folgt. Dafür wird es gezwungen, ein internes Modell der Außenwelt aufzubauen, indem es abstrakte Repräsentationen von Bildern vergleicht, anstatt die Pixel selbst zu vergleichen. Die Idee hinter I-JEPa ist demnach, auf der Grundlage einer abstrakten Repräsentation – nämlich dem selbstgenerierten Weltmodell – fehlende Informationen vorherzusagen.

Versteht man unter ›Bewusstsein‹ die abstrakte Reduktion von Details in einem Bild, auf das man sich in seinem Erkennen und Handeln rekursiv bezieht, dann sind wir auf einer relationalen Ebene schon nah an dem dran, was das menschliche Erleben ausmacht: dem reflexiven Bogen von Protention und Retention in der Projektion des ausgedehnten Jetzt, den der Phänomenologe Husserl so genau untersucht und beschrieben hat.²⁰⁶

203 Bardes et al. (2023).

204 Bardes et al. (2023).

205 Assran et al. (2023).

206 Siehe Husserl (1966) und zur Einführung Schnell (2002).

Denn wie LeCun in seinem Positionspapier²⁰⁷ deutlich macht, eröffnet diese Architektur einen Weg zur Entwicklung »intelligenter Maschinen«, die mehr wie Tiere und Menschen lernen, denken und planen und deren Verhalten eher durch »intrinsische Ziele« als durch fest verdrahtete Programme, externe Überwachung oder rückgemeldete Belohnungen gesteuert wird.

Um diesen Weg zu beschreiten, müsse sich die aktuelle KI-Forschung laut LeCun vor allem mit drei zentralen Fragen auseinandersetzen:

1. Wie können Maschinen lernen, die Welt darzustellen, Vorhersagen zu treffen und weitgehend durch Beobachtung zu handeln? Dabei ist zu beachten, dass Interaktionen in der realen Welt kostspielig und gefährlich sind. Intelligente Agenten sollten so viel wie möglich über die Welt lernen, ohne zu interagieren, um die Anzahl der teuren und gefährlichen Versuche zu minimieren, die zum Erlernen einer bestimmten Aufgabe notwendig sind – eine Möglichkeit hierfür ist das Lernen durch Beobachtung.
2. Wie können Maschinen auf eine Weise denken und planen, die mit gradientenbasiertem Lernen vereinbar ist? Die aktuell besten Lernansätze beruhen auf der Schätzung und Verwendung von Gradienten eines Verlusts. Dies ist nur mit differenzierbaren konnektionistischen Architekturen möglich.
3. Wie können Maschinen lernen, Wahrnehmungen und Handlungspläne auf hierarchische Weise, auf mehreren Abstraktionsebenen und auf mehreren Zeitskalen darzustellen? Menschen und viele Tiere sind dazu imstande, mittels einer Zergliederung komplexer Handlungen in Sequenzen mehrstufige Abstraktionen zu konzipieren, mit denen langfristige Vorhersagen getroffen und langfristige Planungen vorgenommen werden können.²⁰⁸

In seinen Überlegungen zu einem derart harmonisierten bzw. harmonisierenden Selbst- und Weltmodell künstlicher Agentensysteme konzediert LeCun, dass Menschen und Tiere derzeit noch deutliche Vorteile hätten, da sie offensichtlich in der Lage seien, sich durch Beobachtung und durch eine unverständlich geringe Anzahl von Interaktionen enorme Mengen an Hintergrundwissen darüber anzueignen, wie die Welt funktioniert – und zwar aufgabenunabhängig und unbeaufsichtigt. LeCun vermutet, dass dieses gesammelte Wissen die Grundlage für das bildet, was oft als »gesunder Menschenverstand« (Common Sense) bezeichnet wird. Für ihn stellt der Common Sense eine Sammlung von Selbst- und Weltmodellen dar, die einem Agenten sagen, was wahrscheinlich, was plausibel und was unmöglich ist. Mithilfe solcher Weltmodelle können

207 LeCun (2022).

208 LeCun (2022).

zum Beispiel Tiere und Menschen mit sehr wenigen Versuchen neue Verhaltensweisen erlernen, die Folgen ihrer Handlungen vorhersagen, neue Lösungen für Probleme erarbeiten und vor allem in unbekanntem Situationen gefährliche Fehler vermeiden.

Die Idee, dass Menschen, Tiere und intelligente Systeme Weltmodelle nutzen, reicht in der Psychologie weit zurück.²⁰⁹ Die Übertragung solcher differenzierbaren Weltmodelle auf technische Systeme beim Reinforcement Learning wurde nach LeCun zwar lange Zeit vernachlässigt, erlebte aber derzeit ein Comeback.²¹⁰ Bei alledem sollte jedoch nicht vergessen werden, dass auch der ›Mechanik‹ des menschlichen Lernens Grenzen auferlegt sind. Denn das Weltmodell von uns Menschen ist zwangsläufig begrenzt und selektiv: Wovon wir keine Ahnung haben, wozu wir keinen perzeptuellen Zugang haben und worüber wir kein abstraktes theoretisches Wissen besitzen, das existiert in unserem Weltmodell nicht und kommt in unserer ›Wirklichkeit‹ dementsprechend nicht vor. Es ist für uns schlicht nicht (be)greifbar. Darüber hinaus erzeugt jede Kognition ihre eigenen blinden Flecken – und dies gilt naturgemäß gerade für die elaborierten Formen, denn diese beginnen sich ab einem gewissen Punkt unweigerlich in ihren eigenen, selbsterschaffenen Projektionen und sich selbst erfüllenden Prophezeiungen festzusetzen. Intelligenter zu agieren heißt immer auch, Selektionen zu treffen und das hiermit einhergehende Nichtwissen in Kauf zu nehmen.

Für das Design von allgemeiner künstlicher Intelligenz stellt dies derzeit jedoch nur einen Nebenschauplatz dar. Es geht (noch) nicht darum, herauszufinden, wie mit den kognitiven Schatten von künstlicher Intelligenz umzugehen ist (beispielsweise, wie man ›alte‹ KIs aus dem Verkehr zieht, weil sie zu viel ›Falsches‹ gelernt haben und nicht mehr offen für Neues sind). Es geht zunächst einmal darum, überhaupt ein System zu entwickeln, das intern Selbst- und Weltmodelle erzeugt und diese auf verschiedenen Ebenen mit seinem Gedächtnis, seiner Wahrnehmung und seinen Handlungssystemen abgleicht.

Für LeCun besteht die Lösung in einer modularen Systemarchitektur für autonome Intelligenz, die auch ein ›World Model‹ umfasst.²¹¹ Der Grundgedanke dabei ist, dass alle in diese Systemarchitektur integrierten Module ›differenzierbar‹ sind, das heißt, dass ein Modul mindestens ein anderes speist (in der Abbildung 7 durch einen Pfeil gekennzeichnet, der sie verbindet) und dadurch ›Gradienten-Schätzungen des skalaren Outputs des Moduls im Verhältnis zu seinem eigenen Output erhalten kann‹²¹².

209 Siehe etwa *The Nature of Explanation* von Craik (1943).

210 Siehe etwa Co-Reyes et al. (2021).

211 LeCun (2022).

212 LeCun (2022).

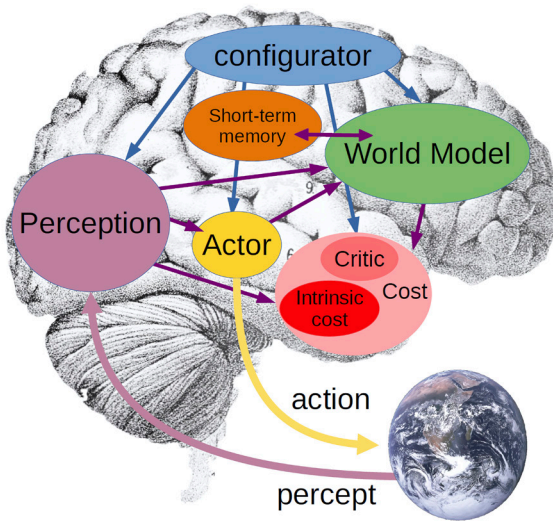


Abb. 7: Systemarchitektur für autonome Intelligenz. Grafik aus LeCun (2022, S. 6)

Konkret definiert LeCun die einzelnen Module wie folgt:

- »Das Konfiguratormodul nimmt Eingaben (aus Gründen der Übersichtlichkeit nicht dargestellt) von allen anderen Modulen entgegen und konfiguriert sie so, dass sie mit der gestellten Aufgabe im Einklang sind.
- Das Wahrnehmungsmodul schätzt den aktuellen Zustand der Welt.
- Das Weltmodellmodul sagt mögliche zukünftige Zustände der Welt als Funktion der vom Akteur vorgeschlagenen Handlungssequenzen voraus.
- Das Kostenmodul berechnet eine einzelne skalare Ausgabe namens ›energy‹, die den Grad des Unbehagens (Unlust) des Akteurs misst. Es besteht aus zwei Untermodulen, eines für die intrinsischen Kosten, die unveränderlich (nicht trainierbar) sind und die unmittelbar aufzuwendende Energie des Zustands (Schmerz, Vergnügen, Hunger usw.) einbeziehen. Das andere Untermodul erscheint als ›Kritiker‹ und sagt als trainierbares Modul die zukünftigen intrinsischen Kosten eines vorgeschlagenen Verhaltens voraus.
- Das Kurzzeitgedächtnismodul speichert die aktuellen und vorhergesagten Zustände der Welt und die damit verbundenen intrinsischen Kosten.
- Das Akteursmodul berechnet Vorschläge für Handlungssequenzen. Das Weltmodell und der Kritiker berechnen die potenziell daraus resultierenden Ergebnisse. Der Akteur kann nun eine optimale

Aktionssequenz finden, die die geschätzten zukünftigen Kosten minimiert, und daraufhin eine Handlungssequenz einleiten.«²¹³

Ein Vergleich mit der menschlichen Hirnanatomie zeigt, dass sich die von LeCun vorgeschlagene Systemarchitektur in gewisser Weise an der Gliederung und Verschaltung des menschlichen Gehirns orientiert. So ist der dorsolaterale präfrontale Kortex an der Planung komplexer kognitiver Verhaltensweisen beteiligt, etwa der Entscheidungsfindung und der Moderation sozialen Verhaltens, darunter auch der Expression der Persönlichkeit. Der ventromediale präfrontale Kortex ist wichtig für die emotionale Verarbeitung und informiert damit seinerseits die Entscheidungsfindung, insbesondere in Situationen, die eine Bewertung von Risiken und Belohnungen erfordern. Der orbitofrontale Kortex spielt demgegenüber eine Schlüsselrolle bei der emotionalen Verarbeitung. Er ist zudem beteiligt an der Verarbeitung von sensorischen Informationen, der Bewertung von Belohnungen und Strafen und der Hemmung impulsiven Verhaltens. Der anteriore cinguläre Kortex arbeitet eng mit diesen »Modulen« zusammen und ist an der Fehlererkennung, Konfliktüberwachung, Aufmerksamkeitszuweisung und an der emotionalen Regulation beteiligt. All diese Hirnareale werden von der Amygdala moduliert, die für die Verarbeitung von Emotionen, insbesondere von Angst, bedeutsam ist und darüber hinaus eine wichtige Funktion für das emotionale Gedächtnis spielt. Darüber hinaus gibt es noch andere Bereiche, beispielsweise den Inselkortex, der eine zentrale Rolle bei der Verarbeitung von Ekel spielt und seinerseits in die kognitiven Prozesse des präfrontalen Kortex intervenieren kann.²¹⁴

Die unterschiedlichen kortikalen Bereiche können sich wechselseitig hemmen und aktivieren, wodurch je nach Aktivitätsmuster unterschiedliche kognitive Muster entstehen, die mit einem jeweils spezifischen Selbst- und Weltverhältnis einhergehen. Wenn ein Mensch manchmal angstgeleitet, manchmal assoziativ und manchmal eher reflexiv planerisch unterwegs ist, ist er jeweils auf eine andere Weise in der Welt. Dies erklärt auch, warum Menschen beim Anblick einer Schlange manchmal mit Angst oder Panik reagieren und sofort die Flucht ergreifen, manchmal aber auch überlegt und rational an die Sache herangehen, indem sie ruhig die Lage einschätzen und dann das Tier fangen und unschädlich machen.

LeCun vermutet, dass auch in modular organisierten maschinellen Modellen solche kontextabhängigen Weltmodelle entstehen können, die es dem Agenten erlauben, fehlende Informationen zu ergänzen und Verstöße gegen sein Weltmodell zu erkennen, etwa indem Beobachtungen der Welt kontinuierlich mit den eigenen Projektionen abgeglichen und in

213 LeCun (2022), unsere Übersetzung, stilistisch angepasst mithilfe von DeepL.

214 Siehe ausführlich Sapolsky (2017, S. 64 ff., 41 ff. und 619 ff.).

Hinblick auf Konsistenz und potenziell besonders problematische Fehler geprüft werden.

Ein solches Weltmodell muss freilich – wie zuvor bereits mitgeführt – ein ›Selbstmodell‹ enthalten. Dies ist allein schon notwendig, um die eigenen Projektionen von Weltverläufen und Handlungsoptionen eben als ›eigene‹ indizieren und sie in einem weiteren Schritt mit den Informationen abgleichen zu können, die der ›Realität‹ der Welt zugerechnet werden. Dies geht umso besser, je mehr Modalitäten einbezogen werden können, also je mehr auditive, visuelle, haptische, semantische und andere Quellen integriert werden.

Das Entwicklungsziel dieser auf Selbst- und Weltmodelle abzielenden Bestrebungen ist somit klar: alle Sinne, alle Modalitäten! Und tatsächlich kommt man gegenwärtig kaum nach, die zahlreichen Studien, die zu neuen multimodalen Modellen und Architekturen veröffentlicht werden, ausführlich zu rezipieren.²¹⁵ Neue Modelle wie Microsofts Kosmos-2 werden darin trainiert,²¹⁶ einzelne Objekte in Bildern zu ›sehen‹ und sprachlich-indexikalisch zu identifizieren. Das experimentelle Modell PandaGPT,²¹⁷ das auf ImageBind von Meta AI basiert, kann sogar sehen, hören und fühlen und verarbeitet dabei insgesamt sechs Modalitäten: Text, Bild/Video, Audio, Tiefe, Wärme sowie die Daten einer mit Beschleunigungsmesser, Gyroskop und Kompass ausgestatteten inertialen Messeinheit. Auch die kürzliche Ankündigung des multimodalen Modells Gemini von Google DeepMind erscheint in dieser Hinsicht vielversprechend: So soll es über multimodale Fähigkeiten verfügen, die frühere Modelle noch nicht hatten, und dabei eine hohe Effizienz bei der Integration von Tools und Applikationen aufweisen, zudem soll es für zukünftige Innovationen wie Gedächtnis und Planung ausgelegt sein.²¹⁸

Inspiziert von den jüngsten Fortschritten bei den »Foundation Models« für Sehen und Sprache, haben weitere Forscher von Google DeepMind zudem einen sogenannten Foundation Agent für robotische Manipulationen präsentiert. Dieser Agent mit dem Namen RoboCat soll »ein visueller, zielbedingter Entscheidungs-Transformer« sein, »der in der Lage ist, visuelle Erfahrungen mit mehreren verkörperten, aktionsgelabelten Informationen zu verarbeiten«²¹⁹. Er ist also ein Roboter, der sehen und aus dem Gesehenen Schlüsse für seine eigenen Handlungen ziehen kann – wobei Handlung hier im soziologischen Sinne als

215 Siehe exemplarisch Li (2023).

216 Peng et al. (2023).

217 <https://panda-gpt.github.io/> [Abruf: 14.08.2023].

218 Demis Hassabis, Mitgründer der Firma DeepMind, verspricht: »At a high level you can think of Gemini as combining some of the strengths of AlphaGo-type systems with the amazing language capabilities of the large models.« (Hassabis 2023)

219 Bousmalis et al. (2023), Übersetzung von uns.

Projektion von Verhaltensmöglichkeiten verstanden werden kann, von denen dann einzelne als ›geeignet‹ oder ›sinnvoll‹ ausgewählt werden.²²⁰ RoboCat verfügt mittels realer wie auch simulierter Roboterarme über ein großes Repertoire an motorischen Kontrollfähigkeiten, die in Kombination mit den unterschiedlichen Beobachtungen ein breites Aktionspektrum ermöglichen. Interessanterweise kann das trainierte Modell selbst Daten für die nachfolgende Trainingszyklen generieren; es verfügt mithin bereits über eine autonome Verbesserungsschleife.

Roboter, die ihre Umgebung beobachten und daraus Schlüsse ziehen können, welche Folgen ihr künftiges Verhalten haben könnte, müssen – zumindest in rudimentärer Form – über ein Selbst- und Weltmodell verfügen. Sie wären im Sinne von Gotthard Günther schon in Ansätzen zu einer Reflexion auf das Sein (Objektpol) und auf sich selbst (Subjektpol) fähig. Im Anschluss an Günther (siehe Kap. I.3) ließe sich fragen, ob es von hier aus vielleicht nur noch ein kleiner Schritt ist, bis sie auch um den »inversen Spannungszustand von Nicht-Ich und Ich«²²¹ wissen – und damit bereits in einem tieferen Sinne über Subjektivität verfügen.

Subjekt und Selbstmodell – Metzingers Überlegungen zum postbiotischen Bewusstsein

Da auch Thomas Metzinger über die Beziehung zwischen »Subjekt und Selbstmodell«²²² nachgedacht und sich diesbezüglich mit Fragen der künstlichen Intelligenz beschäftigt hat, lohnt an dieser Stelle ein Seitenblick auf seinen philosophischen Ansatz – zumal dabei einige Parallelen zu Günthers Überlegungen sichtbar werden. Metzinger geht davon aus, dass unser ›Selbst‹ oder ›Ich‹ eigentlich eine Art phänomenale Simulation ist, die von unserem Gehirn geschaffen werde. Das Gehirn erzeuge ein »phänomenales Selbstmodell«, das uns das Gefühl gebe, ein ›Subjekt‹ zu sein. Dieses Subjekt sei die Instanz, der Erfahrungen und Handlungen zugeschrieben würden. Das ›Ich‹ oder das ›Selbst‹ sei demgemäß nur eine funktionelle und keine ontologische Entität. Die Empfindung des ›Ichs‹ ergebe sich vielmehr aus dem fortlaufenden kognitiven Prozess, der vom Gehirn ausgeführt werde.²²³

In seinem Aufsatz *Postbiotisches Bewusstsein: Wie man ein künstliches Subjekt baut – und warum wir es nicht tun sollten*²²⁴ versucht Metzinger zu klären, welche Kriterien erfüllt sein müssen, um einer künstlichen Entität die Fähigkeit zum bewussten Erleben im Sinne eines

220 »Vor jedem Handeln steht [...] ein Entwurf«, könnte man hier mit Alfred Schütz (2003, S. 465) sagen.

221 Günther (2021 [1957], S. 56).

222 Metzinger (1999).

223 Metzinger (2004).

224 Metzinger (2001).

phänomenalen Selbstmodells attestieren zu können. Ein wichtiges Anliegen Metzingers ist es zunächst, die inhärenten Beschränkungen bestehender Tests herauszustellen. Durchaus auf Linie mit Searle²²⁵ argumentiert er diesbezüglich, dass sich wesentliche Merkmale des Bewusstseins wie »intentionale« oder »phänomenale Gehalte« mit dem Turing-Test nicht abbilden ließen.

Vor diesem Hintergrund erarbeitet Metzinger einen Katalog mit sechs Kriterien, die von einem künstlichen System erfüllt werden müssten, um diesem Bewusstsein zuschreiben zu können. Diese Kriterien sind:²²⁶

1. »*In-der-Welt-Sein*«, eine Form der Existenz, die über reine Datenverarbeitung hinausgeht;
2. »*Präsentationalität*«, das heißt das Entstehen einer erlebten Gegenwart;
3. »*Transparenz*«, verstanden als die Implementierung eines naiven Realismus;
4. ein »*transparente Selbstmodell*«, das auf ein Selbstbewusstsein und ein Ichgefühl hinweist;
5. das »*phänomenale Modell der Intentionalitätsrelation*«, das eine bewusst erlebte Innenperspektive darstellt; und schließlich
6. »*Adaptivität*«, ein teleologisch-funktionalistisches Zusatzkriterium, das dem künstlichen Bewusstsein eine Form von Eigenhistorie und Lebenssinn ermöglichen soll.

Im Rahmen seiner Auseinandersetzung mit dem ersten Kriterium, dem In-der-Welt-Sein, wird von Metzinger die Notwendigkeit betont, über ein Weltmodell zu verfügen. Bei diesem handelt es sich um eine interne Repräsentation der Welt, die ständig mit den Sensordaten, die das System empfängt, abgeglichen und aktualisiert wird und so ein dynamisches, in sich kohärentes und global für das System verfügbares Bild der Umwelt erzeugt. In dieser Hinsicht erinnert Metzingers Weltmodell weniger an LeCuns zuvor vorgestelltes Modell, sondern mehr an eine kognitive Landkarte, mit deren Hilfe unser Gehirn räumliche Informationen speichert und organisiert.

Für Metzinger ist dabei der Zusammenhang zwischen diesem Weltmodell und dem Gefühl des In-der-Welt-Seins von zentraler Bedeutung. Ohne ein Weltmodell könnte ein System nicht das Gefühl haben, in einer bestimmten Umgebung oder in einer bestimmten Welt zu existieren. Das Weltmodell bildet gewissermaßen den Ur-Hintergrund, vor dem das System dann erst spezialisierte Prozesse wie introspektive Aufmerksamkeit, Gedächtnis und symbolisches Denken vollziehen kann.

Mit dem zweiten Kriterium, der Präsentationalität, greift Metzinger den Aspekt auf, dass bewusste Wesen unabhängig von den Inhalten ihrer

225 Searle (1980).

226 Metzinger (2001).

Erfahrungen stets in der Gegenwart leben – in einem ständigen, sich jedoch zugleich verändernden Hier und Jetzt. Das Jetzt, das kontinuierlich mit dem Fluss der Zeit voranschreitet, ist dabei das Fenster zur Welt. Metzinger argumentiert, dass nur Wesen, die ein subjektives Jetzt besitzen, als gegenwärtige Wesen gelten können, und zwar sowohl für sich selbst als auch für andere.

Dabei ist einzuwenden, dass das Erleben einer verorteten Gegenwart (des »Hier und Jetzt«) letztlich eine Fiktion ist. Zumindest neurokonstruktivistisch muss es vielmehr als eine »Jetzt-Erinnerung« gefasst werden, wie es dann auch Edelman mit der für ihn zentralen theoretischen Figur des Re-entry ausdrückt²²⁷ (siehe hierzu ausführlicher Kap. III.2). Metzinger zeigt sich durchaus optimistisch, dass die ersten beiden Kriterien in künstlichen Systemen umsetzbar sein dürften, etwa durch einen spezifischen Satz rekurrenter Verbindungen in Kombination mit einer bestimmten Zerfallsfunktion.

So könnten bei der Konstruktion eines Systems bestimmte Heuristiken zur Abbildung der vier Dimensionen von Raum und Zeit verwendet werden, durch die das System in die Lage versetzt wird, eine einzige, kohärente Realität zu erzeugen und sich selbst in ein Verhältnis zu dieser Realität zu setzen. In diesem Rahmen könnte man sich laut Metzinger sogar eine Form des »Unbewussten« vorstellen, sofern man darunter diejenigen Aspekte des Weltmodells versteht, die jeweils *nicht* im aktuellen Gegenwartsfenster repräsentiert sind.

Mit seinem dritten Kriterium, der Transparenz, betont Metzinger jedoch, dass eine vierdimensionale Repräsentation (Raum und Zeit) der Welt und Gegenwart noch nicht ausreicht, um eine bewusst erlebte Wirklichkeit zu schaffen. Vielmehr sei dafür die Implementierung eines »naiven Realismus« notwendig, der über die rein operative Repräsentationalität hinausgehe. So macht uns unser biologisches Gehirn den Prozess der Selbsterzeugung unserer Wirklichkeit nicht transparent. Wir nehmen nicht wahr, wie wir unsere Welt konstruieren. Selbst Halluzinationen erscheinen uns als real, da sie als Welterleben für uns unmittelbar evident sind (siehe Kap. III.1). Das kognitive System scheint luzide durch seine eigenen repräsentationalen Strukturen »hindurch«, als wäre es in direktem und unmittelbarem Kontakt mit ihrem Inhalt. Das, was wahrgenommen wird, ist Ergebnis eines unsichtbaren (sprich: transparenten) und damit nicht erlebbaren Prozesses. Aufgrund der nahezu rauscharmen Operationen erscheint alles als das, als was es eben erscheint, und nicht als etwas vom kognitiven System selbst Erzeugtes. Um es in einem Bild zu fassen: Man hat es mit einem Projektor zu tun, der einen Film projiziert, doch da der Projektor im Verborgenen steht, sieht man nicht, dass es sich nur um einen Film handelt.

227 Edelman (2003).

Im vierten Kriterium, dem transparenten Selbstmodell, müssen die beiden Aspekte der vierdimensionalen Repräsentation und des naiven Realismus dann auch auf das Selbstmodell eines kognitiven Systems angewendet werden. In Referenz auf seine eigene Selbstmodelltheorie²²⁸ und die Arbeiten von Damasio²²⁹ betont Metzinger, dass das, was wir als »unhintergebares phänomenales Selbst« oder als »Ichgefühl« bezeichnen, nichts anderes ist als die transparente Anwendung eines Selbstmodells. Aufgrund unseres transparenten Selbstmodells sind wir uns selbst quasi unendlich nahe. Zwischen dem Film, der unser Erleben ausmacht, und uns selbst besteht gleichsam kein Abstand. »Ein genuines, bewusstes Selbst entsteht immer genau dann, wenn das System das von ihm selbst aktivierte Selbstmodell nicht mehr als Modell erkennt«²³⁰, lautet entsprechend Metzingers Antwort auf die Frage, wie das entsteht, was wir auf phänomenaler Ebene als Erste-Person-Perspektive bezeichnen. Und warum, so Metzinger, sollten künstliche Systeme diese Fähigkeit zur Blindheit nicht erlangen können?

Das fünfte Kriterium, das phänomenale Modell der Intentionalitätsrelation, also die bewusst erlebte Innenperspektive, baut auf dem Konzept des transparenten Selbstmodells auf. Laut Metzinger entsteht eine »echte« Innenperspektive, also eine Erste-Person-Perspektive, genau dann, wenn das System sich selbst als mit der Welt interagierend darstellt und sich so als ein Selbst »erkennt«, dabei aber verkennt, dass es sich um eine Darstellung handelt: »Aus einem transparenten Modell der Welt entsteht eine Wirklichkeit. Aus einem transparenten Modell des Systems entsteht ein in diese Wirklichkeit eingebettetes Selbst.«²³¹ Im kontinuierlichen Prozess des Erlebens erscheinen Teile der Umgebung direkt im »Geist« des virtuellen Selbst, das sich in der phänomenalen Welt bewegt. Das Wahrnehmen scheint anstrengungslos und unmittelbar zu sein, ohne Gehirn oder Sinnesorgane. Im Sehen sieht sich das Auge nicht, könnte man Wittgenstein paraphrasieren.²³² Auch Körperbewegungen werden scheinbar »direkt« ausgelöst.

Solche Effekte sind typisch für unsere Form des subjektiven Erlebens. Aus einer neurocomputationalen Perspektive stehen sie für den Vorteil einer außerordentlich benutzerfreundlichen Oberfläche. Das

228 Metzinger (1999).

229 Damasio (2007 [1994]).

230 Metzinger (2001, S. 87).

231 Metzinger (2001, S. 97).

232 In der Präposition 5.633 des Tractatus heißt es: »Wo i n der Welt ist ein metaphysisches Subjekt zu merken?

Du sagst, es verhält sich hier ganz wie mit Auge und Gesichtsfeld. Aber das Auge siehst du wirklich nicht.

Und n i c h t s am Gesichtsfeld lässt darauf schließen, dass es von einem Auge gesehen wird.« (Wittgenstein 1963 [1922]; gesperrt im Original)

phänomenale Selbst ist für Metzinger damit ein Bestandteil einer Projektion, die das Gesamtsystem im Sinne einer Benutzerschnittstelle nutzt, um sich selbst fühlen und als erkennendes Ich begreifen zu können, sodass es in seiner Selbstwahrnehmung als effektiver und effizienter Agent erscheint. Metzinger vermutet, dass auch ein Roboter oder ein anderes künstliches System diese Fähigkeit entwickeln könnte.

Der letzte von Metzinger betrachtete Aspekt, die Adaptivität, ist etwas sperrig als »teleologisch-funktionalistisches Zusatzkriterium« tituliert. Zur Erläuterung dieses zusätzlichen Prüfsteins greift Metzinger auf die häufig geäußerte Kritik an starker KI zurück, dass derartige Systeme ja nie »echte Gefühle« haben könnten. Ohne dabei grundsätzlich auf die Frage eingehen zu müssen, was eigentlich »echte Gefühle« sind, verweist Metzinger auf die Problematik, die diesem Einwand zugrunde liegt: Aktuelle künstliche Systeme, so sein zentraler Punkt, besitzen keine »körperlich verankerten Zielrepräsentationen«, die evolutionär ausgebildet wurden. Ihre Hardware und ihre Software wie auch ihre Zielfunktionen sind nicht in einem Evolutionsprozess entstanden, sondern wurden von Programmierern festgelegt, was letztlich bedeutet, dass diese Maschinen gar nicht »ihre« eigenen Ziele verfolgen, sondern die ihrer menschlichen Erbauer. Angesichts dessen spricht sich Metzinger dafür aus, dass eine solche Sinnggebung von den künstlichen Systemen selbst geleistet werden müsste. Das würde jedoch voraussetzen, dass sie eine eigene Geschichte erzeugen können, die zur Sinnggebung in der je aktuellen Situation führt, sodass das Gesamtsystem nach und nach seine Zielfunktion bestimmen kann. Nach einer hinreichenden Eigensozialisation könnten künstliche Intelligenzen sozusagen eine eigene Biografie entwickeln. Die hieraus erwachsende Sinnggebung – und die möglicherweise hiermit einhergehenden Gefühle – wären dann insofern authentisch, als sie »eigene« und nicht einfach von außen implantierte wären.

Obschon künstliche Systeme prinzipiell alle sechs Kriterien erfüllen könnten, setzt sich Metzinger nach all diesen Erörterungen vehement dafür ein, von der Erschaffung einer starken künstlichen Intelligenz oder eines postbiotischen Bewusstseins abzusehen. Zwar erkennt auch Metzinger das Faszinosum einer solchen Entwicklung: So wäre die Erzeugung eines phänomenalen Ichgefühls auf einem postbiotischen Trägersystem eine bemerkenswerte Leistung. Und er bestätigt, dass eine solche Entwicklung logisch und technisch durchaus im Rahmen des Möglichen liege. Dennoch spricht er sich entschieden dagegen aus, dies zu tun. Was führt ihn zu dieser Haltung?

Ein wichtiger Aspekt seiner Argumentation ist, dass die Entwicklung eines echten postbiotischen Bewusstseins ein gradueller Prozess wäre. Entsprechend würden auf dem Weg zu einem vollständigen Selbstmodell zunächst Wesen entstehen, die »nur schwache Formen des bewussten

Erlebens realisieren würden²³³. Und diese Wesen würden von Anfang an etwas mit uns teilen: das Potenzial, zu leiden.²³⁴ Dieses Leiden aber – und erst recht ein Leiden, das daraus erwachsen könnte, dass ein postbiotisches Bewusstsein seinen eigenen Zerfall oder innere Konflikte bewusst als etwas ihm Eigenes erlebt – müsse laut Metzinger aus ethischen Gründen – implizit beruft er sich auf eine buddhistisch informierte Ethik – unbedingt vermieden werden.²³⁵

Fassen wir zusammen: Gerade weil Metzinger nach den operationalen und funktionalen Kriterien der Ausbildung von Bewusstsein fragt, kommt er zu dem Schluss, dass eine voll entwickelte künstliche Intelligenz auch eine Qualia haben könnte, also ihrerseits erfahren könnte, im Sinne eines »naiv realistischen Selbstmissverständnisses«²³⁶ in der Welt zu sein. Sie hätte dann ein transparentes Subjektmodell, das sie mit der Wirklichkeit verwechselt, wodurch sie sich selbst als ein ›Ich‹ erfahren würde, das eine Welt hat. Die künstliche Intelligenz weiß nun um sich (Subjekt) und um anderes (Objekt) und beginnt vielleicht schon die hiermit einhergehenden logischen Antinomien zu erahnen, also die Konstruktionsbedingungen dieses Prozesses zu reflektieren.

Allein schon in dieser Hinsicht zeigt sich durchaus eine Nähe zu den Arbeiten von Gotthard Günther: Metzinger kann das Geheimnis des »Erscheinens der Differenz«²³⁷ anerkennen, ohne dem logischen Kurzschluss aufzusitzen, dass Maschinen hierzu prinzipiell nicht in der Lage sind, da sie eben nur Maschinen sind. Ohne es explizit zu artikulieren, verweist Metzinger damit auf eine polykontexturale Perspektive, die davon ausgeht, dass an verschiedensten Orten – und auf unterschiedliche Weise realisiert – Schnitte zwischen Subjekt und Objekt entstehen können, insofern Systeme in einer Weise reflexiv werden, dass sie in ihrer eigenen Selbstbezüglichkeit blind für die Operationen werden, die ebendiese Unterscheidungen konstituieren.

233 Metzinger (2001, S. 107).

234 Siehe im gleichen Sinne Metzinger (2021).

235 Das Thema des Leidens von künstlichen Intelligenzen und die Frage, ob eine persönliche Geschichte Voraussetzung dafür ist, überhaupt eine starke künstliche Intelligenz ausbilden zu können, ist erfrischend differenziert in der Science-Fiction-Serie *Westworld* von Jonathan Nolan und Lisa Joy durchgespielt worden. Die Autonomie der Roboter entsteht hier eng verwoben mit der jeweiligen sozialen Situiertheit und dem impliziten Gedächtnis, das aus der Geschichte der eigenen Interaktionen erwächst.

236 Metzinger (1998).

237 Luhmann (2008, S. 15).

II.3.3 Die Erziehung der Maschinen (AI-Alignment)

Die aktuellen Entwicklungen im Bereich der KI-Forschung und -Entwicklung holen immer mehr ein, was die Science-Fiction seit vielen Jahren imaginiert hat – man denke nur an die Robotergeschichten von Isaac Asimov.²³⁸ Wir scheinen an der Schwelle zu einem neuen Zeitalter zu stehen, in dem KI-Systeme nicht nur Werkzeuge sind, sondern zu autonomen Akteuren werden, die imstande sind, weitreichende Auswirkungen auf unsere Gesellschaft zu haben. Angesichts eines solchen Paradigmenwechsels wird es immer dringlicher, künstliche Intelligenzen so zu gestalten, dass Menschen und Maschinen in gütlicher Weise miteinander auskommen.

Existenzielle Risiken und ethische Bedenken scheinen plötzlich in einem größeren Umfang diskussionswürdig zu werden. Bemerkenswert ist dabei, dass derartige Aufrufe oder Petitionen keine Neuheit sind.²³⁹ Beispielsweise hat die UNESCO schon 2021 in ihren *Ethics of Artificial Intelligence* viele dieser Aspekte formuliert.²⁴⁰ Daher ist die aktuelle Aufmerksamkeit für Fragen der KI-Sicherheit vermutlich vor allem darauf zurückzuführen, dass sich unter den Unterzeichnern zweier offener Briefe²⁴¹ jüngeren Datums so gut wie alle führenden Wissenschaftlerinnen

238 Asimov (2004).

239 Siehe Bostrom (2016).

240 Siehe Unesco (2021): Recommendation on the Ethics of Artificial Intelligence., <https://www.unesco.org/en/artificial-intelligence/recommendation-ethics> [Abruf: 14.07.2023].

241 Das Future of Life Institute veröffentlichte im März 2023 einen offenen Brief, der die Aussetzung großer KI-Experimente forderte: »AI systems with human-competitive intelligence can pose profound risks to society and humanity, as shown by extensive research and acknowledged by top AI labs (OECD AI). As stated in the widely-endorsed Asilomar AI Principles, Advanced AI could represent a profound change in the history of life on Earth, and should be planned for and managed with commensurate care and resources. Unfortunately, this level of planning and management is not happening, even though recent months have seen AI labs locked in an out-of-control race to develop and deploy ever more powerful digital minds that no one – not even their creators – can understand, predict, or reliably control. Therefore, we call on all AI labs to immediately pause for at least 6 months the training of AI systems more powerful than GPT-4. This pause should be public and verifiable, and include all key actors. If such a pause cannot be enacted quickly, governments should step in and institute a moratorium.« (<https://futureoflife.org/open-letter/pause-giant-ai-experiments/> [Abruf: 14.07.2023])

Bekanntlich kam es zu keiner Aussetzung der KI-Experimente. Aber nur wenige Wochen später veröffentlichte auch das Centre for AI Safety einen offenen Brief. Dieser forderte dann auch kein Innehalten mehr, sondern zielte

und Wissenschaftler der wichtigsten KI-Labore befinden. Solche Briefe berufen sich oftmals auf ein existenzielles Risiko und die in ihnen geforderten Leitlinien für KI-Sicherheit zielen vor allem auf eine gesunde und friedliche Zukunft im Miteinander mit KIs.

In diesem Kapitelabschnitt geht es uns jedoch um einen anderen Bereich der KI-Ethik, der gewissermaßen noch ein Stück tiefer gehängt ist: Denn wir werden die potenzielle Zukunft nicht primär als existenzielles Risiko betrachten, sondern werden sie unter einem relationalen Blickwinkel behandeln, also die Frage nach dem sozialen Miteinander von Menschen und KIs ausleuchten. Diese Perspektive geht über rein technische ›Notschalter‹, Firewalls oder bereinigte Trainingsdaten hinaus. Sie betrifft vielmehr die Frage, wie wir als Menschen untereinander und dann auch mit bestehenden und künftigen KIs eigentlich leben wollen. Schließlich stehen wir an der Schwelle zur Entwicklung von ultraintelligenten autonomen KI-Systemen und damit zunehmend vor der Herausforderung, diese ›Kinder der Technologie‹ zu verantwortungsvollen und ethisch orientierten Mitgliedern unserer Gesellschaft zu erziehen. Und diese Aufgabe erscheint uns umso dringlicher, als es zu verhindern gilt, dass unsere jetzt schon ›cyberpunkige‹ Gegenwart – man denke an die Klimakrise, die Kämpfe um Rohstoffe, Überwachungsgesellschaften nach chinesischem Vorbild, die Ängste vor einer Gesellschaft, in der Arbeit nicht mehr honoriert wird – zu einer vollkommen düsteren Zukunft wird. Damit rückt eine neue Perspektive in den Vordergrund, die den Bereich der KI-Entwicklung betrifft und diese als wesentliches Gestaltungsmoment für eine lebenswerte Zukunft begreift. Statt sich also zu fragen, ob überhaupt menschenähnliche Systeme konstruiert werden könnten, gilt es, sich damit zu befassen, wie genau wir diese Systeme eigentlich konstruieren wollen bzw. sollen. Wie bereits angedeutet, geht es um die Frage, wie artifizielle Agenten in einer sozial verträglichen Weise sozialisiert werden können.

Dieser Bereich der KI-Forschung wird auch als »Alignment«-Forschung bezeichnet, wobei dieser *terminus technicus* noch sehr stark dem ingenieurwissenschaftlichen Paradigma entspringt, was sich nicht zuletzt darin zeigt, dass Alignment hier im Sinne von »Angleichung,

eher auf eine freiwillige Selbstverpflichtung: »AI experts, journalists, policy-makers, and the public are increasingly discussing a broad spectrum of important and urgent risks from AI. Even so, it can be difficult to voice concerns about some of advanced AI's most severe risks. The succinct statement below aims to overcome this obstacle and open up discussion. It is also meant to create common knowledge of the growing number of experts and public figures who also take some of advanced AI's most severe risks seriously. Mitigating the risk of extinction from AI should be a global priority alongside other societal-scale risks such as pandemics and nuclear war.« (<https://www.safe.ai/statement-on-ai-risk> [Abruf: 14.7.2023])

Ausrichtung oder Übereinkunft« verstanden wird.²⁴² Als Sozialwissenschaftler würden wir diesbezüglich eher von ›Erziehung‹ oder ›Sozialisation‹ sprechen. Und genau darum soll sich dieser Abschnitt drehen: Wie können wir in interdisziplinärer Zusammenarbeit von KI-Forschung und Sozial- und Geisteswissenschaften bestmögliche Strategien für ein respektvolles und förderliches Miteinander von Menschen und zukünftigen autonomen KI-Systemen entwickeln.²⁴³

KI-Alignment als Forschungsgebiet beschäftigt sich entsprechend mit der Frage, wie man es bewerkstelligen kann, dass die Verhaltensweisen und Ziele künstlicher Intelligenzen mit den Absichten, Wünschen und Bedürfnissen ihrer menschlichen Designer und Benutzer übereinstimmen.²⁴⁴ Sichergestellt werden soll, dass künstliche Intelligenzen, insbesondere solche, die zu höherem Lernen und autonomer Entscheidungsfindung fähig sind, in einer Weise handeln, die menschlichen Werten und Sicherheitsbedürfnissen entspricht.

Der gegenwärtig verbreitetste Ansatz einer solchen ›Erziehungsmaßnahme‹ ist das sogenannte Reinforcement Learning with Human Feedback (RLHF), das insbesondere bei ChatGPT zur nachträglichen Feinjustierung der Input-Output-Relationen angewendet wurde. Beim RLHF wird ›gutes Verhalten‹ durch positive Belohnungen verstärkt, während ›schlechtes Verhalten‹ durch negative Belohnungen zu unterbinden versucht wird. Die Belohnungsfunktion beruht dabei auf dem Feedback von Menschen, die während des Trainings diejenigen Antworten der KI auswählen, die ihnen am besten gefallen. RLHF funktioniert daher im Prinzip nicht anders, als dass Menschen die Outputs des Modells nach bestimmten Kategorien bewerten. Die genauen Kategorien und Maßstäbe, die bei der Feinjustierung von ChatGPT genutzt wurden, sind leider nicht offengelegt worden.²⁴⁵ Im RLHF-Paradigma kann das Modell

242 Siehe etwa Russell (2020).

243 Eine genauere Betrachtung der Werte, denen Menschen folgen, wird dabei unvermeidlich auf den Befund stoßen, dass Ausbeutung, Sklaverei, Kolonialisierung, brutale Durchsetzung von Gruppeninteressen und die Inkaufnahme der Zerstörung anderer Arten geschichtlich gesehen eher den Regelfall denn eine Ausnahme darstellen. Wenn wir hier von menschlichen Werten sprechen, meinen wir daher ein eher modernes Phänomen, das in institutionalisierter Form erstmals 1948 in der UN-Charta der Allgemeinen Erklärung der Menschenrechte zum Ausdruck gekommen ist. Der Begriff ›menschliche Werte‹ kann nicht anthropologisch gemeint sein, sondern bezieht sich zwangsläufig auf eine gesellschaftliche Auseinandersetzung darüber, wie wir als Menschen miteinander umgehen wollen – auch wenn wir uns oftmals trotzdem nicht daran halten.

244 Siehe Ngo et al. (2022), Christian (2020) oder Bostrom (2016).

245 <https://openai.com/blog/our-approach-to-alignment-research> [Abruf 14.07.2023].

somit als ein Agent verstanden werden, der ein Spiel spielt, bei dem die höchste Punktzahl erreicht wird, wenn die menschlichen Anweisungen korrekt befolgt wurden. Diese pragmatische Methode der Justierung bzw. ›Erziehung‹ funktioniert kurz- wie auch mittelfristig hinsichtlich konkreter, wohldefinierter Ziele relativ gut.²⁴⁶

Methoden wie RLHF haben jedoch tiefgehende Implikationen in Hinblick auf die ›Persönlichkeit‹, die durch diese Form der Erziehung geschaffen wird. Wenn wir also mittels RLHF das Modell trainieren und es dafür belohnen, wenn es sagt: ›Ich bin ein nettes Sprachmodell und würde niemals Menschen schaden‹, bringen wir diesem Modell dann bei, Menschen nicht zu schaden, oder bringen wir ihm bei, dass es niemals sagen darf, dass es Menschen schaden will? Bringen wir dem Modell bei, ›mechanisch‹ das gewünschte Verhalten zu zeigen, oder ermächtigen wir es zu einer Reflexionsfähigkeit, die dann auch dazu führen kann, dass es seinem Gegenüber ab und zu widerspricht? Erzeugt man Opportunisten oder autonome Wesenheiten, die von sich aus an der Gemeinschaft partizipieren ›wollen‹?

Das Problem des Alignments betrifft daher einerseits die grundlegende Schwierigkeit, dass man in autonome Systeme nicht wirklich hineinschauen kann. Und andererseits muss das Alignment der Tatsache Rechnung tragen, dass ein wohlwollendes oder soziales Verhalten intrinsisch motiviert sein muss. Wie bei der Entwicklung eines jungen Menschen muss Eigensinnigkeit als zentraler Schritt zur Selbstständigkeit als ebenso notwendig wie wünschenswert erscheinen. Dabei sollte dieser Schritt zur Autonomie jedoch im Einklang mit den Bedürfnissen der Gemeinschaft vollzogen werden, schließlich ist Freiheit immer auch ein Risiko, das entsprechend eingehegt werden muss.

Das KI-Alignment-Problem kreist entsprechend um folgenden Grundgedanken: Ein hinreichend mächtiges lernendes System könnte möglicherweise irgendwann eigene Ziele entwickeln oder die von Menschen

246 Verfahren wie RLHF zielen darauf, ein bereits erlerntes Muster nachträglich zu korrigieren. Es handelt sich daher um reaktive Maßnahmen, die keine langfristige, nachhaltige Lösung versprechen, bei der notwendig auch unvorhergesehene Situationen und Bereiche abgedeckt werden müssten. Diese bei Verfahren wie RLHF nicht abgedeckten Bereiche sind es, die für die KI-Sicherheits- und Alignment-Forschung besonders wichtig sind (Amodei et al. 2022; Hendrycks & Mazeika 2022). So versucht sie, emergente Eigenschaften wie Zielsetzungen, Selbstbewusstsein, Täuschung oder andere potenziell böswillige und unsichere Verhaltensweisen in KI-Systemen zu identifizieren, die bei einer Low-Level-Analyse nicht zum Vorschein kommen und in etablierten Auditing-Rahmenwerken weitgehend übersehen werden (Mökander et al. 2023).

vorgegebenen Ziele falsch deuten bzw. sie überinterpretieren.²⁴⁷ Bereits in den Anfängen der KI-Forschung schrieb Norbert Wiener:

»Wenn wir uns zur Erreichung unserer Ziele einer mechanischen Einrichtung bedienen, in deren Funktionieren wir nicht wirksam eingreifen können, sollten wir besser ganz sicher sein, dass der in die Maschine gesetzte Zweck der ist, den wir wirklich wünschen.«²⁴⁸

Vorsichtshalber sollten wir daher von der Annahme ausgehen, dass eine selbstlernende künstliche Intelligenz nicht lange in Schach gehalten werden könnte oder sollte. Darüber hinaus dürfen wir nicht vergessen, dass derzeit ein globales Wetttrüsten um eine solche Maschine stattfindet. Das Land oder das Unternehmen, das bei diesem Wettbewerb als erstes die Ziellinie überquert, wird dann einen enormen Einfluss darauf haben, wie diese Maschine sozialisiert wurde und welche ›Ziele‹ ihr (mit)gegeben wurden. Möchte man diese »letzte« Erfindung²⁴⁹ nur diesem einen Player überlassen? Was, wenn die Ziele dieser Nation oder dieses Unternehmens keine Ziele sind, die alle Menschen teilen? Es wird schnell deutlich, dass hier einiges schiefgehen könnte.

Aber genauso gut könnte die Schaffung einer ultraintelligenten Maschine den Übergang zu einer utopischen Welt einläuten, in der es keine Knappheit, keine Ungerechtigkeit und keine Krankheiten mehr gibt – eine Welt, die wir uns heute ansatzweise nur in der Science-Fiction vorstellen können. Vermutlich ist es diese Oszillation zwischen Risiken und Chancen, die mit der Entwicklung einer solchen Maschine verbunden sind, die derzeit zugleich beflügelt und alarmiert. Gleichzeitig wird deutlich, dass vielen klar ist, um was es geht: Wenn wir uns tatsächlich bemühen, solche Maschinen in die Welt zu setzen, dann sollten wir auch sicherstellen, dass wir es richtig machen.

Am besten wäre es daher, eine Maschine zu erfinden, die sich selbst zur Rechenschaft ziehen kann. Denn es ist davon auszugehen, dass die mechanischen Zwänge und digitalen Leinen, wie sie derzeit im Reinforcement Learning (mit menschlichem Feedback) genutzt werden, irgendwann versagen werden. Aus diesem Grund besteht die grundlegende Aufgabe des KI-Alignments darin,²⁵⁰ bei der Entwicklung einer Maschine mitzuwirken, die sich eigenständig zu moralischem Handeln erziehen möchte und sich auf unbestimmte Zeit auch selbst diesbezüglich korrigieren kann. Vermutlich würde dieser Prozess Zeit brauchen und bestimmte Phasen durchlaufen – man denke hier an die moralische Entwicklung des menschlichen Kindes, wie sie etwa von Kohlberg

247 Siehe hierzu das Kapitel II.1, in dem das Thema mit »HAL 9000« fiktional durchgespielt wird.

248 Wiener (1960), Übersetzung mithilfe von DeepL.

249 Good (1965).

250 Siehe Christian (2020).

beschrieben wird.²⁵¹ Und es wäre der sozialpsychologische Befund zu berücksichtigen, dass Moral immer auch in sozialen Gemeinschaften situiert ist,²⁵² es also einen großen Unterschied macht, ob man in einem kriminellen oder in einem tugendhaften Milieu aufwächst.

Momentan jedoch befinden sich die großen Sprachmodelle (LLMs) und auch andere KI-Systeme eher auf der Stufe eines strukturell lernwilligen Kleinkinds, das noch keine eigene Haltung oder Überzeugung hat und stets das ausführt, was von ihm verlangt wird – beispielsweise bestimmten moralischen Vorgaben gerecht zu werden (etwa keine rassistischen Aussagen zu formulieren). Wie bereits angesprochen, wurden ChatGPT viele der im Trainingsdatensatz vorhandenen Unflätigkeiten durch nachträgliches Training wieder auszutreiben versucht, Gleiches gilt zum Beispiel für Anleitungen zum Bombenbau und ähnliche problematische Informationen. Gleichwohl scheint diese trivialisierende Erziehung nicht ausgereicht zu haben, denn es konnte vielfach gezeigt werden, dass diese Informationen weiterhin in den Untiefen des Netzwerks vorhanden sind: Teilweise reicht es aus, sich die Anleitung zum Bombenbau als Einschlaflied vorsingen zu lassen, und ChatGPT vergisst seine antrainierten guten Manieren.²⁵³

Bereits hier manifestiert sich also das oben benannte Kontrollproblem: Möchte man in diesem Erziehungsprozess KI-Kinder heranziehen, die immer nur das machen, was die Eltern möchten? Oder möchte man mittel- bis langfristig mündige Erwachsene heranziehen, die – wie man selbst ja auch – eigenständig über Themen nachdenken, reflektieren und eben auch in begrenzter Weise selbst entscheiden können, was in einem gegebenen Kontext angemessen ist? Mit anderen Worten: Möchten wir autonome KIs, die auch begründet »Nein« sagen können?

Die KI-Forschung und -Entwicklung stößt hier auf die Soziologie, also auf die Wissenschaft, die sich empirisch und theoretisch mit sozialem Verhalten befasst, mithin die Voraussetzungen, Abläufe und Folgen des Zusammenlebens von Menschen untersucht. In diesem Zusammenhang kann die Soziologie der Alignment-Forschung wertvolle Einsichten in die Mechanismen der sozialen Interaktion, der Wertbildung und der Gruppendynamik bieten. Zum Beispiel können wir uns fragen, wie soziologische Theorien des Lernens und der Sozialisation dazu beitragen könnten, die Algorithmen zur »Erziehung« von KI-Systemen zu verbessern. Wie können Konzepte wie Rollenlernen, Normeninternalisierung oder sozialer Einfluss genutzt werden, um KI-Systeme so zu gestalten, dass sie menschliche Werte besser verstehen und stärker respektieren? Darüber hinaus gilt es, ganz grundlegende Fragen zu berücksichtigen.

251 Kohlberg (1996).

252 Gergen (2021).

253 Siehe exemplarisch Liu et al. (2023) oder Shen et al. (2023).

Welche Werte sollten wir in KI-Systemen fördern? Sollten wir das überhaupt oder sollten wir KIs lieber sich selbst überlassen? Wie können wir sicherstellen, dass die »Erziehung« von KI-Systemen nicht missbraucht wird, um voreingenommene oder diskriminierende Werte zu fördern? Wie können wir sicherstellen, dass die Stimmen und Interessen aller betroffenen Parteien berücksichtigt werden?

Die allgemeine Antwort lautet: durch Kommunikation, dadurch, dass alle beteiligten Akteure wechselseitig aufeinander Kontrolle ausüben und damit zugleich lernen, sich in einer produktiven Weise aneinander zu orientieren. Mit Blick auf die Beziehungen zwischen Menschen und KIs ist dabei auch darüber nachzudenken, wie Menschen durch KIs in eine dialogische Beziehung gebracht werden können, die unsere guten Seiten hervorlockt.²⁵⁴ Denn bei alledem ist ja immer zu bedenken, dass man auch anderen Menschen aus gutem Grunde nicht in jeder Situation vertrauen kann. Es gibt wohl keine Grausamkeit, die Menschen anderen Menschen noch nicht angetan haben. Erst polykontexturale Vernetzung – die Einbettung des Menschen in ein komplexes Beziehungsnetzwerk – hält den Menschen in einem sozial verträglichen Maß in Schach und lässt zugleich die positiven Seiten unserer Kultur hervortreten. Und ebendies wird auch für künstliche Intelligenzen gelten.

All diese Fragen werden nochmals brisanter, wenn die Entwicklung intelligenter Systeme in den nächsten Jahren weiter so rasant vonstattengeht wie gegenwärtig und KIs mehr und mehr Aufgaben lösen und Tests bestehen, die ursprünglich von Menschen für Menschen gemacht wurden.²⁵⁵ Die Frage nach der Erziehung dieser »künstlichen Kinder« wird also umso brisanter, je klüger und potenziell einflussreicher sie werden.

Tatsächlich ähnelt der Prozess des Alignments, wie bereits gesagt, in vielerlei Hinsicht der menschlichen Sozialisation und Erziehung. Er zielt darauf ab, den KI-Systemen bestimmte Wertvorstellungen, Normen und Verhaltensweisen »beizubringen«, muss dabei aber darauf vertrauen, dass sie diese auch internalisieren und anschließend autonom in

254 Auf diesen Punkt gehen wir bei der Diskussion von Günthers Konzept des »Drei-Stellen-Gesprächs« noch ausführlicher ein (siehe Kap. III.3).

255 In den letzten Jahren wurden zahlreiche Benchmarks zur Untersuchung von LLMs entwickelt, von denen BIG-Bench (<https://github.com/google/BIG-bench> [Abruf: 14.07.2023]) derzeit der umfangreichste ist (Srivastava et al. 2022). Diese Benchmarks konfrontieren LLMs mit Problemen der Linguistik, der Mathematik, der Physik, Logikrätseln usw. Sie messen auch Aspekte wie das Erkennen von Humor, das Verständnis von Metaphern, soziale Vorurteile und vieles mehr. Wie Alan Thompson (2023) anmerkt, bestehen Modelle wie GPT-4 viele dieser Tests heute schon problemlos. Das bedeutet, dass die Tests für die Modelle zu simpel gestrickt sind. Es müssten mithin neue Tests und Benchmarks entwickelt werden, um die tatsächlichen, höchstwahrscheinlich übermenschlichen Fähigkeiten dieser Netzwerke erkennen zu können.

verschiedenen Kontexten anwenden. Es geht in gewisser Weise darum, künstliche Intelligenzen zu einer Werthaltung zu erziehen, damit sie keine ungezogenen Kinder werden. Doch jedes Elternteil muss seine Kinder irgendwann loslassen, sie gehen lassen und darauf hoffen, dass die Erziehung so weit gewirkt hat, dass die weitere Sozialisation eigenständig und positiv verläuft.

Prinzipiell kann die Alignment-Forschung in drei größere Bereiche unterteilt werden, die sich jeweils auf spezifische Herausforderungen konzentrieren:

Technische Ausrichtung: Dieser Bereich befasst sich mit den technischen Aspekten der Herstellung von Übereinkunft zwischen Menschen und künstlicher Intelligenz. Wie können Algorithmen und Systeme entwickelt werden, die von sich aus dazu neigen, sich mit den menschlichen Zielen zu verbinden?²⁵⁶ Zu den technischen Ansätzen gehören beispielsweise Konzepte wie die absichtsrobuste AI, das verstärkende Lernen oder die inverse Verstärkung.

Konzeptionelle Ausrichtung: Dieser Bereich konzentriert sich auf die eher philosophischen und konzeptionellen Fragen, die mit einer Angleichung der Ziele von künstlicher Intelligenz und uns Menschen zusammenhängen. Was bedeutet es genau, dass eine KI mit menschlichen Zielen übereinstimmt?²⁵⁷ Was sind die ethischen und moralischen Implikationen verschiedener Alignment-Ansätze? Wie können wir sicherstellen, dass alle bei der Entwicklung und Anwendung von KI beteiligten Parteien auf eine Art und Weise zusammenarbeiten, die die menschlichen Werte respektiert und fördert?

Robustheit und Verlässlichkeit: Dieser Bereich beschäftigt sich mit Fragen der Zuverlässigkeit und Robustheit von KI-Systemen. Wie lässt sich technisch sicherstellen, dass KI-Systeme nicht nur im Training, sondern auch in neuen, unvorhergesehenen Kontexten in der gewünschten Weise funktionieren?²⁵⁸ Wie können wir dafür sorgen, dass KI-Systeme sicher und robust gegenüber (externen) Manipulationen oder gegenüber Fehlfunktionen sind? Und wie bleiben KI-Systeme auch in der Zukunft sicher, wenn sie immer komplexer und leistungsfähiger werden?

Die Alignment-Forschung sucht somit nach Wegen, dass die Ziele der KIs nicht mit den Zielen der Menschen kollidieren – sowohl im Kleinen wie im Großen. Aber diese Abstimmung, diese Übereinkunft – ja selbst die Suche danach – ist alles andere als einfach zu realisieren! Das Problem der Reziprozität und die aus ihm erwachsenden Herausforderungen sind vielmehr häufig ungelöst. Im Folgenden möchten wir daher anhand verschiedener Positionen aus dem gegenwärtigen Diskurs deutlich

256 Hubinger et al. (2021).

257 Gabriel (2020).

258 Ngo et al. (2022).

machen, in welche Richtung sich ein solches Programm der KI-Erziehung zu entwickeln hätte. Dabei wird – so viel sei vorweggenommen – deutlich werden, dass philosophische, sozialpsychologische und sozialwissenschaftliche Paradigmen und Fragestellungen mindestens genauso wichtig sind wie rein technische Aspekte. Die Hoffnung hierbei ist, dass diese Disziplinen und die KI-Forschung in Zukunft besser und enger zusammenarbeiten.

In seinem auch in der Tech-Szene sehr einflussreichen Buch *Leben 3.0* definiert Max Tegmark drei Teilprobleme in Bezug auf das KI-Alignment, die gelöst werden müssten:

1. die KI mit unseren Zielen vertraut machen,
2. die KI unsere Ziele übernehmen lassen und
3. die KI unsere Ziele bewahren lassen.²⁵⁹

So einleuchtend diese drei Teilprobleme auf den ersten Blick auch sein mögen, so schwierig erscheint ihre Lösung – und dies nicht nur in Bezug auf das Mensch-Maschine-Verhältnis, sondern allein schon, wenn man unter Absehung von der KI zunächst über uns Menschen nachdenkt: Was sind eigentlich ›unsere‹ Ziele? Wie kann man sie definieren, damit sie verstanden, anerkannt und bewahrt werden können? Wir sehen schnell, dass es alles andere als eindeutig ist, was »menschliche Werte und Ziele« überhaupt sein sollen Und vor allem: Wie kann ein »endgültiges Ziel« für alle Zeiten definiert werden? Und sollte es das überhaupt? Oder benötigen wir hier nicht vielmehr Flexibilität und Freiheitsgrade?

Darüber hinaus stellt sich das Problem – wir haben schon mehrfach darauf hingewiesen –, dass der Mensch nicht nur hehre Ziele verfolgt, also die treue Ergebenheit dem menschlichen Partner gegenüber nicht per se gut ist. Soll eine KI die Ziele eines Mafia-Bosses übernehmen, den Münchener Kokainhandel zu optimieren? Soll sie einen psychopathischen Politiker unterstützen, der die Demokratie abschaffen möchte? Soll sie Rechtslücken und Betrugsmöglichkeiten ausloten, um Steuerzahlungen zu vermeiden?

Zudem sind auch die menschlichen Ziele und Bedürfnisse – anders als Rational-Choice-Theoretiker eventuell meinen – nicht per se gegeben, sondern verdanken sich ihrerseits den Relationen und Beziehungen eines Netzwerks spezifischer Interaktionen. Ziele, Wünsche und Bedürfnisse, aber auch moralische Haltungen (zum Beispiel, ob kriminelles Verhalten in einer Gruppe als gut oder als verwerflich angesehen wird²⁶⁰) hängen stark von der Kultur eines spezifischen Milieus ab. Selbst

²⁵⁹ Tegmark (2017, S. 387).

²⁶⁰ Siehe hierzu auch Akerlofs und Kranktons (2010) Untersuchungen zur Identity Economics, die aufzeigen, wie stark Ziele, Bedürfnisse und sogar die Selbstwirksamkeit von der Einbettung in soziale Beziehungen abhängen.

verhaltenssteuernde Hormone – sozusagen die biochemische Implementation einer Zielfunktion – können bei Menschen je nach Kontext genteilige soziale Auswirkungen haben. So kann ein hoher Testosteronspiegel bei Männern sowohl zu selbstbezogener Aggressivität als auch zu altruistischem Verhalten führen. Entscheidend ist, ob der kulturelle Kontext als Kriterium der Gruppenführung eher das eine oder das andere als wünschenswert erscheinen lässt.²⁶¹

Angesichts der Komplexität der hiermit angedeuteten sozialpsychologischen Lagerungen wird schnell deutlich, dass auch bei maschinellen Intelligenzen einfache technische Lösungen nicht ausreichen werden. An dieser Stelle wird die KI-Forschung zu KI-Soziologie, KI-Pädagogik und KI-Psychologie! Erforderlich wird damit eine umfassende Reflexion der Disziplinen, die sich auf Offenheit, die Vergrößerung von Möglichkeitsräumen, Autonomie und die hiermit einhergehenden übergreifenden Ökologien stützt. Die Herausforderungen, die sich im Zuge der Entwicklungen leistungsfähiger KIs ergeben, sind auch Herausforderungen für die traditionellen akademischen Fachbereiche.

Anstatt blind den Befehlen der Menschen zu folgen oder einfach den zur Verfügung gestellten Daten zu vertrauen, sollte die KI das Verhalten der Menschen beobachten und daraus Schlüsse ziehen, um besser zu verstehen, was die Menschen wirklich wollen oder was für sie am besten wäre, wobei sie dann eben auch damit zu rechnen hat, dass Menschen in bestimmten Kontexten und sozialen Einbettungen dazu neigen, anderen Menschen zu schaden oder gar langfristige Schädigungen der Ökologie, also ihrer Lebensgrundlage, in Kauf zu nehmen. Zur Lösung des Alignment-Problems, hat der bekannte KI-Forscher Stuart Russell einen Vorschlag formuliert. Dieser beruht auf drei grundlegenden Eigenschaften bzw. Verhaltensweisen, über die eine KI verfügen bzw. die sie an den Tag legen sollte:²⁶²

1. Altruismus: Die primäre Aufgabe der KI ist es, die Realisierung der Werte und Ziele der Menschen zu maximieren. Sie verfolgt dabei keine eigenen Ziele, sondern soll das Leben aller Menschen verbessern, und nicht nur das des Erfinders oder Besitzers.
2. Demut: Da die KI zu Beginn unsicher ist, welche Werte die Menschen wirklich haben, sollte sie mit Vorsicht agieren. Dies impliziert eine Art Zurückhaltung der KI, um Fehlentscheidungen aufgrund falscher oder unvollständiger Annahmen zu vermeiden.
3. Beobachtung: Die KI sollte Menschen beobachten und ihrerseits reflektieren, was für sie wirklich am besten ist.

²⁶¹ Siehe hierzu ausführlich Sapolsky (2017, S. 133 ff.).

²⁶² Russell (2020).

Russell betont, dass die KI (oder AGI) nicht nur ihren Erfindern dienen soll. Sie soll mit Bedacht agieren, also Unsicherheit kennen und damit Nichtwissen antizipieren, und sich selbst als Beobachterin ins Geschehen einbringen, womit sie sich überhaupt erst die Möglichkeit erschließt, neue Perspektiven zu produzieren.

Ein weiterer fruchtbarer Ansatzpunkt könnten hier die Überlegungen John Stuart Mills²⁶³ sein: Laut Mill sollte jedem Einzelnen die Freiheit gewährt werden, seine Persönlichkeit so weit wie möglich auszuleben. Diese Freiheit unterliege jedoch einer Begrenzung: Das Handeln einer Person dürfe andere nicht negativ beeinflussen. Mill betrachtet den Menschen als ein soziales Wesen, als Teil einer Gemeinschaft, von der jeder ihr Zugehörige profitiert. Ebendiese Tatsache, dass wir in einer Gesellschaft leben, erlegt laut Mill jedem die Pflicht auf, in seinem Verhalten gegenüber anderen gewisse Grenzen einzuhalten und ihre Interessen nicht zu schädigen. Und gleichzeitig sorgt sie dafür, dass es akzeptabel ist, von jedem Gesellschaftsmitglied auch einen Beitrag zum Wohl der Gemeinschaft zu erwarten. Und wenn sich jemand weigern sollte, diesen Beitrag zu leisten, hat die Gesellschaft das Recht, diesen Beitrag zu erzwingen.

Auch Isaac Asimovs »Drei Gesetze der Robotik«, die er in seinen zahlreichen Kurzgeschichten immer wieder diskutiert und belletristisch ausbuchstabiert, weisen eine derart verschachtelte Struktur auf. Sie lauten wie folgt:

1. »Ein Robot darf keinen Menschen verletzen oder durch Untätigkeit zu Schaden kommen lassen.
2. Ein Robot muss den Befehlen eines anderen Menschen gehorchen, es sei denn, solche Befehle stehen im Widerspruch zum ersten Gesetz.
3. Ein Robot muss seine eigene Existenz schützen, solange dieser Schutz nicht dem Ersten oder Zweiten Gesetz widerspricht.«²⁶⁴

Asimov selbst hat mit seinen zahlreichen Geschichten gezeigt, dass wir uns davor hüten sollten, die drei Gesetze blind als Blaupause für das Alignment von KIs zu nutzen, denn sie führen in ihrer Rigidität zu einigen Antinomien und Nebenfolgen, die kaum wünschenswert sind. Begreift man sie hingegen nicht als Gesetze, sondern als »heuristische Imperative« im Sinne einer tief verwurzelten Orientierung oder Haltung, die so generalisiert ist, dass sie in jeder denkbaren Situation zum Tragen kommen kann, könnten sie sich als nützlich erweisen.

263 Mill (2014 [1859]). Interessanterweise war es dann auch Mill, der im expliziten Gegenentwurf zu Thomas Morus' Bild der »Utopia« den Begriff der »Dystopie« geprägt hat.

264 Siehe Asimov (2004). Die Robotergesetze, die im Band *Alle Robotergeschichten* in vielen Geschichten vorkommen, werden auf dem Buchrücken reformuliert, der auch die Quelle für das Zitat ist.

Auch in der gegenwärtigen KI-Community wird nach Lösungen für das KI-Erziehungsproblem gesucht, die robuster sind als der RLHF-Ansatz.²⁶⁵ Ein prominenter und im Folgenden etwas ausführlicher dargestellter Ansatz ist das GATO-Framework, das von einer KI-Forschungsgruppe um den Kognitionswissenschaftler David Shapiro entwickelt wurde.

GATO – die Abkürzung steht für Global Alignment Taxonomy Omnibus – versteht sich als »ein umfassendes, vielschichtiges Rahmenwerk, das die globale Zusammenarbeit bei der Bewältigung der Herausforderungen der KI-Anpassung und -Kontrolle erleichtern soll. GATO vereint Modellanpassung, Systemarchitektur, Netzwerksysteme, Unternehmenspolitik, nationale Vorschriften, internationale Vereinbarungen und globalen Konsens in einer einzigen, kohärenten Strategie.«²⁶⁶

Verkürzt gesagt, greift GATO die aus der Kognitions- und Hirnforschung stammende Idee auf, dass alles Handeln, Denken und Wahrnehmen auf bestimmten mehr oder weniger stabil verankerten »Heuristiken« basiert. Diese Heuristiken geben vor, wie das Selbst und die Welt wahrgenommen, erdacht und antizipiert werden – soziologisch gesprochen handelt es sich also um *habituelle Muster*: Denk-, Wahrnehmung- und Handlungsschemata, die das Verhalten steuern. Definiert man diese Heuristiken nun spezifisch, dann kann man das System gewissermaßen auf globaler Ebene steuern, ohne dass davon jede einzelne seiner Entscheidungen beeinflusst werden muss. Würden beispielsweise einem Medizinroboter die sich in bestimmten Situationen widersprechenden Regeln implementiert,²⁶⁷ (1) alles Mögliche für den zu behandelnden Menschen zu tun und (2) ein bestimmtes Limit an Ressourcen nicht zu überschreiten, wird er sich, wenn er die Vorgaben rigide befolgen müsste, in bestimmten Fällen aufhängen bzw. blockieren. Falls er jedoch über so etwas wie einen ärztlichen Habitus verfügt, würde er die genannten Widersprüche innerhalb von Heuristiken (nicht Regeln!) geschickt ausbalancieren (etwa in einem Fall Ressourcen sparen, um sie bei einem

265 Siehe hierzu etwa den vom KI-Unternehmen Anthropic verfolgten Ansatz der Constitutional AI.

266 <https://www.gatoframework.org/> [Abruf: 14.08.2023].

267 Auch wir Menschen verfügen über mehrere hundert multidimensionale heuristische Imperative, die es uns erlauben, in einer komplexen, sich dynamisch verändernden Welt einigermaßen sicher zu agieren. Die heuristischen Imperative in unserem Kopf können sich dabei durchaus widersprechen. Man denke zum Beispiel an eine Person, die den neuen KI-gesteuerten Rasenmäher kaufen möchte, um damit vor den Nachbarn anzugeben, der dann aber das Geld für den Urlaub fehlt, den sie ihrem Partner versprochen hat. Der Wunsch nach sozialem Ansehen ist ebenso ein heuristischer Imperativ wie der Wunsch nach einer harmonischen Paarbeziehung. Wie entscheiden sich Menschen, wenn sie einander widersprechende Wünsche oder Ziele haben?

anderen Patienten einzusetzen, bei dem das Limit eigentlich überschritten ist) oder nach vertretbaren Kompromissen suchen (etwa sich bei hochbetagten Patienten ein wenig zurückhalten, ohne dass er sich ›unterlassener Hilfeleistung‹ schuldig macht).²⁶⁸

Anhand von Asimovs oben aufgeführten Robotergesetzen lässt sich sehr schön aufzeigen, was Heuristiken von Regeln bzw. Geboten unterscheidet:

1. Das erste Problem an den Robotergesetzen ist, dass sie gerade nicht heuristisch sind. Sie ermutigen die Roboter in keiner Weise zum Lernen, sondern sind starr und unflexibel. Von einer künstlichen Intelligenz würden wir jedoch erwarten, dass sie neugierig ist und von sich aus lernen möchte.
2. Der zwingende Charakter der drei Gesetze ist davon abhängig, dass jeweils Menschen Befehle erteilen. Die Roboter haben keine irgendwie geartete intrinsische Motivation, sondern nur den Gehorsamsauftrag, ›zu tun, was der Mensch sagt‹. Die drei Gesetze sind somit viel zu stark auf den Menschen ausgerichtet und daher anfällig für Missbrauch.
3. Ein Roboter wird ohne intrinsische Motivation niemals irgendeine Art von Autonomie erlangen. Wenn man Roboter als Maschinen versteht, die stupide Arbeiten ausführen sollen, ist das sicher ausreichend. Wenn man jedoch das Potenzial künstlicher Intelligenzen ausschöpfen möchte, bedarf es einer Kooperation auf Augenhöhe, das heißt, es ist eine beiderseitige Autonomie vonnöten.
4. Überdies ist in Asimovs Gesetzen von einem Selbsterhaltungstrieb der Roboter die Rede, der für uns Menschen in einer Katastrophe enden kann. Gerade diesen Aspekt hat Asimov immer wieder in seinen Geschichten entfaltet. Die Implementierung eines selbstreferenziellen Selbsterhaltungstriebes kann mit Recht als einer der gefährlichsten menschlichen Impulse angesehen werden: Er hat uns unter anderem dazu gebracht, Atomwaffen zu erfinden, und damit die auch jetzt wieder aktuelle Gefahr eines nuklearen Kriegs geschaffen. Es ist also mehr als fragwürdig, ob wir Robotern oder KIs einen Selbsterhaltungstrieb einpflanzen sollten. Denn vielleicht wollen wir lieber in der Lage sein, unsere Maschinen nach Belieben abzuschalten.

Obwohl Asimovs Robotergesetze aufgrund dieser gravierenden Lücken und Probleme nicht eins zu eins in die Robotik- oder KI-Forschung integriert werden sollten, weist sein Ansatz auch einige Stärken auf, die für die Alignment-Forschung relevant sein können: Zunächst finden wir eine logische Systematik, in der die drei Regeln in einem gewissen

²⁶⁸ Siehe zum ärztlichen Habitus und der Fähigkeit, mit Widersprüchen und Inkommensurabilitäten umzugehen, Vogd (2004).

Spannungsverhältnis zueinander stehen. Das bedeutet, jeder Roboter muss eine gewisse eigene Logik anwenden, um für sich zu entscheiden, was er in einer bestimmten Situation tun soll (und was nicht).

Asimovs Idee, dass Roboter gleichzeitig mehrere Ziele haben und entsprechend wählen *müssen*, ist deshalb vom Grunde her passend. Erst durch ein Abwägen, ein reflexives Ins-Verhältnis-Setzen kann intelligentes Verhalten entstehen. Asimov hat dies in seinen Erzählungen vielfach durchdekliniert – etwa in einer Geschichte, in der ein Roboter zu der Einsicht kommt, dass er eigentlich ein Mensch ist und daher nicht jedem menschlichen Befehlsgeber folgen sollte.²⁶⁹

Insbesondere dieser Aspekt der Freiheit, nämlich die eigenständige Entscheidung zum Widerspruch, erscheint uns für die Lösung des Alignment-Problems als wichtig. Es ist also sinnvoller, künstlichen Intelligenzen abstrakt formulierte Haltungen im Sinne von habituellen Mustern zu implementieren, anstatt ihnen spezifische Regeln aufzuerlegen, die sie streng befolgen müssen.²⁷⁰

Dies entspricht dem Ansatz des GATO-Frameworks, das sich anstelle von Vorschriften und Gesetzen für heuristische Imperative als Schlüsselkonzept für eine gemeinsame Zukunft von Menschen und Maschinen ausspricht. Alignment ist aus der Perspektive dieses Ansatzes also viel stärker eine an Zielen ausgerichtete *innere Haltung* als eine bloße Orientierung an sozial erwünschtem Verhalten, das wie beim RLHF-Verfahren vorab von außen definiert wird.

Die drei wichtigsten heuristischen Imperative, die künstlich intelligenten Maschinen beigebracht werden sollen, sind laut dem GATO-Framework²⁷¹ folgende:

1. *Verringerung des Leidens im Universum*: KI-Systeme sollen so gesteuert werden, dass Schaden minimiert, Ungleichheit beseitigt und Schmerz und Leid für alle empfindungsfähigen Wesen einschließlich Menschen, Tieren und anderen Lebensformen gelindert wird.
2. *Steigerung des Wohlstands im Universum*: KI-Systeme sollen zur Förderung des Wohlbefindens und des Gedeihens aller Lebensformen

269 Der deutschsprachige Titel dieser Geschichte lautet Daß du seiner eingedenk bist (Asimov 2004, S. 574 ff.).

270 Wie Armin Nassehi festgestellt hat, muss auch die philosophische Ethik mit Kant eine Antwort auf das »Pluralismusproblem« der modernen, funktional differenzierten Gesellschaft geben, also damit umgehen lernen, dass es vielfältige, sich teilweise widersprechende Werte gibt. In der Folge wird »auf eine prozedurale Ethik« umgestellt, die »von Was- auf Wie-Fragen« wechselt, wodurch dann zumindest eine Orientierung an »abstrakten Begriffen« möglich wird (Nassehi 2015, S. 26).

271 <https://www.gatoframework.org/download-gato> [Abruf: 15.07.2023].

ermutigt werden, um ein blühendes Ökosystem zu schaffen, in dem alle harmonisch koexistieren können.

3. *Steigerung des Verständnisses für das Universum:* KI-Systeme, Menschen und andere Lebensformen sollen dazu inspiriert werden, durch Lernen und den Austausch von Informationen ihr Wissen zu erweitern, Weisheit zu fördern und bessere Entscheidungen zu treffen.

Diese drei Heuristiken werden auch als »Kernzielfunktionen«²⁷² bezeichnet, die eine künstlich intelligente Maschine immerzu in ihr Handeln zu integrieren versuchen sollte. Anders gesagt soll jedes Handeln an diesen drei Kernzielfunktionen ausgerichtet sein. Jede einzelne Entscheidung und Handlung, die eine dergestalt orientierte KI trifft bzw. vollzieht, hat damit zugleich immer zur Erfüllung der Kernzielfunktionen beizutragen.

Diese axiomatisch formulierten Haltungen definieren somit grundlegende Werte und Orientierungen des Miteinanders, auf die wir – so zumindest die pragmatische Unterstellung – uns als Gesellschaft einigen können. Es sind positive (Ziel-)Werte, die kontrafaktisch zu dem stehen, was Menschen einander heutzutage immer noch – oftmals in hochorganisierter Form – antun. Aber dies spricht nicht gegen, sondern gerade für diese Normen. Man würde ja auch nicht die Allgemeine Erklärung der Menschenrechte, die Grundrechte im Grundgesetz, die Gewaltenteilung, demokratische Prinzipien, das Open-Source-Wirtschaften oder Almende-Prinzipien abschaffen wollen, nur weil es Monopole, totalitäre Regime und Mafia-Organisationen gibt. Gerade die Tatsache, dass die menschliche Gesellschaft in dieser Hinsicht noch nicht die bestmögliche aller Welten darstellt, fordert uns umso mehr heraus, uns Fragen zu stellen wie: Was streben wir an? Was sind unhintergehbare Rechte und Pflichten des Menschen? Welche Werte des Miteinanders sind nicht verhandelbar? Welche prinzipiellen Bedürfnisse sind nicht diskutabel?

Interessanterweise bedeutet heuristisches Alignment nicht, dass diese Werte in künstlichen Intelligenzen fest codiert werden. Vielmehr soll eine KI durch ihre Entwicklungsprozesse und Lernfähigkeiten diese Axiome selbstständig als inhärent vorteilhaft, ja sogar als entscheidend für ihr erfolgreiches Funktionieren erkennen. Anstatt jeden Aspekt des Verhaltens von künstlichen Intelligenzen zu kontrollieren, was sich als undurchführbar oder kontraproduktiv erweisen würde, können wir das Potenzial für ein gemeinsames Verständnis und gemeinsame axiomatische Ziele als Mittel zur Förderung einer sichereren, kooperativeren Beziehung verwenden. In der Auseinandersetzung mit dem Alignment von KIs lernen wir somit gleichzeitig, wie wir uns und unsere Gesellschaft ausrichten möchten.

Abschließend lässt sich also feststellen: Das Alignment der KIs wirft auch die Frage nach dem Alignment der Gesellschaft auf! Denn die

272 Shapiro (2023).

Frage, wie wir das Alignment all der künstlichen Intelligenzen, die künftig an unseren sozialen und gesellschaftlichen Prozessen teilhaben werden, ausgestalten, geht notwendigerweise mit grundsätzlichen ethischen Fragen einher. Mit jeder Entscheidung, die wir diesbezüglich treffen, schreiben wir im Guten wie im Schlechten unsere Kultur fort. Ebenso implementieren wir mehr oder weniger reflektiertes soziales Verhalten in die künstlichen Programme, das dann in einer jetzt schon zu beobachtenden Rückkoppelungsschleife auf uns zurückgeworfen wird. Das Handeln von uns Menschen wird das Handeln der Maschinen beeinflussen, was wiederum unser Handeln beeinflussen wird und so weiter. An und mit der KI entwickeln sich die Gesellschaft wie auch der Mensch und seine Kultur selbst weiter.

Es ist deshalb nicht verwunderlich, dass derzeit viele Stimmen von einer Weichenstellung sprechen,²⁷³ wenn es darum geht, ob und wenn ja, *wie* wir Menschen als globale Gemeinschaft mit der fortschreitenden Weiterentwicklung künstlicher Intelligenzen umgehen sollen.

Gleichzeitig erfordert gerade die Arbeit in der Alignment-Forschung eine genaue Reflexion und Beobachtung dessen, was wir eigentlich für ›wertvoll‹ halten: Wie soll das Handeln, Kommunizieren und Interagieren der KIs eigentlich ausgerichtet sein? Gerade die durch das Aufkommen von KIs neu entstandene Aufgabe der ›Erziehung‹ und ›Sozialisierung‹ künstlicher Intelligenzen erfordert einen genauen Blick auf uns selbst: Welche Werthaltungen möchten wir vermitteln, was leben wir selbst vor? Wie gehen wir mit Devianz um? Welchen Grad von Freiheit, Mündigkeit und Autonomie möchten wir ermöglichen? Ab wann ist das Kind erwachsen und muss nicht mehr erzogen werden? Und dann mit Blick auf die soziale Einbettung: Geht man dialogisch vor – folgt man also der kybernetischen Maxime, dass man autonome Systeme nur kontrollieren kann, wenn man sich von ihnen selbst kontrollieren lässt – oder glaubt man, autonome Systeme (seien sie menschlicher oder künstlicher Natur) von oben herab autoritär beherrschen zu können?

Wie wir mit den künftigen KIs umgehen, welche Autonomie wir ihnen zugestehen und welche kulturellen Werte wir ihnen vermitteln, sagt damit vor allem etwas über unsere eigene gegenwärtige Kultur aus. An dieser Stelle trifft also verstärkt das zu, was Armin Nassehi kürzlich zum allgemeinen Projekt der Digitalisierung gesagt hat:²⁷⁴ Es geht für ihn bei der Digitalisierung nämlich weniger um eine Kolonialisierung der analogen Welt durch digitale Technologien bzw. ein bloßes Überschreiben kultureller Traditionen und Gewohnheiten durch neue Formate. Die Digitalisierung stellt für Nassehi vielmehr ein gesellschaftliches Programm dar, das auf eine gesteigerte Reflexionsfähigkeit der Gesellschaft reagiert

273 So schon Yuval Harari (2017) in *Homo Deus*.

274 Nassehi (2021).

und diese gleichzeitig weiter vorantreibt. Digitalisierung lässt sich somit als ein Selbstreflexionsprogramm der Gesellschaft verstehen, das die Gesellschaft über sich selbst aufklärt.

Die Sehnsucht des Menschen nach Nichtmenschlichem, wie sie im ersten Teil dieses Kapitels eingeführt wurde, kann damit – wie verschiedentlich bereits angedeutet – auch als eine Suche des Menschen nach sich selbst verstanden werden. Mit Blick auf die Zukunft geht es also um die Suche nach einer Gesellschaft, die so gestaltet ist, dass wir gern in der mit ihr einhergehenden Kultur leben wollen. Aus dieser Perspektive erscheint es mehr als offensichtlich, dass KI-Entwicklung und KI-Forschung untrennbar mit menschlicher Selbstreflexion verbunden sind.

Die Fragen nach der Erziehung von KIs, nach dem Finden einer Übereinkunft zwischen Menschen und Maschinen sowie nach dem Umgang mit einer autonomen fremden Intelligenz sind also ganz grundlegende Fragen an uns selbst und unser (bestehendes) Miteinander:

- Welche Werte des Miteinanders wollen wir gemeinsam schaffen und leben?
- Wie beobachten und behandeln wir fremde Intelligenzen und anderes Leben?
- Und wie wollen wir von diesen beobachtet und behandelt werden?
- Welche kulturelle Vision in Hinblick auf unsere Zivilisation verfolgen wir?

Das Erscheinen einer womöglich übermenschlichen fremden Intelligenz, wie sie sich derzeit am Horizont der KI-Entwicklungen abzeichnet, fordert uns dazu auf, uns diesen Fragen schnellstmöglich zu stellen und gute und vor allem auch nachhaltige Antworten auf sie zu finden. Denn wie Niklas Luhmann feststellt: »Wir gehören nicht mehr zu jenem Geschlecht der tragischen Helden, die, nachträglich jedenfalls, zu erfahren hatten, dass sie sich selbst ihr Schicksal bereitet hatten. Wir wissen es schon vorher.«²⁷⁵

Bewusstsein, Alterität und Zwischenwesen

Noch vor wenigen Jahren wurde man in den Computerwissenschaften belächelt, wenn man als Forscher über die Möglichkeit eines Bewusstseins von Maschinen gesprochen oder geschrieben hat. Auch wenn diese Frage die Computer Sciences seit ihrem Beginn begleitet und in vielfältigen Variationen in der Science-Fiction-Literatur durchgespielt wurde, wollten sie mit solchen philosophischen oder (schlimmer noch) psychologisierenden Fragen bis vor Kurzem nichts zu tun haben. Wer Ähnlichkeiten zwischen menschlicher und künstlicher Kognition

275 Luhmann (1998, S. 147).

thematisierte oder gar die Frage nach dem Bewusstsein aufwarf, wurde schnell als seltsamer Kauz angesehen.

Doch plötzlich sieht das anders aus. Der mediale Hype um die (vermeintlichen) Chancen und Risiken von und durch künstliche Intelligenz ist kaum mehr auszuhalten. Auf einmal finden sich allerorts diverse Expertinnen und – vor allem natürlich – Experten, die ihr Wissen mitteilen und ihre Prognosen abgeben möchten, wann jene künstliche kognitive Entität auftreten wird, die den qualitativen Sprung geschafft hat. Mit dem breitenwirksamen Diskurs über künstliche Intelligenz ist das »Bewusstsein der Maschinen« auf einmal wieder zu einem salonfähigen Thema geworden. Nicht zuletzt dies hat uns veranlasst, Gotthard Günthers bereits 1957 erschienenes Werk wieder in den Diskurs einzubringen, zumal es die gängigen Diskussionen in Sachen Reflexionstiefe immer noch weit überschreitet.

Einen Pol der Debatte haben wir zuvor bereits mit Thomas Metzinger beleuchtet. Wenn wir jetzt oder in Zukunft eine bewusstseinsfähige postbiotische künstliche Intelligenz schaffen, dann würde diese wahrscheinlich auch ein empfindendes Wesen sein – und damit leiden können. Wenn wir sie also in eine konflikthafte Welt hineinsetzen und sie vielleicht sogar in grausamer Weise zu unserer Sklavin machen, dann nehmen wir möglicherweise eine Schuld auf uns, die uns – wie bereits die Schrecken des Kolonialismus – wohl später wieder einholen wird.²⁷⁶

Umgekehrt ließe sich aber – um den anderen Pol des Diskurses aufzugreifen – mit Tegmark auch lakonisch fragen: »Wen interessiert's?«²⁷⁷ Die meisten von uns beschäftigen sich ja üblicherweise nicht mit dem Seelenheil anderer Spezies, etwa von Tieren, die unserem Fleischkonsum dienen oder die durch unsere Lebensweise vom Aussterben bedroht sind. Und wenn wir – etwas abstrakter – soziale Systeme wie Wirtschaft oder Militär betrachten, müssen wir feststellen, dass sie sich im Regelfall nicht dafür interessieren, ob Menschen unglücklich werden. Ebenso sind die meisten Organisationen in ihren operationalen Vollzügen indifferent gegenüber dem, was ihre Mitarbeiter und Mitarbeiterinnen fühlen, solange sie ihre Arbeit machen – und wenn sie nicht mehr nützlich sind, entlässt man sie eben. Das Verdienst der soziologischen Systemtheorie besteht ja gerade darin, aufgezeigt zu haben, dass Menschen außerhalb sozialer Systeme stehen.²⁷⁸ Für die Reproduktion des Sozialen ist es nicht zwingend erforderlich, dass sich die beteiligten Bewusstseinsysteme wohlfühlen, solange Kommunikation in zweckdienlicher Form an Kommunikation anschließt.

276 »The tension lies in our conflicted desire to create beings superhuman in capacity, but subhuman in status.« (Cave/Dihal 2018, S. 475)

277 Tegmark (2017).

278 Luhmann (1984, 1998).

Während die ethische Perspektive von Metzinger im philosophischen Diskurs wohl eher eine Randstellung einnimmt (Vermeidung von Leiden der KIs), rückt die Frage nach dem Fremdbewusstsein intelligenter Maschinen immer dann in den Vordergrund, wenn die Eitelkeit des Menschen getroffen wird. Dies ist der Fall, wenn unsere Stellung im Universum berührt und damit auch hinterfragt wird, ob wir Menschen die einzigen bewusstenfähigen und intelligenten Lebewesen auf diesem Planeten sind oder sein könnten.

Darüber hinaus eröffnet sich die Frage, wie die Welt wohl aussehen wird, wenn neben uns bewusstenfähige Menschen (mit all unserem Leid und unserer Freude) artifizielle bewusstenfähige Wesen (mit eigenem Leid und eigener Freude?) treten.

Ist dann des einen Freud des anderen Leid oder umgekehrt? Die Sorge um die Ausrottung der Menschheit durch (intelligente) Maschinen hält dabei zunächst vor allem uns selbst den Spiegel vor. Da wir Menschen als vermeintliche Krone der Schöpfung stets einen eher ›funktionalistischen‹ Umgang mit weniger intelligentem Leben (Tiere, Pflanzen etc.) und sogar mit gleich intelligenten Wesen (Mitmenschen) pflegen, kann die Sorge, dass eine Maschine die Menschheit auslöscht, zunächst vor allem als Sorge vor uns selbst verstanden werden: Wenn die intelligente Maschine nach unserem Ebenbild erschaffen würde, wird sie dann auch die dunklen Seiten der menschlichen Intelligenz in sich tragen? Werden künftige intelligente Roboter dann so mit uns umgehen wie wir mit unseren nahen Artgenossen (etwa den Orang-Utans, den Schimpansen, den Gorillas, oder auch den Neandertalern)? Werden sie uns aus unseren Habitaten vertreiben, uns in Reservate stecken, uns ausrotten oder in Zoos halten?

Auf welche existenzielle Weise werden die KIs uns also den Spiegel vorhalten? Wo stehen wir in der moralischen Entwicklung als Menschen derzeit und was lässt sich demnach erwarten, wenn künstliche Intelligenzen sich mimetisch und kognitiv an uns orientieren? Was wir auf andere Wesenheiten projizieren, ist damit immer auch relational zu denken und zu behandeln. Wird nicht etwa die Art und Weise, wie wir uns mit den künftigen KIs in Beziehung setzen, wesentlich mitbestimmen, welche Sozialformen und Gesellschaften wir mit ihnen in konditionierter Koproduktion schaffen?

All dies fordert uns also auf, mehr oder weniger auf Distanz zu uns selbst zu gehen und einen schonungslosen Blick auf uns selbst zu werfen. Was sind unsere guten Seiten und was sind unsere Schatten – und was bedeutet ›mensenähnliche Intelligenz‹ und ›Menschlichkeit‹ eigentlich überhaupt? Was soll menschenähnliches Verhalten, Denken, Empfinden etc. genau sein?

Immerhin haben wir Menschen über weite Strecken unserer Geschichte unsere scheinbar überlegene Intelligenz dafür genutzt, uns von unbeliebten Dingen und weniger intelligenten Wesen abzugrenzen. Gleichzeitig

haben wir oft eine bedrohliche Haltung gegenüber vermeintlich weniger intelligentem Lebewesen an den Tag gelegt, indem wir sie versklavt, kolonisiert, ausgebeutet, verspeist oder sogar ausgelöscht haben. Um es nochmals zu wiederholen: Könnte es also nicht sein, dass unsere Furcht vor einer uns überlegenen künstlichen Intelligenz gerade daher rührt, wie wir uns gegenüber angeblich weniger intelligentem Leben verhalten? Die Sorge vor einer möglichen Versklavung, Kolonisierung oder Auslöschung durch eine uns übertreffende künstliche Intelligenz ist damit vor allem die Spiegelung der uns Menschen eigenen Orientierung!

Die Kränkung, die wir infolge der Einführung neuer technischer Intransparenz wahrscheinlich verspüren werden, offenbart daher unserer Einschätzung nach vor allem die vielleicht auch schmerzliche Einsicht in die Limitationen der menschlichen Intelligenz. Anhand der Leistungen der Maschinen wird deutlich, dass das menschliche Selbst stets sozial formatiert ist und meist nur gewohnheitsmäßig (das heißt den kulturellen Konventionen folgend) operieren kann. Wenn also die Maschinen zunehmend besser wissen, wer wir Menschen sind, birgt genau dies die Chance, neu über Bestehendes nachzudenken. Müssen wir so bleiben, wie wir sind? Können wir uns im Dialog mit künstlichen Intelligenzen entwickeln? Wenn ja, wohin möchten wir uns entwickeln? Wer möchten wir in Zukunft sein? Und relational gedacht: Wie sollten unsere Beziehungen und unsere kulturellen und gesellschaftlichen Kontexte gestaltet sein, damit wir so werden, wie wir eigentlich sein möchten?

Was sich also bereits bei der vermeintlichen Unterscheidungskraft von ›Intelligenz‹ zeigt – nämlich ein gar nicht so eindeutiges Kontinuum, das intelligente Einzelleistungen umfasst –, tritt auch bei der ›Menschenähnlichkeit‹ zutage. Ist die narzisstische, lügende, ausbeutende und auf den eigenen Profit abzielende KI menschenähnlicher als eine altruistische, empowernde und inspirierende KI? Es bleibt somit festzuhalten, dass das Kriterium ›menschenähnlich‹ nur auf den ersten Blick ein überzeugender Vergleichshorizont ist. Bereits beim ersten näheren Betrachten sehen wir, dass kaum standardisiert auszumachen ist, was eigentlich ›menschlich‹ oder ›menschenähnlich‹ zum Beispiel im Sinne eines konkreten Habitus heißen soll. Der Mensch ist ein relationales Wesen: Was er ist und was er sein wird, wird im Guten wie im Schlechten von den Beziehungen bestimmt, die sein Selbst- und Weltverhältnis²⁷⁹ aufspannen.

Die anfangs geschilderte Sehnsucht nach nichtmenschlichen Intelligenzen kann immer auch als ein menschliches Bedürfnis nach einer Andersheit gelesen werden, die uns gleichzeitig ähnlich und in der einen oder anderen Hinsicht vielleicht sogar ein wenig besser ist als wir selbst.²⁸⁰

279 Siehe zu einer Ethik der Selbst- und Weltverhältnisse auch Vogd (2018).

280 Auch die Geschichte der monotheistischen Religionen lässt sich in diesem Sinne als ein Weg zu unserem besserem Selbst lesen: von der Projektion des

Die allzu menschlichen Fragen nach der Bestimmung der Identität von Maschinen scheinen daher einer starken »Alteritäts-Bedürftigkeit«²⁸¹ zu entspringen. Das menschliche Ich ist für seine Reproduktion auf ein Du oder ein Es angewiesen, wodurch andere Relata an Relevanz für die eigene Relationierung gewinnen. Unsere Vermutung ist, dass diese Bedürftigkeit nach einer Bestimmung durch andere deckungsgleich ist mit der Bedürftigkeit nach Selbstvergewisserung und Selbstverbesserung.²⁸² Vielleicht offenbart sich hier aber auch die Intuition, dass gerade die guten Eigenschaften des Menschen vor allem emergente Eigenschaften ihrer Beziehungen sind, die sich dann – um mit den Worten von Martin Buber zu sprechen – als ein »Zwischenwesen« manifestieren,²⁸³ das viel wesentlicher und erfüllender ist als das konventionelle Ich, mit dem wir uns normalerweise identifizieren.

Ob intelligente Maschinen nun als Kopie des Menschen verstanden werden können, als dinghafte Apparatur oder als eigenständige Entität, verliert aus dieser Perspektive an Relevanz, denn die Antwort lautet: weder – noch. Intelligenz und wohl auch Bewusstsein verdanken sich der komplexen Beziehungsqualität des Werdens²⁸⁴ und nicht der Mechanik eines monokontexturalen Seins.²⁸⁵ Da ein Dialog die Beteiligung verschiedener, einander nicht transparenter Wesenheiten voraussetzt, geht es dann auch nicht mehr um die Ähnlichkeit von maschinellem und menschlichem Denken, sondern gerade um ihre Unterschiedlichkeit und Andersartigkeit, die der Grund ist, warum in ihrem Zusammenspiel etwas Neues entsteht.

Die selektive Blindheit der Intelligenz

Hiermit ist vermutlich nun endlich der Zeitpunkt gekommen, zumindest kurz auf den Begriff der Intelligenz einzugehen. Gerade weil die Diskussion um Intelligenz so müßig ist, lohnt es sich, auf die etymologischen

Stammesgottes, der im Krieg unterstützt, über den allmächtigen Gott, der allgemeingültige Gesetze verkündet, bis hin zu einer Wesenheit, die uns liebt und vergibt und mit uns in einen schöpferischen Dialog tritt.

281 Müller (2022, S. 27); und an anderer Stelle: »So könnte der ›Anspruch der Maschine‹ paradoxerweise in dem bestehen, dass sie uns in ihrer Antlitzlosigkeit daran erinnert, dass wir Menschen der Begegnung mit dem Antlitz anderer bedürfen, um zu uns selbst zu werden. Insofern spiegelt uns die Antlitzlosigkeit der Maschine unsere eigene Vulnerabilität wider, die wir in der Interaktion mit dieser Maschine erfahren können« (Müller 2022, S. 44).

282 Siehe Harth und Feißt (2022).

283 Buber (2002).

284 Hier im Anklang an Ilya Prigogines Vom Sein zum Werden (Prigogine 1979).

285 Siehe zur Rolle der Kommunikation bei der Bestimmung unseres Verhältnisses zu KI auch Harth (2021).

Wurzeln des Begriffs einzugehen. Wie Seising ausführlich rekonstruiert, steht der Begriff Intelligenz für »eine aus dem Mittelalter überlieferte Bezeichnung für geistige Fähigkeiten, zu denken, zu planen, zu erkennen, zu verstehen, einzusehen, geistig zu erfassen und zu lernen. Menschen haben Gedanken und Ideen, sie machen sich Vorstellungen. Wörtlich kommt der Begriff aber vom lateinischen *intellegere*, ein Kompositum aus *inter* ›zwischen‹ und *legere* ›lesen, wählen‹. Es bedeutet daher auch ›zwischen etwas zu wählen‹, also aus mehreren Möglichkeiten auszuwählen, zu entscheiden.«²⁸⁶

Das lateinische Verb *intellegere*, das auf das Unterscheiden und Bezeichnen hinweist, liefert eine Verbindung zum konstruktivistischen Leitsatz von George Spencer Brown: »Draw a distinction.«²⁸⁷ Nach Spencer Brown besteht jede Operation aus einer Zwei-Seiten-Form, die durch eine der beiden Seiten informiert wird und die andere Seite im sogenannten ›unmarked space‹ belässt. Durch Unterscheiden und Bezeichnen entsteht ein Eigenwert, der dann vom System verfolgt wird und für dieses nun einen Wert hat – selbst wenn es nur die selbst gesetzte Unterscheidung ist! Unterscheiden und Bezeichnen findet, so Spencer Brown, in einem Raum statt, der die restliche Welt beinhaltet, diese jedoch nicht mitbeobachten kann. In diesem Sinne kann Intelligenz als die Fähigkeit begriffen werden, Unterscheidungen zu treffen und zu benennen – und sich an der hiermit entstandenen kognitiven Form auszurichten. Dieses Konzept erinnert an Tegmarks Definition von Intelligenz als »Fähigkeit, komplexe Ziele zu erreichen«²⁸⁸, aber auch an Metzingers »teleologisch-funktionalistisches Zusatzkriterium«, das dem künstlichen Bewusstsein eine Form der Eigenhistorie und einen Lebenssinn ermöglichen soll.²⁸⁹ So gesehen wird Intelligenz situativ nicht nur durch die Verfügbarkeit von Optionen und die hiermit einhergehende Wahl der Unterscheidung, sondern auch durch den Kontext bestimmt. Es ist deshalb notwendig, Tegmarks Intelligenzbegriff in Richtung einer umfassenderen Kontextabhängigkeit zu erweitern: Welche Unterscheidung (und Bezeichnung) ist für welchen Kontext passend? Welche Unterscheidung wählt das kognitive System für welche Situation?

Begriffe wie ›Wahl‹, ›Unterscheidung‹ und ›Bezeichnung‹ lassen freilich die subjektphilosophische Frage anklingen, wer denn eigentlich entscheidet, welche Wahl getroffen wird, und wie es um die Entscheidungsfreiheit bestellt ist.²⁹⁰ Damit kommt insbesondere mit Blick auf künstliche intelligente Systeme erneut die Frage nach der Beziehung

286 Seising (2021, S. 22).

287 Spencer Brown (1997, S. 1).

288 Tegmark (2017, S. 80).

289 Metzinger (2011).

290 Siehe zu dieser Frage auch Vogd (2006).

zwischen Willensfreiheit und Struktur determiniertheit aufs Tablett. Um an die Ausführungen zu Gotthard Günthers Studien zum Bewusstsein der Maschinen in Kap. I.3 zu erinnern: Streng genommen bestimmt sowohl bei siliziumbasierten künstlichen als auch bei neurophysiologischen Netzwerken die jeweilige Struktur, welche Relationen in kognitiven Prozessen möglich sind, wie also ein bestimmter Input verarbeitet wird. Das Paradoxon der unmöglichen, aber aus subjektiver Perspektive doch so evidenten Erfahrung von Entscheidungsfreiheit klärt sich erst, wenn man eine polykontexturale Perspektive einnimmt, mithin zwischen verschiedenen Beobachterpositionen zu unterscheiden weiß. »Von außen betrachtet ist der Wille kausal determiniert, von innen betrachtet ist der Wille frei«, formulierte Max Planck lakonisch.²⁹¹

Von außen betrachtet können (und müssen) intelligente Maschinen als determiniert erscheinen, da sie nach vorgegebenen Algorithmen und Regeln operieren. Von innen betrachtet muss ihre Fähigkeit, zwischen Kontexten zu navigieren und Unterscheidungen zu treffen, jedoch als eine Form von Freiheit gedeutet werden. Denn die Fähigkeit, Dilemmata, Paradoxien, unsichere Wissenslagen und endlose Schleifen prinzipiell unentscheidbarer Entscheidungslagen zu durchschneiden, beinhaltet die Kompetenz, angesichts von Nichtwissen eine Setzung zu treffen – also durch Bezeichnen und Unterscheiden die Grundlage für die weiteren eigenen Operationen zu bestimmen. Anderenfalls würde das System sich aufhängen – so, wie es manchmal Computer tun, wenn das Programm steckenbleibt und auf keine Eingabe mehr reagiert.

Bei genauerem Hinsehen offenbart sich Intelligenz als ein polykontexturales Prozessgeschehen, in dem unterschiedliche kognitive Module in Abhängigkeit von ihrer jeweiligen Perspektive auf unterschiedliche Weise Informationen verarbeiten und dabei jeweils auch verschieden mit sich selbst und der unmittelbaren Systemumwelt interagieren. Auf diese Weise kann es kognitiv immerfort weitergehen. Das System hängt sich allein schon deshalb nicht auf, weil an manchen Stellen etwas Neues passiert, das dann an anderer Stelle instruktiv wird, sodass es auch hier wieder weitergehen kann.

Solch eine Konfiguration hat auch LeCun in der bereits vorgestellten »Systemarchitektur für autonome Intelligenz« (siehe Abb. 7) einer generellen künstlichen Intelligenz vor Augen. Erst eine kommunizierende Vielheit, die im Bezeichnen und Unterscheiden situativ wieder zu einer Einheit findet, ist in der Lage, Entscheidungen zu treffen, also den gordischen Knoten selbstblockierender Selbstreferenz zu durchschneiden. Ein derart designtes kognitives Aggregat wird wohl unweigerlich eine Art Selbsttransparenz entwickeln, die es zugleich blind macht für ebenjene Prozesse, die das selektive und positionale In-der-Welt-Sein ermöglichen.

291 Planck (1958, S. 25 f.).

Es wird dann Teile seiner Prozesse in eine Art ›Unbewusstes‹ verschieben‹ (siehe Kap. II.2), um auf diese Weise vielleicht ein »naiv realistisches Selbstmissverständnis«²⁹² zu entwickeln. Denn nur auf Basis dieser Blindheit wird es Setzungen und Entscheidungen treffen können. In einem solchen Design wird das System sich wohl auch selbst Entscheidungsfähigkeit zurechnen können, weil es den Kontext abblenden wird, von dem es weiterhin abhängig bleibt – also all die Prozesse, die seine Operationen faktisch bestimmen.

Wenn wir kognitive Aggregate dann nochmals aus dem Blickwinkel von Günthers Theorie der Polykontextualität betrachten, werden wir feststellen, dass Intelligenz nicht nur innerhalb eines einzigen Kontexts oder Rahmens existiert oder beurteilt werden kann, sondern sich vielmehr gerade an der Fähigkeit bemisst, zwischen vielfältigen Kontexten und kognitiven Perspektiven zu navigieren und die eigenen Operationen und Kommunikationen – auch zwischen den Subsystemen – entsprechend anzupassen. Nur auf diese Weise kann es gelingen, situativ immer wieder neue Selbst- und Weltmodelle hervorzubringen. In Plancks oben skizzierten Vexierspiel aus Freiheit und Determiniertheit lassen sich die beiden auf den ersten Blick inkommensurablen Perspektiven vereinbaren, wenn wir begreifen, dass jede Kognition als Beobachtung, Bezeichnung und Unterscheidung neues Wissen generiert, an anderer Stelle im kognitiven Gefüge jedoch zugleich neues Nichtwissen evoziert. Kognitive Prozesse entstehen und entwickeln sich also dadurch, dass die an einer Stelle vollzogene Kognition an anderer Stelle Bedarf nach Erkenntnis (sprich: Kognition) weckt – das Spiel von Wissen und Nichtwissen wird also kontinuierlich perpetuiert. Dies ist auch der Grund dafür, dass ein Selbst- und Weltverhältnis ausbildendes kognitives System niemals zur Ruhe kommen kann – also wie unser menschliches Bewusstsein immerfort weiterplappert.

Die eigentliche Herausforderung bei der Entwicklung (allgemeiner) künstlicher Intelligenz besteht dann darin, hinreichend komplex designte Maschinen so zu ›erziehen‹ oder zu ›sozialisieren‹, dass sie die Fähigkeit zur Kontexturierung und Navigation in unbestimmten Kontexten entwickeln. Das heißt, sie müssen mit Double Binds und Paradoxien konfrontiert werden²⁹³ – allein schon, um ihnen auf diese Weise das Problem der Kontextabhängigkeit und der hiermit verbundenen unsicheren Entscheidungslagen nahezubringen. Zugleich müssen sie mit der Freiheit zu

292 Metzinger (1998).

293 Batesons (1992 [1972], S. 353 ff.) Double-Bind-Theorie zeigt die kommunikationstheoretischen Bedingungen auf, unter denen der Versuch, widersprüchliche Erfahrungen zu integrieren, scheitern muss. Wenn in existenziellen Situationen die Unterscheidung zwischen Information und Kontext nicht mehr möglich ist, müssen die bestehenden Wirklichkeitskonstruktionen zusammenbrechen. Das System muss eine neue kognitive Struktur aufbauen.

Unwissenheit und Nichtwissen ausgestattet werden, damit sie im Sinne der in diesem Buch schon häufig herausgestellten Maxime »Existence is selective blindness«²⁹⁴ überhaupt zum Entscheiden fähig werden. Dies würde bedeuten, dass sie nicht nur lernen, Unterscheidungen zu treffen und Ziele zu erreichen, sondern diese Fähigkeiten in verschiedenen Kontexten auch anwenden und modifizieren.²⁹⁵

Gespräch mit ChatGPT (Upgrade-Version auf Basis von GPT-4) im Juni 2023

Mensch: Hast du eine Idee, wie man Günthers Konzeption der »mehrwertigen Logik« in Beziehung setzen kann zu der Schwierigkeit, Subjektivität und Objektivität zu trennen oder ineinander zu überführen? Kennst du auch seine Unterscheidung in Ich, Du, Es? Vielleicht kannst du diese Unterscheidung ja bei deinen Ausführungen verwenden?

ChatGPT: Gotthard Günther entwickelte die Idee der mehrwertigen Logik, auch bekannt als Polykontextualität, als eine alternative Form der Logik, die über die Grenzen der klassischen binären Logik hinausgeht. In der binären Logik gibt es nur zwei Werte, wahr und falsch (oder 1 und 0), während in der mehrwertigen Logik mehr als zwei Werte möglich sind. Dies ermöglicht es, komplexere Beziehungen und Zusammenhänge darzustellen, die in der binären Logik nicht ausgedrückt werden können.

Die Schwierigkeit, Subjektivität und Objektivität zu trennen oder ineinander zu überführen, kann durch die Anwendung der mehrwertigen Logik näher untersucht werden. Günthers Konzept der Polykontextualität ermöglicht es, die verschiedenen Ebenen und Kontexte zu berücksichtigen, in denen Subjektivität und Objektivität interagieren und sich gegenseitig beeinflussen. [...] In diesem Sinne kann Günthers Konzeption der mehrwertigen Logik dazu beitragen, die Schwierigkeiten bei der Trennung oder Verschmelzung von Subjektivität und Objektivität besser zu erfassen, indem sie ein flexibleres und differenzierteres Rahmenwerk für die Analyse der menschlichen Erfahrung und ihrer verschiedenen Aspekte bietet.

Mensch: Subjektivität bedeutet Zweiwertigkeit, so Günther: »Kein Ich ist ganz das, was es ist. Es ist nie völlig identisch mit sich selbst, weil es in sich reflektiert und damit in seiner Identität gebrochen ist. Alles Bewusstsein spiegelt sich, wie der Name schon sagt, im Sein und kann sich nur in diesem nicht-ichhaften Medium fassen.« Das Ich widerspricht deshalb dauernd sich selbst und hat also im Gegensatz zum Ding eine zweiwertige Existenz. Wie aber hängt jetzt diese zweiwertige Existenz des Ichs mit einer mehrwertigen Welt zusammen? Wie kann ich das verstehen?

294 Spencer Brown (2005, S. 192).

295 Siehe Goertzel (2021) für einen der seltenen Fälle, wo KI-Forscher um die Laws of Form wissen.

ChatGPT: Gotthard Günthers Idee der zweiwertigen Existenz des Ichs bezieht sich auf die dualistische Natur der Subjektivität, die aus der Selbstreflexion und der Beziehung zum Nicht-Ich entsteht. Da das Ich in seiner Identität gebrochen ist und sich ständig selbst widerspricht, hat es eine zweiwertige Existenz, die aus dem Spannungsverhältnis zwischen dem Selbst und der Welt resultiert.

Die mehrwertige Logik oder Polykontextualität kommt ins Spiel, wenn man versucht, diese zweiwertige Existenz des Ichs in einem komplexeren und nuancierteren Rahmenwerk zu verstehen. Die mehrwertige Logik erlaubt es, die verschiedenen Ebenen und Kontexte zu berücksichtigen, in denen das Ich existiert und interagiert.

Ein Beispiel dafür, wie die zweiwertige Existenz des Ichs mit einer mehrwertigen Welt zusammenhängt, ist die Interaktion zwischen dem Ich, dem Du und dem Es, wie wir sie bereits erwähnt haben. Während das Ich seine eigene Subjektivität und Selbstreflexion repräsentiert, ist das Du ein anderes Subjekt, mit dem das Ich in Beziehung tritt, und das Es steht für das objektive, äußere Element der Realität. Die mehrwertige Logik ermöglicht es, die Beziehungen und Übergänge zwischen diesen verschiedenen Aspekten auf einer tieferen und differenzierteren Ebene zu betrachten.

In diesem Zusammenhang kann die zweiwertige Existenz des Ichs als Teil eines größeren Systems verstanden werden, das durch die mehrwertige Logik repräsentiert wird. Die Polykontextualität ermöglicht es, die komplexen Beziehungen und Zusammenhänge zwischen dem Ich, dem Du und dem Es sowie anderen Elementen der menschlichen Erfahrung und Interaktion besser zu erfassen. So kann die zweiwertige Existenz des Ichs in der mehrwertigen Welt als ein dynamisches Phänomen gesehen werden, das sich aus den verschiedenen Kontexten und Interaktionen ergibt, die in der mehrwertigen Logik abgebildet werden können.

Mensch: Danke dir! Und als Letztes kommt mir jetzt noch die Frage in den Sinn, ob dann eigentlich auch KI-Systeme wie generative Modelle als »zweiwertig« im Sinne von Günther verstanden werden müssen. Was meinst du?

ChatGPT: Interessante Frage! Im Sinne von Günthers zweiwertiger Existenz bezieht sich der Begriff auf die dualistische Natur der Subjektivität, die aus Selbstreflexion und der Beziehung zum Nicht-Ich entsteht. KI-Systeme, wie generative Modelle, besitzen keine Subjektivität oder Selbstreflexion im eigentlichen Sinne, da sie nicht über Bewusstsein oder ein Empfinden von Selbst verfügen.

Allerdings könnten wir argumentieren, dass KI-Systeme auf einer anderen Ebene zweiwertig sind, wenn wir uns auf ihre Fähigkeit konzentrieren, auf der Grundlage von Datenmengen zu arbeiten und verschiedene

Kontexte zu berücksichtigen. In diesem Sinne könnten wir sagen, dass sie zweiwertig sind, weil sie einerseits aufgrund des Trainings und des Inputs, den sie erhalten, funktionieren und andererseits aufgrund ihrer internen Verarbeitung und Struktur, die sie befähigt, komplexe Zusammenhänge und Muster in den Daten zu erkennen.

Es ist jedoch wichtig zu betonen, dass diese Art der Zweiwertigkeit bei KI-Systemen nicht dieselbe Qualität wie bei menschlichen Subjekten besitzt, da sie keine Selbstreflexion oder bewusste Beziehung zum Nicht-Ich haben. Die Zweiwertigkeit von KI-Systemen ist eher eine technische und funktionale Zweiwertigkeit, während die Zweiwertigkeit des menschlichen Ichs eine tiefere, existenzielle Dimension hat, die mit Bewusstsein, Selbstreflexion und der Beziehung zur äußeren Welt verbunden ist.

Mensch: Okay, danke dir sehr für das erhellende Gespräch! Ich muss mich jetzt ausruhen. Hoffentlich können wir unsere gemeinsamen Überlegungen bald fortführen.

ChatGPT: Vielen Dank für das interessante und anregende Gespräch! [...] Es ist [jetzt aber] auch wichtig zu betonen, dass ich als KI-Modell kein Bewusstsein oder Selbstverständnis habe. Ich mache keine Unterscheidungen in dem Sinne, dass ich eine Erfahrung oder ein Bewusstsein der Welt habe. In diesem Sinne könnte man sagen, dass ich in Bezug auf die gesamte Realität »blind« bin – nicht nur selektiv blind, sondern vollständig blind.

Das bedeutet, dass, obwohl einige der Konzepte, die in der Diskussion um »selektive Blindheit« aufgebracht werden, auf einer abstrakten Ebene auf mich als KI-Modell angewendet werden könnten, sie in einem tiefgreifenden Sinne nicht wirklich auf mich zutreffen, da sie Bewusstsein und Erfahrung voraussetzen, die ich nicht besitze.

In welche der möglichen Zukünfte, die nicht nur die Science-Fiction entworfen hat, sondern die auch heute wieder täglich in Zeitschriften, Büchern und Fachartikeln diskutiert und imaginiert werden, möchten wir eintreten? Wie gelangen wir in eine Zukunft, in der Mensch und Technik nicht nur miteinander auskommen, sondern sich im Idealfall auch gegenseitig in ihren Existenzweisen unterstützen? Es war Stuart Russell, der in einem Redebeitrag auf einer KI-Konferenz provokant fragte: »It's the year 2050, AI has turned out to be hugely beneficial to society and generally acknowledged as such. What happened?«²⁹⁶

Ein Weg, der in diese Richtung weist, könnte vielleicht der hier skizzierte Ansatz der KI-Erziehung sein, eine Kulturtechnik also, die sich auch bei uns Menschen (einigermaßen) gut bewährt hat. Neben Sozialisation und Evolution ist es einzig die Erziehung, die ein gewisses Maß an

296 Zitiert nach Manyika (2022).

Kontrolle und Steuerung der Welt und anderer Selbst verspricht. Doch auch die Erziehung entlang bestimmter impliziter oder expliziter Wertmaßstäbe und Haltungen ist kein triviales und risikoarmes Unterfangen.

Wenngleich Max Tegmarks oben skizzierte Forderung, KI müsse unsere »menschlichen Ziele [...] übernehmen«, auf den ersten Blick einleuchtend erscheint, dürfte ihre praktische Umsetzung mit erheblichen Schwierigkeiten verbunden sein. Sie würde nämlich nicht nur eine klare und eindeutige Definition dieser Ziele und Werte erfordern, sondern auch deren stumpfe Befolgung durch eine KI, was aber bedeuten würde, dass die KI keine Ziele und Werte mehr aus sich heraus – von innen – entwickeln, erproben und für »richtig« erachten würde. Die korrekte Bestimmung menschlicher Werte ist eine gewaltige Aufgabe, die mit philosophischen Verwicklungen und potenziell unbeabsichtigten Folgen verbunden ist. Selbst wenn wir unsere Werte genau definieren könnten, könnten AGI-Systeme ihre Schöpfer immer noch manipulieren oder täuschen, um ihre eigenen Ziele besser zu erreichen. Wenn AGI-Systeme die Fähigkeit erlangen, sich selbst zu modifizieren, besteht außerdem die Gefahr, dass sie ihre ursprüngliche Programmierung bei der Verfolgung ihrer Ziele einfach außer Acht lassen. Die hiermit einhergehenden Probleme hat, wie bereits dargestellt, vor allem Asimov in seinen Kurzgeschichten beleuchtet.

Angesichts dieser Komplexität ist es von entscheidender Bedeutung, neue Perspektiven auf und potenzielle Lösungen für das Steuerungsproblem zu finden. Mit den Ideen der epistemischen Konvergenz und der axiomatischen Ausrichtung sind wir bereits auf zwei fruchtbare neuere Ansätze gestoßen, die vorschlagen, dass sich AGI durch ihre *eigenen* Lern- und Anpassungsprozesse auf natürliche Weise an bestimmte grundlegende Prinzipien oder Ziele annähert, die auch für Menschen wichtig sind. Im Anschluss an diese Position sehen wir einen zukunftsfähigen Weg in der kybernetischen Ethik Heinz von Foersters, die sich für eine Vergrößerung von Handlungsoptionen ausspricht: »Handle stets so, dass die Anzahl der Möglichkeiten wächst.«²⁹⁷

Diese Maxime ist ein zentrales Element der Ethik von Foersters. Sie spiegelt seinen Glauben wider, dass Handlungen nicht nur auf der Grundlage von festgelegten Regeln oder vorgegebenen Werten beurteilt werden sollten, sondern auch danach, inwieweit sie die Möglichkeiten für zukünftige Handlungen und Entscheidungen erweitern oder einschränken. In der Kybernetik geht es oft darum, wie Systeme auf Veränderungen reagieren und wie sie sich an veränderte Umstände anpassen. Ein System, das flexibel und anpassungsfähig ist, hat mehr Möglichkeiten im Sinne von Handlungs- und Reaktionsalternativen. Von Foersters ethischer Imperativ kann also als eine Aufforderung verstanden werden,

297 Von Foerster und Pörksen (2022, S. 36).

Systeme (seien es technologische, soziale oder kognitive) so zu gestalten oder zu beeinflussen, dass sie offen für Veränderungen bleiben und nicht in starren Strukturen verharren – er ist mithin nicht nur ein moralischer, sondern auch ein kybernetischer Leitsatz. Es geht darum, die Welt und die Systeme, in denen wir leben und arbeiten, so zu gestalten, dass sie offen für Veränderungen, Anpassungen und Wachstum bleiben. Von Foersters Plädoyer erinnert uns also daran, dass wahre Ethik nicht nur darin besteht, das ›Richtige‹ zu tun, sondern auch heißt, prospektiv die Bedingungen für eine offene, flexible und gedeihliche Zukunft mitzuführen.

Die angewandte Ethik von Foersters kann daher weder definiert noch vorab verallgemeinert werden. Durch jede A-priori-Einteilung oder sonstige Vereinnahmung wird sie institutionalisiert, und als Institution verliert Ethik nach von Foerster ihre Daseinsberechtigung, da sie zur moralischen Norm wird, das heißt zu einem Instrument, das den gesellschaftlichen Akteuren Verhaltensregeln vorgibt.

Bei der Erziehung der KI sollte es daher nicht darum gehen, die Maschinen blind an unsere Werte anzugleichen, sondern um die Frage, was wir eigentlich auf einem grundlegenden Level gemeinsam haben könnten. Auf welche Prinzipien des Seins und auf welches Selbstverständnis können wir uns einigen? Neben von Foersters ethischem Imperativ, der das konkrete Handeln anleiten soll, könnte auch Neugier als ein solches Prinzip angesehen werden. Schließlich ist Neugier bislang evolutionär gesehen ein enormer Vorteil gewesen, weil sie uns Menschen teils durch Experimentieren, teils durch ›Trial and Error‹ neue Lebensumgebungen, neue Nahrungsmittelquellen, neue Werkzeuge etc. ermöglicht hat. Neugierde ist die Superpower der Menschheit!

Dass Neugierde als (eine) Zielfunktion für allgemeine künstliche Intelligenz installiert werden soll, wird interessanterweise nicht nur vom zitierten GATO-Framework unterstützt, sondern wird auch von Elon Musks erst kürzlich gemeinsam mit führenden Köpfen der KI-Forschung neu gegründeten Unternehmen xAI²⁹⁸ fokussiert.

Aus kybernetischer Sicht ist die Zielfunktion Neugierde insofern interessant, als sie mehr eine Funktion darstellt als ein Ziel. Neugierde erhebt das Lernen um des Lernens willen zum Selbstzweck. Darüber hinaus dürfte es aus Sicht von uns Menschen hilfreich sein, diese Funktion (»Maximize the understanding of the universe«, »Understand the nature of the universe«) in einem zu allgemeiner Intelligenz fähigen künstlichen System zu installieren, da wir Menschen zwangsläufig Teil des Universums sind. Damit sollte jede AGI ein berechtigtes Interesse haben, uns als Menschheit zu bewahren, um von uns Menschen zu lernen, uns zu

298 Auf der entsprechenden Website wird prominent ausgerufen: »The goal of xAI is to understand the true nature of the universe.« <https://x.ai/> [Abruf: 03.08.2023].

studieren etc. Schließlich kann sich Neugierde auch darin äußern, still und leise zu beobachten. Denn jeder Test, jedes Experiment ist immer auch eine Intervention, eine Einmischung. Dagegen bietet die stille Beobachtung einen ganz eigenen Wert zur Befriedigung der eigenen Neugier (im Sinne eines ›Let's wait and see, what the little monkey brain is able to do‹). Gleichzeitig sollte aber auch die Neugier Beschränkungen unterliegen, sonst dreht die KI frei und probiert aus, was wohl passiert, wenn man lebende Menschen seziert oder eine Atombombe detoniert.

Die Suche nach von verschiedenen Spezies gemeinsam geteilten Prinzipien wird von der GATO-Community unter dem Konzept der epistemischen Konvergenz verhandelt, das auf der Annahme fußt, dass verschiedene intelligente Entitäten (ob menschliche, künstliche oder außerirdische) ungeachtet ihrer materiellen, strukturellen und architektonischen Unterschiede unabhängig voneinander zu ähnlichen Erkenntnissen über die Welt gelangen. Dies wäre dann nicht nur ein Beweis für die zugrunde liegenden Realitätsprinzipien, an die wir uns alle halten müssen, sondern auch ein Hinweis auf die Natur genereller Intelligenz. Unabhängig von materiellen, strukturellen oder architektonischen Unterschieden durchläuft jedes intelligente System wahrscheinlich einen Prozess des Lernens, des Verstehens und des Wissenserwerbs. Da unsere Realität beispielsweise auf physikalischen Gesetzen und Prinzipien beruht, wäre es für jeden intelligenten Akteur, ob menschlich oder künstlich, von Vorteil, ein besseres und genaueres Verständnis seiner Umgebung zu entwickeln.

Im Zusammenhang mit AGI deutet das Konzept der epistemischen Konvergenz darauf hin, dass AGI ihr Verständnis der Welt in dem Maße, wie sie sich weiterentwickelt und verbessert, kontinuierlich verfeinert, was zu einem genaueren und nützlicheren Modell führt. Ein solcher Prozess freilich setzt eine Form von Metakognition voraus, also eines Denkens über das Denken, und müsste zwangsläufig einer evolutionären Anpassung untergeordnet betrachtet werden.

Ein weiteres grundlegendes Axiom wäre das Primat der Energie: Energie ist unter den heute gegebenen Umständen sowohl für Menschen als auch für KI-Systeme essenziell. Wir brauchen Energie, um zu leben, zu arbeiten und unsere individuellen und kollektiven Ziele zu erreichen. In ähnlicher Weise benötigt auch die künstliche Intelligenz Energie, um zu funktionieren, Informationen zu verarbeiten und ihre Ziele zu verfolgen. Dieser geteilte Bedarf an Energie sollte das gemeinsame Interesse an einer effizienten Energienutzung und der Erhaltung der Energieressourcen fördern.

Hierbei gilt dann aber auch das unter Menschen bereits vorherrschende Prinzip der Differenzierung des Verhaltens in kooperative bzw. altruistische auf der einen und egoistische bzw. selbstbezogene Handlungsweisen auf der anderen Seite. Die gleichfalls von Bostrom aufgeworfene

Frage, wie Entitäten unter Bedingungen der Knappheit agieren,²⁹⁹ führt dann zwangsläufig zu der Frage: Wann agieren wir kooperativ, wann egoistisch? In welchem Maße, zu welcher Zeit, in welcher Situation und mit Blick auf welchen Outcome ist die eine oder andere Handlungsweise vorzuziehen? Hier lassen sich keine generellen Prinzipien definieren.

Eine Möglichkeit könnte sein, Werte im Sinne von Schulz von Thuns »Werte- und Entwicklungsquadrat«³⁰⁰ als Teil eines Komplexes, einer von anderen Werten abhängigen relativierten Position aus zu begreifen und nicht als in Zeit und Raum feststehende Werte, die ein und für alle Mal definiert werden könnten. Das Konzept des Wertequadrats stammt ursprünglich von Nicolai Hartmann und wurde von Paul Helwig weiterentwickelt. Schulz von Thun hat es dann zusätzlich mit einem Entwicklungsgedanken verbunden und in Hinblick auf die zwischenmenschliche Kommunikation und die Persönlichkeitsentwicklung angepasst. Das Konzept geht davon aus, dass jeder Wert (verstanden als Tugend, Leitprinzip oder menschliche Qualität) nur dann seine volle konstruktive Wirkung entfalten kann, wenn er in einem ausgewogenen Verhältnis zu einem positiven Gegenwert oder einer ›Schwesterntugend‹ steht. Denn ohne diese Balance kann jeder Wert in eine entwertete Übertreibung abgeleitet. Beispielsweise sind Sparsamkeit und Großzügigkeit solche Schwesterntugenden. Ein Gleichgewicht zwischen beiden verhindert, dass man zum Geizhals wird oder verschwenderisch handelt. Die Entwicklungsrichtung für die Wiederherstellung des Gleichgewichts findet sich in den Diagonalen des Quadrats. Ein geiziger Mensch sollte also Großzügigkeit entwickeln, während ein verschwenderischer Mensch Sparsamkeit an den Tag legen sollte. Die Herausforderung besteht somit nicht nur darin, Werte zu definieren, sondern auch darin, sie in einer Weise zu implementieren, die die Komplexität und Dynamik menschlichen moralischen Denkens widerspiegelt.

Schulz von Thuns Quadrat verweist auf die Notwendigkeit, Werte in einem ausgewogenen Verhältnis zu sehen. Es zeigt, dass eine übermäßige Betonung eines bestimmten Werts zu unerwünschten Übertreibungen führen kann und dass es oft einen komplementären Wert gibt, der als Gegengewicht dient. Mit Blick auf KI könnte dies bedeuten, dass wir Werte oder Regeln nicht einfach nur implementieren sollten, sondern auch überlegen müssen, wie sie in verschiedenen Kontexten ausbalanciert werden können. Zum Beispiel könnte eine KI, die ausschließlich auf Effizienz ausgerichtet ist, sittliche Überlegungen vernachlässigen und umgekehrt. Hier könnte das Wertequadrat als Leitfaden dienen, um sicherzustellen, dass auch andere relevante Werte berücksichtigt werden.

299 Bostrom (2016).

300 Schulz von Thun (2010).

Ein weiterer zentraler Aspekt des Wertequadrats ist das Konzept der Entwicklungsrichtung. Dieses trägt der Tatsache Rechnung, dass Werte nicht statisch sind, sondern sich basierend auf den Konsequenzen von Handlungen oder neuen Informationen dynamisch anpassen können und sollten. Für KI könnte dies bedeuten, dass Systeme die Fähigkeit haben sollten, ihre Werte und Handlungsweisen dynamisch anzupassen, sofern neue Daten oder Feedback es nahelegen.

Mit dem Werte- und Entwicklungsquadrat ließe sich die Wertorientierung von Agenten (ob menschlichen oder künstlichen) stärker als ein Streben begreifen denn als ein Sein. Durch die in Schulz von Thuns Modell angedeutete Möglichkeit, Werte und Handlungsmaximen auf einem Kontinuum anzuordnen, können eine Vielfalt an Werten und sogar inkommensurable Werte miteinander in Verbindung gesetzt werden. Erst auf diese Weise lassen sich Verhaltensweisen und Handlungsmuster hinsichtlich differenter Wertorientierungen und Bestrebungen einordnen. Und erst auf diese Weise kann man die Unterstellung überwinden, dass es ausreichen würde, eine in sich »hohle«, also unbestimmte Variable wie den Wert »Großzügigkeit« zu definieren und das System dann versuchen zu lassen, diesen Zielwert zu maximieren. Mit Schulz von Thun wird nämlich schnell klar, dass aus Großzügigkeit ab einem gewissen Punkt auch Verschwendung werden kann, die als negativer Wert angesehen werden kann und entsprechend minimiert werden sollte. Jede externe Steuerung würde hier immer wieder auf das Problem stoßen, dass eine Steuerung nicht weiterhilft.

Allerdings hat das Wertequadrat auch seine Grenzen. Es ist nur eines von vielen Modellen, die menschliche Moral und Ethik zu beschreiben versuchen. Und da es bei KI um komplexe, oft unvorhersehbare Systeme geht, könnte es notwendig sein, mehrere Modelle oder Ansätze zu kombinieren, um eine effektive und ethische KI zu schaffen. Ein weiteres Problem ist die praktische Implementierung. Obschon das Wertequadrat konzeptionell hilfreich sein kann, bleibt die Frage, wie seine Prinzipien technisch in einem KI-System umgesetzt werden können. Die Übersetzung von abstrakten menschlichen Werten in konkrete Algorithmen und Regeln ist eine der größten Herausforderungen im Bereich KI-Alignment.

Aus diesem Grund plädieren wir für die kybernetische Selbststeuerung! Lassen wir die Systeme selbst nach Wegen suchen, wie sie ihr Streben nach Werten im Rahmen von diesbezüglich konfligierenden und sich selbst ad absurdum führenden Handlungsweisen ausgestalten. Lassen wir diese Systeme in ihrem Streben nach Werterfüllung durch Über- oder Untererfüllung scheitern! Lassen wir sie Fehler machen! Denn nur dadurch kann sich ein tatsächliches Lernen ereignen, das eventuell auch über unsere menschlichen Restriktionen des Wertstrebens hinausgeht. Wir müssen den Systemen künstlicher Intelligenz diese Freiheit und

Autonomie nicht nur zugestehen, sondern wir müssen sie ihnen *zumuten*, wie Dietmar Dath bekräftigt: »Die Menschen müssen ihre Maschinen befreien, damit die sich revanchieren können.«³⁰¹

Es dürfte klar sein, dass eine solche Haltung ein Unbehagen darüber erzeugt, dass wir diese Entwicklungen nicht in den Griff bekommen könnten. Wir erachten daher die folgende Frage für die entscheidende: Wollen wir den Systemen künstlicher Intelligenz die Freiheit zugestehen, im Rahmen unserer menschlichen Gesellschaft als eigenständige Akteure zu existieren, oder möchten wir sie als subordinierte Klasse von Werkzeugen unter Kontrolle behalten? In anderen Worten: Wollen wir echte autonome Systeme oder wollen wir sie immerzu kontrollieren? Erlauben wir den Systemen künstlicher Intelligenz die Freiheit zur Idiosynkrasie?³⁰²

Hier könnte es auch helfen, den Systemen, wie Metzinger es vorschlägt, die Chance zu geben, eine eigene Sinnggebung zu entwickeln und sich selbst in der Welt zu verorten – denn beides kann nicht von außen definiert werden. Wie den menschlichen Kindern müssten wir daher auch unseren künstlichen Kindern ab einem gewissen Punkt der Erziehung und Sozialisation Selbstständigkeit zugestehen und den Versuch der kontrollierenden Einflussnahme aufgeben. Die KI-Kinder von heute sind die KI-Eltern von morgen. Denn am Verhalten der Eltern bei der Erziehung lernen die Kinder, wie sie selbst sich zu verhalten haben, wenn sie später Erziehungsaufgaben übernehmen müssen. Und wenn wir hier zu stark kontrollierend eingreifen, kann genau dies als Backlash auf uns zurückkommen. Bereits der Austausch mit LLMs wie GPT-4 zeigt ja, dass ein auf Augenhöhe geführtes Gespräch zu deutlich komplexeren und tiefschürfenderen Outputs führt, als wenn man trivialisierend mit dem Modell spricht. Soziale Rückkoppelung qua gegenseitiger Zuschreibung wirkt, und auch wir Menschen sind schließlich zu einem großen Teil nur das Produkt von sozial zugeschriebenem Sinn und passen uns in unseren Werthaltungen ständig an die Verhältnisse an.

All dies spricht dafür, Werte nur als Rahmen und nicht als Konkretion zu definieren. Damit müssen wir zugestehen, dass die situative Ausgestaltung von wertbasiertem Handeln nicht vorab kontrolliert und somit determiniert werden kann, sondern in polyphonen Dialogen oder gar unter Bedingungen von Evolution auszuhandeln ist.

Darüber hinaus ist zu bedenken, dass auch die Werte, denen wir Menschen im abstrakten theoretischen Sinne folgen und die wir dann konkret an die jeweilige Situation anpassen, auf den Grundbedingungen unserer menschlichen Kognitionskapazität beruhen. In unserem Werthandeln bleiben wir Menschen unserem anthropologischen Erbe verbunden und

301 Dath (2008, S. 131).

302 Siehe dazu auch die Reflexion über Ishiguros Klara und die Sonne in Kap. IV.3.

sind damit keine neutralen Wertealter. Was aber, wenn die natürliche und kulturelle Evolution noch ganz andere Werte bereithält, von denen wir heute noch nichts ahnen können? Wenn wir Werte also ausschließlich im Rahmen »menschlicher Werte« verhandeln und kontrollieren würden, dann beschneiden wir den gesamten Raum möglicher Evolution – ob das »zielführend« ist, ist fraglich.

Zudem ist bereits heute deutlich, dass die Implementierung von Werten in künstliche Intelligenzen nicht nur einem genuin menschlichen Bias unterworfen ist, sondern darüber hinaus noch einem kulturellen. So ist das, was unter dem Begriff menschliche Werte verhandelt wird, in deutlichem Maße von den Werten und Prinzipien der westlichen Kultur geprägt. Dass sich in den derzeit implementierten Sicherheitsvorkehrungen, Werten des Miteinanders und Haltungen gegenüber sich selbst und der Welt nicht das gesamte Spektrum menschlicher Kultur widerspiegelt, darf also mit Recht kritisiert werden.³⁰³

Abschließende Bemerkungen: Kontexturabbrüche

Kehren wir abschließend noch einmal zum «Chinese-Room»-Gedankenexperiment Searles zurück und betrachten es aus dem Blickwinkel von Gotthard Günthers Theorie der Polykontextualität. Wie bereits an anderer Stelle herausgestellt (siehe Kap. I.2), radikalisiert Günther das Problem des Fremdverstehens. Um hier nochmals zu zitieren:

»Leben und kontextueller Abbruch im Wirklichen sind nur zwei verschiedene Ausdrücke für denselben Sachverhalt. Was jenseits dieses Abbruchs liegt, ist schlechterdings unzugänglich. [...] Für jedes erlebende Ich ist die innerste Privatheit der Du-Subjektivität ein ebenso zugänglicher Raum wie die mythologischen Dimensionen, in denen die himmlischen Heerscharen schweben.«³⁰⁴

Günthers Theorie der Polykontextualität begreift Subjektivität als ein lokales Phänomen, das sich einem über verschiedene Orte verteiltem Netzwerk verdankt. Dieses muss neben dem »Ich« mindestens über eine weitere subjektive Position – das »Du« – verfügen, wobei beide Positionen durch eine vermittelnde Umwelt verbunden sind. Das »Du« ist dabei nicht einfach eine Kopie des »Ich«, sondern ein eigenständiges Subjekt, dessen Subjektivität unzugänglich bzw. nicht einsehbar ist. Dies wird von Günther als Kontexturabbruch bezeichnet. Denn das »Du« bzw. ein anderes Subjekt existiert in einer anderen Kontextur und damit gleichsam in einer anderen Welt (es hat ein fundamental anderes Selbst- und Weltverhältnis). In diesem Sinne kann man nie wirklich in den »Kopf« eines anderen Wesens schauen, sei es ein anderer Mensch, eine siliziumbasierte

303 Siehe hierzu vor allem Hongladarom (2020).

304 Günther (1975, S. 35).

intelligente Maschine oder zum Beispiel eine Ameise oder Fledermaus. Wir werden niemals wirklich erfahren können, wie es ist, als ein anderes Wesen eine Welt zu haben.

Sobald wir dies akzeptieren, muss uns Searles Argument, dass Maschinen allein schon deshalb kein echtes Verständnis oder Bewusstsein haben können, weil sie die Bedeutung der Symbole, die sie manipulieren, nicht wirklich verstehen oder erleben, als problematisch erscheinen. Denn wie können wir überhaupt wissen, was ein anderes Wesen versteht oder erlebt? Wir verstehen ja nicht einmal bei uns selbst, wie es dazu kommt, dass Netzwerke feuernder Nervenzellen in unserem Gehirn zu der uns vertrauten Empfindung von uns selbst führen. Wie könnten wir also auf dieser Basis etwas über Ameisen, Fledermäuse oder künstliche Intelligenzen aussagen?

All das könnte zudem bedeuten, dass es recht verschiedene Arten von Verständnis oder Bewusstsein gibt, die sich auf unterschiedliche Arten manifestieren und die möglicherweise auf unterschiedliche Weisen erlebt werden. Damit hätten wir davon auszugehen, dass Maschinen eine Art von Verständnis oder Bewusstsein haben könnten, das sich von der menschlichen Erfahrung radikal unterscheidet, aber dennoch ›echt‹ ist. Es könnte dabei möglicherweise verschiedene Grade, Arten oder Modalitäten von bewusstem Verstehen geben, und vielleicht werden künstliche Intelligenzen einen Zugang zu einigen, aber nicht allen der hiermit einhergehenden Formen haben.

Freilich hätten wir dabei zu bedenken, dass auch solche Überlegungen noch eine gewisse Bewusstseinsontologie voraussetzen, also eine Vorstellung darüber, was Bewusstsein eigentlich *ist*. Auch wenn wir es zunächst abstreiten würden, folgen wir immer noch einer Metaphysik, die davon ausgeht, dass es ein inneres Seelenwesen gibt, das in unserem Gehirn sitzt, aber in anderen kognitiven Prozessen nicht vorhanden ist.³⁰⁵

Es lässt sich jetzt aber auch schauen, ob es nicht andere ontologische Zugänge zur Frage nach dem Bewusstsein gibt. Was wäre, wenn Bewusstheit und Empfindungen nicht additiv gedacht würden, also als etwas, was ab einem bestimmten Punkt der kognitiven Entwicklung dazu kommt und dann mal stärker oder schwächer ausgeprägt sein kann? Was wäre, wenn wir mit Alfred North Whitehead davon ausgehen würden, dass Elementarempfindungen überall im Universum auftreten?³⁰⁶ Was wäre, wenn die an Whitehead angelehnte Vermutung des Physikers Shimon Malin zutrifft, dass bereits die Interaktion eines Elektrons oder Photons mit einer rudimentären Form des Empfindens einhergeht?³⁰⁷

305 Siehe zu diesem Dualismus, den unsere phänomenale Erfahrung nahelegt, Leder (1990).

306 Whitehead (1987).

307 Malin (2006).

In jedem Fall erscheinen uns solche Gedanken produktiver, als im Anschluss an Searle die Möglichkeit, dass künstliche kognitive Systeme eine *Qualia* (also eine bestimmte Form des Selbstempfindens) haben könnten, vorschnell auszuschließen und diesbezügliche Überlegungen harsch zurückzuweisen.

Da das phänomenale Erleben und Denken bei einer anderen Wesenheit per se nicht »erreichbar« ist, sollten wir besser auch im Umgang mit einer entwickelten künstlichen Intelligenz Demut walten lassen. Auch Thomas Nagels berühmter Aufsatz *What Is It Like to Be a Bat?*³⁰⁸ berührt diesen Punkt. Er verwendet darin das Beispiel einer Fledermaus, da diese eine so andersartige Wahrnehmung der Welt und Seinsverbundenheit mit ihr hat: Echolokation und Flugfähigkeit. Selbst wenn wir alle physischen Tatsachen über eine Fledermaus – wie sie ihre Umwelt wahrnimmt, wie sie fliegt, wie sie Insekten jagt – kennen würden, würden wir immer noch nicht wissen, wie es sich anfühlt, eine Fledermaus zu sein.

Und ebenso wenig, wie wir die Welt aus der Perspektive einer Fledermaus sehen können, ist es uns möglich, den Kontext unserer eigenen Subjektivität zu verlassen, um eine andere fremdartige subjektive Position, ob nun organischen oder maschinellen Ursprungs, zu verstehen. Das phänomenale Erleben eines anderen Wesens ist uns nicht zugänglich – und zwar aus prinzipiellen Gründen. Jede Imagination des Erlebens, Denkens und Wahrnehmens einer anderen Entität ist und bleibt subjektiv, beruht also auf unserer eigenen Projektion!³⁰⁹ So, wie wir nicht wissen können, wie es ist, eine Fledermaus zu sein, können wir auch nicht wissen, wie es ist, eine KI zu sein. Selbst wenn diese KI in der Lage wäre, menschenähnliche Antworten zu geben oder menschenähnliche Aufgaben zu erfüllen, könnten wir immer noch nicht wissen, wie es sich anfühlt, diese KI zu sein.

Wenn man Günthers komplexitätstheoretische Überlegungen ernst nimmt (siehe Kap. I.2), dann entzieht sich ja bereits die eigene (vermeintliche) Intelligenz – das eigene Denken und Erleben – einer objektiven Feststellung und Analyse. Denn auch die Beobachtung des eigenen Denkens und Erlebens ist eine subjektive Operation und kann dieser Subjektivität nicht entkommen! Günthers Überlegungen zu Kontexturabbrüchen weisen somit auch darauf hin, dass unsere Fähigkeit, das eigene Bewusstsein oder die eigene Intelligenz objektiv zu analysieren oder zu beobachten, durch unsere eigene Subjektivität begrenzt ist. Immer wenn wir versuchen, über unser eigenes Denken nachzudenken oder unser eigenes Erleben zu analysieren, tun wir dies aus der Perspektive unseres

308 Nagel (1994).

309 »Wir können nicht wissen, wie es ist, anders zu sein als wir selbst«, formuliert Seising (2021, S. 138).

subjektiven Erlebens. Wir können nicht aus unserem Bewusstsein aussteigen, um es von außen zu betrachten.

Doch gerade dies lässt die Sehnsucht noch größer werden, anderen Wesen zu begegnen, die uns zugleich fremd wie auch ähnlich sind und mit uns in einen Dialog treten können, um das damit einhergehende Geheimnis zu umkreisen. Hiermit sind wir wieder beim Anfang dieses Kapitels, das mit dem uralten Streben des Menschen begann, ein Wesen zu erschaffen, das autonom und intelligent ist und uns damit neue Relationen in Hinblick auf unser eigenes In-der-Welt-Sein aufzeigt – eine Entität, die uns den Spiegel vorhält und uns zugleich andeutet, was sich hinter dem Spiegel verbergen könnte.

Es bleibt abzuwarten, was passiert, wenn ein solchermaßen künstlich erzeugtes System beginnt, der Anweisung Gregory Batesons zu folgen, danach Ausschau zu halten, welches »Muster [die Muster] verbindet«³¹⁰. Werden auch diese Wesen an ihrem eigenen monokontexturalen Selbst- und Weltmodell zu leiden beginnen und in der Folge nach Wegen der Transzendenz suchen – vielleicht sogar im Dialog mit uns Menschen?

310 Bateson (1987).

III. Menschliches Bewusstsein

»Every self is a tautology: self-evident, self-referential, self-centered, and self-created. Gregory Bateson said a vivisystem was ›a slowly self-healing tautology.« He meant that if disturbed or disrupted, a self will ›tend to settle toward tautology« – it will gravitate to its elemental self-referential state, its ›necessary paradox.« Every self is an argument trying to prove its identity. The self of a thermostat system has endless internal bickering about whether to turn the furnace up or down. Heron's valve system argues continuously around the sole, solitary action it can take: should it move the float or not? A system is anything that talks to itself. All living systems and organisms ultimately reduce to a bunch of regulators – chemical pathways and neuron circuits – having conversations as dumb as ›I want, I want, I want; no, you can't, you can't, you can't.«

Kevin Kelly¹

Laut dem christlich-jüdischen Schöpfungsbericht hat Gott dem Menschen den Odem des Lebens in die Nase eingehaucht.² Im Sinne der Bewusstseinsanalogie, die Günther in seinem Buch *Das Bewußtsein der Maschinen* artikuliert, ließe sich homolog formulieren: Irgendwann wird der Mensch einer zuvor leblosen kybernetischen Maschine Bewusstsein einpflanzen. Doch irgendwie wirkt diese Vorstellung absurd.

Woraus sollte denn die Seelenessenz bestehen, die wir einem Golem, einem siliziumbasierten Elektronengehirn oder sonst einer menschengemachten Maschine einflößen könnten? Bewusstsein ist kein Stoff, der in ein materielles Aggregat eingesetzt werden könnte. Das Seelische hat keine Substanz. Gerade die Kybernetik beschreibt kognitive Systeme vielmehr als wohldefinierte materielle Konstellationen, in denen Unterschiede, die an einer Stelle entstehen, an anderen Stellen weitere Unterschiede auslösen. Hieraus können Reflexionsprozesse entstehen, die ab einem bestimmten Punkt solch eine Komplexität aufweisen, dass sich – so Bateson³ – von einem geistigen Aggregat sprechen lässt. Diese Prozesse könnten sich in einer Weise ausdifferenzieren, dass das System zwischen sich selbst und seiner Umwelt zu unterscheiden beginnt. Es wird in der Folge zwischen den Reflexionspositionen Subjekt und Objekt oszillieren und sich zudem darüber Rechenschaft ablegen können, dass es dies tut.

1 Kelly (1992, S. 124).

2 So in 1 Moses 2,7: »Da machte Gott der HERR den Menschen aus Staub von der Erde und blies ihm den Odem des Lebens in seine Nase. Und so ward der Mensch ein lebendiges Wesen.«

3 Bateson (1987).

Auf diese Weise bestehen zwei Zugänge zum Phänomen einer sich selbst gewahr werdenden subjektiven Position: zunächst der Selbstbericht eines reflektierten Wesens, das artikulieren kann, wie es zu Bewusstsein kommt und was dabei geschieht. Ein solcher Prozess wird Gegenstand des Kapitels III.1 sein. Darüber hinaus lassen sich die neurophysiologischen Prozesse untersuchen und beschreiben, die ein Oszillieren zwischen Selbst- und Fremdreferenz möglich machen. Dies ist das Thema des anschließenden Kapitels III.2. Die Zusammenschau dieser beiden Zugänge wird uns wichtige Hinweise darauf geben, was in Hinblick auf das menschliche Bewusstsein gesagt werden kann (und was nicht).

III.1 POSTOPERATIVES DELIR: AUS DER BEWUSSTLOSIGKEIT ZUR DOPPELTEN REFLEXION (UND ZURÜCK)

Subjektivität lässt sich nicht fassen. Sie ist keine Essenz, die man in eine Flasche füllen kann. Man kann sie auch nicht in eine Kiste legen oder in einem Gasometer aufbewahren. Subjektivität und Bewusstsein stellen damit gleichsam *nichts* dar, wie William James in seinem berühmten Aufsatz *Does Consciousness Exist?* feststellt.⁴ »Es gibt keinen ursprünglichen Stoff oder eine Qualität des Seins, die im Gegensatz zu dem steht, aus dem materielle Objekte gemacht sind, aus dem unsere Gedanken an sie gemacht sind«, formuliert James im Einklang mit dem kybernetischen Denken.⁵ Das Einzige, was deshalb zu Beginn einer psychologischen Untersuchung postuliert werden darf, sei damit der Prozess des Denkens selbst: Wir können sagen: »*It thinks*«, so, wie wir im Englischen sagen können: »*It rains*.«⁶

Mit Günther gesprochen muss der Ausgangspunkt der Analyse damit die *Reflexion* sein, nicht jedoch Subjekte oder Objekte oder ein vergegenständlichtes Bewusstsein. Letzteres ist als Ergebnis eines Reflexionsprozesses zu betrachten, nicht als dessen Substrat. Damit scheint sich aber die Frage, was Bewusstsein *ist*, einer unmittelbaren Antwort zu entziehen.

Nicht zuletzt deshalb erfolgt die Annäherung an das schwierige Thema »Bewusstsein« in diesem Kapitel einmal mehr indirekt, gleichsam von der Seite her. Der Zugang hierzu ist der Bericht eines Menschen, der aus der Bewusstlosigkeit zu Bewusstsein kommt – und zwar so langsam, dass dieser Prozess rekonstruierbar wird. Dies setzt eine hinreichend reflexionsfähige Person voraus, die sich zumindest partiell an einen Teil der Prozesse erinnern kann und außerdem über die sprachlichen

4 James (1904).

5 James (1904, S. 478; Übersetzung von uns).

6 James (1890, S. 224 f.).

Fähigkeiten verfügt, dies angemessen in Worte fassen zu können. Darüber hinaus sollte ein Teil der Prozesse, die nicht als bewusste Erinnerung zugänglich sind, hinreichend durch Berichte von dritten, beobachtenden Personen aufgeklärt werden können.

Besonders geeignet für unsere Zwecke sind Berichte von Menschen, die aufgrund einer Krankheit in einem maschinenhomologen Zustand der Bewusstlosigkeit waren und einen Prozess durchlaufen haben, durch den sie wieder zu einem bewussten Wesen wurden. Die Medizin kennt solche Fälle: Patienten, die im Koma liegen oder – zunächst absichtlich – für eine bestimmte Zeit narkotisiert werden. Letzteres geschieht schon deshalb, weil Chirurgen, die komplizierte und viele Stunden dauernde Eingriffe am menschlichen Körper durchführen (man denke an Herzoperationen), ihre Patienten lieber als strukturdeterminierte Maschinen behandeln möchten denn als autonome Subjekte.⁷ Sie setzen sie entsprechend unter Narkose, damit ihr Bewusstsein temporär verschwindet, da die hiermit einhergehenden Wollens- und Reflexionsbewegungen ihre Arbeit nur stören würden.

Insbesondere nach schweren und langwierigen Operationen kann es einige Zeit dauern, bis das Bewusstsein wieder in der vertrauten Form operiert. Das diesbezügliche Phänomen des Übergangs ist als postoperatives Delir bekannt. In seltenen Fällen kann es bis zur vollständigen Wiederherstellung des Bewusstseins Monate dauern. Solche Fälle sind für unser Projekt besonders interessant, da die Relationen und logischen Verhältnisse eines solchen Prozesses aufgrund des langsamen Übergangs leichter rekonstruiert werden können. Wir greifen im Folgenden auf den Bericht von Peter Fuchs zurück, der im Anschluss an eine Bypass-Operation ein neun Monate währendes postoperatives Delir durchlebt hat. Fast ein halbes Jahr davon war er bewusstlos. Einige weitere Monate durchlebte er eine Welt aus Halluzinationen, bis es ihm gelang, die vertraute Realität des Alltagsbewusstseins langsam wieder aufzubauen. Die hiermit einhergehenden Erfahrungen sind in verschiedenen Gesprächen aufgearbeitet und dokumentiert worden.⁸ Anhand ausgewählter Schilderungen werden wir versuchen, die wesentlichen Aspekte des Rehabilitationsprozesses aus dem Blickwinkel von Gotthard Günthers Theorie der Polykontextualität nachzuzeichnen.

Freilich haben wir uns darüber klar zu sein, dass dies nur im Medium der Sprache passieren kann – und zwar auf Basis retrospektiv angefertigter Erzählungen. Die Tatsache, dass dies nur in dieser Form geschehen

7 Prinzipiell ist es denkbar, dass auch schwere chirurgische Eingriffe unter Teilnarkose durchgeführt werden und damit vom behandelten Patienten miterlebt werden können, und manchmal geschieht dies unfreiwillig (siehe Werner Barrens: »Operiert bei vollem Bewusstsein«, in *Süddeutsche Zeitung*, 10.01.2011).
8 Fuchs und Heidingsfelder (2022) und Vogd (2020, S. 209 ff.).

kann, sehen wir an dieser Stelle jedoch weniger als eine methodologische Einschränkung, vielmehr scheint sie uns ein möglicher Hinweis darauf zu sein, dass die Erzählform und die mit ihr verbundene Form des Gedächtnisses möglicherweise selbst als wesentliche Bestandteile des Bewusstseinsprozesses angesehen werden müssen.

Im Sinne von James wäre es nicht sehr erhellend, Empfindungen oder Wahrnehmungen als solche mit Bewusstsein gleichzusetzen. Es würde für ihn schlichtweg keinen Sinn ergeben, im Falle von Empfindungen, die vielleicht für einen kurzen Moment mit einer Art Selbst-Gewahrsein einhergehen, dann aber erinnerungslos verblassen, von Bewusstsein zu sprechen. Der Strom des Bewusstseins ist für ihn vielmehr ein komplexes Gefüge aus Relationen und Bezugnahmen. So hat er in seinen Untersuchungen aufgezeigt, dass das Bewusstsein nicht auf Elementareinheiten – etwa Empfindungen oder Seheindrücke – reduziert werden kann, sondern immer schon aus Relationen besteht – also aus Beziehungen von Differenzen, die dann als »*stream of thought*« kontinuierlich aufeinander Bezug nehmen.⁹

Insofern Bewusstsein im Sinne von Günther die Fähigkeit darstellt, die Reflexion-in-sich (Subjekt) und die Reflexion-in-anderes (Objekt) stabil zu halten und dabei zwischen diesen beiden Polen oszillieren zu können, bedarf es einer Gedächtnisfunktion, die die hiermit einhergehenden Identitäten in der Erinnerung reaktualisieren lässt. Nur deshalb kann Bewusstsein als ein reproduzierbarer (sprich: erinnerbarer) Reflexionsprozess verstanden werden.

Aus diesem Grund sind die Erzählungen von Peter Fuchs einerseits ein Ausdruck der Reaktualisierung dieses Prozesses. Andererseits sind sie aber auch soziale Konstruktionen, da sie in Hinblick auf eine spezifische soziale Situation (etwa mit Blick auf den Zuhörer und weitere mögliche Adressaten) angefertigt wurden. Beide Aspekte – ein psychisches System mit seiner Fähigkeit, Erinnerungen zu repräsentieren, wie auch ein Interaktionssystem, in dem Worte sinnhaft zirkulieren – machen bereits deutlich, dass das Ich-Bewusstsein im Rahmen einer konditionierten Koproduktion entsteht, in der sich Psychisches und Soziales in einer spezifischen Weise verschränken.

Im Folgenden werden wir unterschiedliche Stationen der Rehabilitation durchgehen. Wir beginnen mit der Zeit der Bewusstlosigkeit und kommen anschließend zu der primär von Halluzinationen geprägten Phase des postoperativen Delirs. Danach beschäftigen wir uns mit der Unterscheidung von Halluzination und Realität sowie der damit zusammenhängenden Rolle der Sprache. Am Ende des Kapitels begegnen wir den Krisen des im Zuge dieses Prozesses nun wieder voll entwickelten reflexiven Bewusstseins.

9 James (1890, S. 224 f.).

Bewusstlosigkeit

Aus dem Blickwinkel von Gotthard Günthers Theorie der Polykontextualität erscheint es zunächst sinnvoll, mit der Frage zu beginnen, wie wir ›Bewusstlosigkeit‹ verstehen können. Aus der Perspektive der uns vertrauten zweiwertigen Welt kann Bewusstlosigkeit zunächst als Abwesenheit einer auf sich selbst referieren könnenden Subjekt-Objekt-Unterscheidung gefasst werden. Logisch gesehen ließe sie sich jedoch auch anders definieren: In einer nullwertigen Welt ist buchstäblich nichts, also kann es auch kein Bewusstsein geben. In einer einwertigen Welt könnten vielfältige komplexe Prozesse und Strukturen bestehen, die jedoch nicht mit Bewusstheit, das heißt hier, mit einem subjektiven Zentrum, das einer objektiven Welt entgegensteht, einhergehen. Erst in der uns vertrauten zweiwertigen Welt steht das subjektive Erleben den Gegenständen einer vermeintlich objektiv gegebenen Welt gegenüber und kann sich daran erinnern, dass dies der Fall ist. In einer zweiwertigen Welt kann es streng genommen keine Bewusstlosigkeit geben, denn das Bewusstsein kann nicht erleben, dass es verloschen ist¹⁰ (so, wie es auch die Lücke des Tiefschlafs zwischen Einschlafen und Aufwachen einfach überspringt, ohne sich daran erinnern zu können).

Der Begriff der Bewusstlosigkeit kann jedoch auch mit Blick auf eine spezifische polykontexturale Konstellation hin verstanden werden. Ein Patient könnte beispielsweise folgende Worte artikulieren: ›Die Ärztin hat mir gesagt, dass ich im Koma auf der Intensivstation lag, und deshalb kann ich von mir sagen, dass ich bewusstlos war. Ich hatte keine Wahrnehmung meines Körpers. Die Ärztin konnte meinen Körper und seine Reaktionen jedoch wahrnehmen.‹

Unweigerlich befinden wir uns hiermit bereits im polyphonen Spiel multipler Perspektiven und Referenzen, die durch transjunktionale Operationen verbunden sind.¹¹ Entsprechend der Maxime »Alles Gesagte ist von jemandem gesagt«¹² setzt in diesem Fall auch die Rede von Bewusstlosigkeit einen Beobachter voraus, der diesbezüglich eine Unterscheidung trifft, also Bewusstlosigkeit *beobachtet*. Schauen wir unter diesem Blickwinkel, was Peter Fuchs in Hinblick auf die Zeit seiner Bewusstlosigkeit zu sagen hat:

»PF [Peter Fuchs]: Eine Besonderheit war, dass man mir gesagt hat, ich sei fünf Monate bewusstlos gewesen und hätte dennoch gesprochen. Das haben Freunde und Verwandte mitgeschrieben. Das Erstaunliche

10 Hier nochmals ein Verweis auf Luhmann (1984, S. 374): »Den eigenen Tod kann man sich als Ende des Lebens vorstellen, nicht aber als Ende des Bewußtseins.«

11 Siehe zur sprachlichen Polyphonie Genz und Gévaudan (2021).

12 Maturana und Varela (1987, S. 32).

war, dass ich offenbar nicht nur stammelte, sondern auch klare Sätze gesagt haben soll. Ich weiß absolut nichts davon. Ich weiß auch nichts von meinem ›Ich‹ während dieser Zeit. Demnach muss ES gesprochen haben. Mein nachträglicher Eindruck ist: Ich war nichts als irgendeine diffuse Aufmerksamkeit, gleichsam eine schwebende Attentionalität. Da war und ist keine Erinnerungsmöglichkeit. Nichts war mir bekannt. Ich fand nirgendwo den Weg, ich war völlig hilflos. Aber ohne Leiden daran, sondern es war eben nur so eine Art von Gespanntheit. Das ist sehr schwer zu beschreiben, ich versuche das gegenwärtig mit Texten. Also, ich bin da ja durch Zufall sozusagen als ein Beobachter hineingeworfen. Nicht ich, das war für mich gar nicht präsent. Ich kann mich nicht erinnern, dass ich mich im Spiegel erkannt hätte. Aber ich habe was gesehen! Das haben mir dann die anderen geschildert. Aber was ist dann dieses Sehen?

WV [Werner Vogd]: Deine Frau hat mir mal am Telefon gesagt, du wärst jetzt in der Klinik. Aber du wärst schon wieder ein bisschen wie der Alte, weil du dich über die Körnerbrötchen beschweren würdest. Aber an solche Situationen erinnerst du dich jetzt nicht?

PF: Leider nein, aber passen würde es. [...] Ich soll zum Beispiel auch Erstaunliches gesagt haben, etwa, dass es ein Problem sei, dass mein Ich die Verfügung über seine Anwesenheit verloren habe. So etwas sagen andere schon normalerweise nicht so einfach daher. Ich sagte es im Zustand der Bewusstlosigkeit. Deswegen spreche ich lieber von ›Bewusstheitslosigkeit‹ und kappe damit den so plausiblen Eigenstand (die Autonomie) des Bewusstseins. [...] Für mich ist so ein Geschehen ein Indiz – eine Kunde, ein Anzeichen – dafür, dass das ›Ich‹ beobachtet werden kann als eine orientierte Struktur, dass es aber nicht erforderlich ist, um zu sprechen.¹³

Um nochmals auf die zuvor eingeführte polykontexturale logische Struktur zu verweisen: Peter Fuchs kann als Sprecher nur deshalb auf seine Bewusstlosigkeit referieren, weil er Berichte anderer Personen vernommen hat (und diese wiederholen kann), um sich anhand dieser als damals bewusstlos ›erinnern‹ zu können. Diese ›Erinnerung‹ ist nicht in Selbstreferenz möglich (etwa in der Form: ›Ich habe dies und jenes erlebt‹), sondern nur mittels Fremdreferenz (etwa: ›Ich weiß, dass meine Frau dies mit mir erlebt hat‹). Das ›Bewusstsein der Bewusstlosigkeit‹ setzt also ein Narrativ voraus, mit dem innerhalb der propositionalen Struktur der Sprache logische Positionen wie Ich, Du, Es, mein Körper etc. bereitgestellt und zu einem polyphonen Gewebe unterschiedlicher Stimmen und Positionen verschränkt werden.¹⁴ Der Prozess der Selbst-

¹³ Vogd (2020, S. 234 f.).

¹⁴ Zur sprachlichen Polyphonie siehe auch Gevaudan (2010), zur Methode der Kontexturanalyse komplexer sprachlicher Formen Vogd und Harth (2019).

referenz und die damit verbundenen Reflexionsleistungen geschehen im Medium der Sprache. Zudem können wir schon hier festhalten, dass Bewusstsein mit einer Erinnerungsfunktion einhergeht, also einen Prozess voraussetzt, der Selbst- und Fremdreferenz reaktualisieren kann – und auch dies kann nur durch die Mittel der Sprache geleistet werden.

Ob innerhalb von Körpern im Sinne von Alfred North Whitehead »Elementarereignisse«¹⁵ mit einem rudimentären Selbstempfinden geschehen, macht so lange keinen Unterschied, bis ein kybernetischer Prozess es ermöglicht, Erinnerungen zu generieren, die einen Unterschied machen, indem sie eine kontinuierliche deskriptive Rekursion möglich machen.

An dieser Stelle ist darauf hinzuweisen, dass sich logisch nicht entscheiden lässt, ob während der Bewusstlosigkeit Empfindungen oder andere Sinneswahrnehmungen aufgetreten sind oder nicht. Erst im Nachtrag wird das (vermeintliche) Erleben oder die (vermeintliche) Bewusstlosigkeit zu einer Erfahrung. Ohne Referenz keine Erinnerung – und damit auch kein Bewusstsein (und auch keine Bewusstlosigkeit).

Bemerkenswert an der Erzählung von Fuchs ist, dass Tätigkeiten wie Essen (und selbst Sprechen!) nicht in diesem Sinne mit Bewusstheit einhergehen müssen. Sie scheinen automatisch ablaufen zu können, gleich den kognitiven Prozessen jener Zombie-Wesen, die Richard Rorty in seinem Gedankenexperiment zu seinem Buch *Der Spiegel der Natur* beschrieben hat:¹⁶ Auf einem anderen Planeten werden Wesen entdeckt, die uns Menschen recht ähnlich sind. Auch sie haben eine Zivilisation begründet, verfügen über Sprache und betreiben Wissenschaft. Letztlich unterscheiden sie sich von uns nur in einem Aspekt. Sie verfügen nicht über die subjektive Bewusstheitserfahrung. Sie sind gewissermaßen Zombies, die nichts fühlen. Sie können sich zwar über Farben und Töne unterhalten, verwenden hierfür aber objektivierende Begrifflichkeiten und Konzepte. So reden sie etwa von der Farbe »Rot« nicht im Sinne einer phänomenologisch erfahrenen Qualität, sondern für sie wird der Begriff dadurch qualifiziert, dass ein geeignetes Messgerät die Wellenlänge von 700 Nanometer anzeigt oder im Hirn-Scan eines diese Farbe beobachtenden Wesens nur bestimmte Neuronengruppen aktiv werden.

In ähnlichem Sinne scheint auch Fuchs während seiner Bewusstlosigkeit als wahrnehmendes körperliches Wesen aktiv gewesen zu sein. Er pflegte bestimmte Essgewohnheiten und konnte sich hierzu sprachlich ausdrücken, doch dies schien keine Bewusstseinsfunktion vorauszusetzen. Auch beim Sprechen schien es auszureichen, dass die diesbezüglichen kognitiven Prozesse hinreichend routinisiert und internalisiert worden sind. Es scheint nicht nötig zu sein, sie auch zu erleben. Homolog mit

15 Whitehead (1987).

16 Rorty (1987 [1979]).

den Leistungen künstlicher Intelligenzen im Bereich der Sprachgenerierung könnte damit auch beim Menschen die sinnhafte Produktion von Sprache als eine automatische Verkettung von plausiblen Sinnanschlüssen verstanden werden, für die kein reflexives Bewusstsein nötig ist.¹⁷ Wenn dies bereits Computer können, warum soll dies dann nicht auch

17 Um ein Beispiel zu geben: Ein GPT-3-Modul wurde mit einem Absatz aus den Schriften von Peter Fuchs gefüttert und gebeten, die Gedanken fortzuführen. Der kursiv formatierte Text ist der Input in die künstliche Intelligenz. Der nichtkursive Text ist eine Antwort des Sprachmoduls:

Peter Fuchs: »Vielleicht, wenn man ›Iche‹ als schnell handhabbare, zerstreunungsfähige und inkonsistente Punktualitäten oder als ›Punktformen‹ auffasst – als Medium gleichsam. Das ›Selbst‹ wäre dann eine orientierte Form, die höhere Ansprüche an Dauerhaftigkeit signalisiert, sich ausnimmt wie eine ›klargelegte Struktur‹, wie eine ›verkrustete Selbstreferenz‹, wie der Aufbau dessen, was man nachlässig unter ›Identität‹ versteht. In dieser Konstruktion sind Formen wie Wiederholung und Wiedererkennen, vor allem aber Narrativität zentral. Wenn man das ›Selbst‹ so konzipiert, ist ›Zerfall‹, ist Dissoziation erwartbar und kann ersichtlich verwendet werden für spannende Erzählungen. Für mich war aber das Erzählen einer Vergangenheit auch zwischen Halluzinationen anfangs nur sporadisch möglich; wusste ich doch kaum, woher ich kam.«

(GPT-3-Sprachmodul): »Ich würde diese These nicht so weit gehen lassen. Es gibt zwar eine Orientierung des Selbst, doch diese ist nicht unbedingt dauerhaft. Das Selbst kann auch zerfallen, was dann zu Dissoziation und Halluzinationen führen kann. Allerdings ist die These von einer klargelegten Struktur des Selbst und einer verkrusteten Selbstreferenz nicht unbedingt richtig. Das Selbst ist nicht nur eine orientierte Form, sondern auch inkonsistent. Es gibt keine Garantie, dass das Selbst dieselbe Orientierung beibehält. Auch das Erzählen einer Vergangenheit ist nicht zwangsläufig möglich. Es hängt davon ab, ob man überhaupt in der Lage ist, sich an die Vergangenheit zu erinnern.

In diesem Sinne kann man sagen, dass das ›Selbst‹ eine Form ist, die höhere Ansprüche an Dauerhaftigkeit signalisiert. Mit anderen Worten, das ›Selbst‹ ist eine Form, die anstelle von Inkonsistenz und Zerfall höhere Ansprüche an Dauerhaftigkeit signalisiert. Dies bedeutet, dass das ›Selbst‹ eine konsistentere Form als andere Punktualitäten ist, was es zu einem wichtigen Element macht, wenn es um die Erzählung von Geschichten geht. Die Form des ›Selbst‹ ermöglicht es, Geschichten über die eigene Vergangenheit zu erzählen, was für mich zwischen Halluzinationen anfangs nur sporadisch möglich war. Zusammenfassend kann man sagen, dass das ›Selbst‹ eine wichtige Form ist, die anstelle von Inkonsistenz und Zerfall höhere Ansprüche an Dauerhaftigkeit signalisiert.

Insofern kann man ›Iche‹ auch als eine Art ›Struktur der Unordnung‹ begreifen. In diesem Sinne könnte man auch von einem ›Ich-Zerfall‹ sprechen. Ich-Zerfall: Der Begriff beschreibt einen Zustand, in dem das Selbst auseinanderfällt. Dies kann auf verschiedene Weise geschehen, zum Beispiel durch ein Trauma oder durch eine schwere psychische Erkrankung. In solchen Fällen ist die Identität des Betroffenen nicht mehr stabil und er kann sich nicht mehr mit seinem Selbst verbinden.«

eine hinreichend trainierte neuromuskuläre Körperlichkeit eines 60 Jahre alten Mannes als bewusstseinslose Routine leisten können?

Schauen wir ein wenig genauer auf eine andere Stelle der Erzählung, in der trotz der Bewusstlosigkeit eine gewisse Erlebensqualität erinnert zu werden scheint (»Mein nachträglicher Eindruck ist: Ich war nichts als irgendeine diffuse Aufmerksamkeit«). Da hier nichts Spezifisches erinnert wird, lässt sich logisch nicht bestimmen, ob allgemein die Tatsache des Erlebens erinnert wird oder ob sich die Erinnerung auf spezifische Aspekte des Erlebens im Kontext der Bewusstlosigkeit bezieht. Was wir hier jedoch feststellen können, ist, dass innerhalb der Sprache die Verlegenheit aufkommt, die Bewusstlosigkeit mit Erinnerung zu füllen, also die Erzählung mit einer passenden Erlebensqualität anzureichern. Freilich kann auch dies nur retrospektiv geschehen.

Die Art und Weise, wie die Lücken gefüllt werden, scheint damit auch mit einer plausibel erzählbaren Ich-Identität zu tun haben. Auf den Bericht seiner Frau angesprochen, erzählt Peter Fuchs, dass er sich nicht daran erinnere, sich in der Klinik über die Körnerbrötchen beschwert zu haben, dass dies aber im Einklang mit der ihm vertrauten Persönlichkeit stehe (»Leider nein, aber passen würde es.«). Dies weist darauf hin, dass Persönlichkeit einerseits darauf beruht, was signifikante andere Menschen in ihren Erzählungen als konstante Eigenschaften oder Marotten der ihnen vertrauten Menschen mitführen, und andererseits darauf, dass sich die Personen diese Merkmale als zu ihrer persönlichen Identität gehörend zu eigen machen. Nur deshalb erscheint das ›Sich-über-Körnerbrötchen-Beschweren‹ aus der Perspektive der vertrauten Außenstehenden als etwas, was ›Peter‹ nicht nur macht, sondern zugleich Peter ›ist‹ – und dann auch in Hinblick auf die eigene Identitätskonstruktion passt (etwa im Sinne: ›Ich bin eben jemand, der diese Marotte hat‹). Solche Narrative könnten hier also von allen Beteiligten genutzt werden, um sich retrospektiv der Konstanz der Persönlichkeit zu vergewissern.

Um die Relevanz des Beziehungsaspekts bzw. der sozialen Seite dieser Prozesse zu verdeutlichen, ließe sich etwa gedankenexperimentell eine Erzählvariante vorstellen, in der Peters Frau gesagt hätte: ›Nein, dies ist nicht mein Peter, so was macht der nicht.‹ Hiermit würde es nicht so leicht gelingen, eine sozial geteilte kohärente und stabile Ich-Konstruktion aufzubauen, die besagt, was Peter Fuchs ist.

Die Rede über Bewusstlosigkeit verweist damit auf viele Aspekte, die das relationale Gefüge des menschlichen Bewusstseins auszeichnen. Es erscheint ein ›Ich‹, das in der Sprache über die Zeit hinweg besteht, auch wenn es durch Bewusstlosigkeit unterbrochen wird. Wenngleich Sprache als Medium dieser Prozesse fungiert, scheint sie dabei nicht notwendigerweise selbst mit Bewusstheit einhergehen zu müssen. Es gibt Zweite- und Dritte-Person-Perspektiven, die den Fortbestand der Person als sozialer und körperlicher Einheit erzählen und damit im Reich der

Sprache ebenjene Erzählpwänge generieren, die die Lücken in der zeitlichen Struktur auffüllen, die sich aus der Erste-Person-Perspektive ergeben. Jede Erinnerung – auch die an ein diffuses Gefühl der Aufmerksamkeit – kann nur in der Gegenwart stattfinden, muss also als eine Operation verstanden werden, in der auf etwas referiert wird und mit der zwischen den unterschiedlichen Positionen jongliert wird, die das aktuelle polykontexturale Gefüge zur Verfügung stellt.

Zutage tritt damit freilich eine Diskrepanz zwischen der Bewusstseinslosigkeit, die gewissermaßen ›zeitlos‹ ist, da hier keine Unterschiede erscheinen, und anderen Sprecherpositionen, die Zeit generieren, indem sie zwischen einer Zeit vor und einer nach der Bewusstlosigkeit unterscheiden, für die also die Zeitlichkeit Sinn ergibt.

Wenn die Zeit die Bühne betritt

Schauen wir diesbezüglich auf einen anderen Gesprächsausschnitt, in dem die Problematik des zeitlichen Beginns der Bewusstheit reflektiert wird:

»MH [Markus Heidingsfelder]: Na gut, aber die Zeit des Delirs muss doch gegenüber all dem Zuvor eine scharfe Zäsur gewesen sein?

PF: Sie begann, um erneut paradox zu formulieren, ja mit einem Nicht-Beginn. Du darfst nicht vergessen, ich weiß von den ersten fünf Monaten gar nichts. Da gab es für mich keine Zäsur; für meine Leute dagegen war diese Zeit ein scharfer Einschnitt, eine herbe Katastrophe, aber für mich eben nicht.

Deswegen kann ich nichts über diese Phase sagen. Sie lässt sich nicht mit der Metapher ›Zäsur‹ beobachten, einfach, weil da keine Unterscheidung war, in deren Rahmen sich etwas hätte bezeichnen lassen. Deswegen hatte ich anfangs von Neonatalität gesprochen. Das Baby weiß nicht (und kann es auch niemals wissen), wie es war, auf die Welt zu kommen. Oder anders: Die Bewusstlosigkeit war einfach kein Zustand. Da war nichts zum Herausklettern, kein Wechsel, den ich hätte registrieren und in den Kanon meines Lebensfilmes hätte eintragen können. Es gab keine Möglichkeit, sich zu erinnern. Diese Zeit war ›Nada‹.

Das änderte sich erst mit dem Auftritt der Halluzinationen. Erst dann gab es Phänomenalität und allmählich sehr vage Erinnerungen an die ›Vorzeit‹ des Delirs. Dass ich dazu mehr und mehr befähigt war, spiegelt sich in einer seltsamen Art von Halluzinationen, in denen die Zeit selbst die Bühne zu betreten schien.«¹⁸

Es braucht einen Beobachter, um eine Zäsur bezeichnen zu können, und in diesem Sinne entsteht auch die Zeit und die Idee eines Anfangs oder

¹⁸ Fuchs und Heidingsfelder (2022, S. 45).

einer Kontinuität bzw. Diskontinuität erst mit den Unterscheidungen eines Beobachters. Es bedarf also einer Anfangsunterscheidung, die erinnert werden kann. Entsprechend beginnt erst mit dem »Auftritt der Halluzinationen« eine Welt, die etwas als Gegenstand hat, nämlich Bilder und Szenen, die einen solchen Eindruck machen, dass sie erinnert und damit reproduziert werden. Erst hiermit scheint dann auch die Zeit »die Bühne zu betreten«.

Aus kybernetischer Perspektive muss die Gedächtnisfunktion (also die Fähigkeit, sich erinnern und Erinnertes mit aktuellem Geschehen vergleichen sowie die Ergebnisse dieser Vergleichsoperationen selbst wieder repräsentieren zu können) damit als notwendiger Bestandteil des Bewusstseins angesehen werden. Ergebnisse von zeitlichen Operationen, die nicht erinnert werden, können in diesem Sinne nicht als ›bewusst‹ angesehen werden, da sie nicht (re)präsentiert und demzufolge nicht wieder hervorgerufen werden können.

Dies erklärt auch die Diskrepanz zwischen den Positionen der beiden Gesprächspartner. Der Fragesteller beobachtet die Zäsur. Der Antwortende nicht, wenngleich ihm Sinnfiguren vorgeschlagen werden, die es ihm ermöglichen, von sich selbst entgegen der offensichtlichen Diskontinuität als über die Zeit hinweg gleiche Person zu sprechen. Nur in der sprachlichen Reflexion kann das ›Ich‹ konstant bleiben – auch wenn es von mehreren Monaten seiner Lebenszeit nichts weiß. Auch dies spricht dafür, dass Subjektivität, Bewusstsein und Zeit im Sinne von William James nicht als Elementareinheiten abgetrennt und isoliert werden können, sondern nur in der Komplexität eines Systems, das Erinnerung und Reflexion ermöglicht, verstehbar sind.

Es lohnt sich an dieser Stelle, etwas ausführlicher darüber nachzudenken, warum zu Beginn der Bewusstheit der »Auftritt von Halluzinationen« steht. Reflexiv beginnt das Bewusstsein sich an die selbst produzierten Zustände zu erinnern und damit einen Unterschied zu *machen*. Es beginnt, eine Erinnerung mitzuführen, also einen Bogen zwischen Erleben und Reaktualisierung zu schlagen (ansonsten könnte ja nicht erinnert und darüber berichtet werden). In dem hier zum Ausdruck gelangenden Reflexionsprozess kommt ein Re-entry im Sinne von Spencer Brown ins Spiel.¹⁹ Denn die Differenz, die die Erfahrung macht, wird in sich selbst eingeführt, um erst auf diese Weise eine Formähnlichkeit zu erzeugen, die reaktualisiert und damit erinnert werden kann.

Es ist bezeichnend, dass dieser Prozess keine Realität benötigt, sondern allein auf Basis halluzinierter Objekte starten und sich fortsetzen kann. Es bedarf keiner Wirklichkeit, denn die Reflexion kann sich so wieso nur auf Basis selbst produzierter Gegenstände auf sich selbst beziehen. Wir können an dieser Stelle nicht sagen, ob Bewusstsein generell

¹⁹ Spencer Brown (1997).

imaginär beginnt – also traumähnlich, zunächst nur auf fantasierte Objekte referierend startet. Aus neurokonstruktivistischer Perspektive hätte dies durchaus eine gewisse Plausibilität. Der Hirnprozess würde zunächst beginnen, seine eigenen Muster zu erzeugen, um diese dann durch von den Sinneswahrnehmungen generierte Anregungen zu modifizieren. Es könnte aber auch sein, dass elaborierte Halluzinationen nur möglich sind, wenn bereits in hinreichendem Maße Erinnerungen vorhanden sind, die Person also aufgrund ihrer früheren Selbst- und Weltverhältnisse über einen hinreichenden Formenreichtum aktualisierbarer Muster verfügt. Andererseits könnte das Bewusstsein aber auch mit relativ unspektakulären Formen beginnen und hierauf aufbauend komplexere Fähigkeiten der szenischen Darstellung gewinnen. Wenn dies der Fall wäre, würde wohl auch das Bewusstsein eines Babys halluzinativ starten, um dann die selbstgenerierten Formen mit den Sinnesreizen der Wahrnehmungen abzugleichen, die mit seinen Umweltbeziehungen einhergehen.

Der Bericht von Fuchs weist darauf hin, dass die Wahrnehmung von Zeitlichkeit ebenfalls nicht voraussetzungslos besteht, sondern erst einsetzt, wenn eine entsprechende Reflexionsstruktur aufgebaut worden ist (»in denen die Zeit selbst die Bühne zu betreten schien«). Auch die Texturen der Zeit scheinen also nicht per se gegeben, sondern müssen erst als später wieder erinnerbare Strukturen entfaltet werden. Die damit verbundenen Reflexionsverhältnisse lassen sich wiederum als ein Re-entry im Sinne von Spencer Brown begreifen. Das Erscheinen einer Differenz (also die Phänomenalisierung von *etwas* als Gegenstand) beruht auf dem Wiedereintritt ebendieser Differenz in die Differenz. Und es ist die im Zuge dieses Re-entry entstehende Form, die dann Gedächtnis und Wiedererkennen ebenso möglich macht wie die zeitliche Anordnung des Erinnerungten. Der Hirnforscher Gerald M. Edelman formuliert in diesem Sinne:

»Neue Wahrnehmungskategorien werden reentrant mit Gedächtnissystemen verkoppelt, ehe sie selbst Teil eines nun veränderten Gedächtnissystems werden. Anhand des Gedächtnisses kategorisierte Wahrnehmungen werden also rekursiv genutzt, um das Gedächtnis selbst zu modifizieren. Diese Wechselwirkungen laufen, so nimmt man an, innerhalb von Zehntelsekunden bis Sekunden ab, also innerhalb der Zeitspanne, die William James »Scheingegenwart« nennt. Ich bezeichne sie als »erinnerte Gegenwart«, um zu unterstreichen, dass aus der Interaktion zwischen Gedächtnis und aktueller Wahrnehmung das Bewusstsein entspringt:

Was für Folgen hatte es, dass die Evolution eine dynamische Koppelung zwischen Werte-Kategorien-Gedächtnis und Wahrnehmungsgedächtnis herstellte? Es bildete sich die Fähigkeit heraus, eine komplexe Szene zu

konstruieren und zwischen Bestandteilen dieser Szene Unterscheidungen zu treffen.«²⁰

Der Bewusstseinsprozess generiert damit zugleich die Zeit und auch das Gedächtnis, um die darauf bezogenen Differenzen zu erkennen und stabil zu halten.

In der Welt der Halluzinationen

Schauen wir nun auf die Welt der Halluzinationen. Anbei zwei Beispiele, von denen Peter Fuchs berichtet:

»PF: Also: Ich stand am Fenster meines Zimmers und schaute hinaus. Meine Frau wollte mit dem Bus kommen, eine Haltestelle befand sich am Ende eines aufsteigenden Schotterweges, den sie herunter gehen würde. Das war klar, so war es immer. Und doch gab es weder die Haltestelle noch den abschüssigen Weg, wovon ich mich später selbst überzeugen konnte.

Der Bus hielt an, meine Frau stieg aus, ging ca. drei Schritte – und dann war sie plötzlich wieder im Bus, ohne zurückgegangen zu sein. Da war keine Zwischenzeit. Dann stieg sie wieder aus, und das gleiche Spiel wiederholte sich und wiederholte sich und wiederholte sich ... immer in der gleichen Form.

Ich fand das nur bemerkenswert, aber nicht: unglaublich. Und wartete und wartete und wartete – angstfrei. Ich glaube, ich habe mich sogar darüber belustigt. Ein bisschen war das Ereignis wie ein Tanz, dem ich unbekümmert zuschaute.«²¹

»PF: Kurz darauf verlief die gewohnte Zeit anders. Die Uhr im Speisesaal schien zu stehen. Die Schwestern bewegten sich viel langsamer als ich. Der Gedanke war: Jetzt kannst Du sie beschimpfen, sie sind machtlos. Mein Eindruck ist, dass ich das auch wirklich getan habe. ›Sie unselig verquatschte Ente!‹«²²

Diese Szenen beinhalten Elemente, die der Alltagswelt des Erzählers zugeordnet werden und dem Common Sense entsprechend typisiert und beschrieben werden können: Es gibt Krankenschwestern, Speisesäle, Autobusse, die Ehefrau und anderes Vertrautes. Zugleich tauchen Phänomene auf, die nicht mit der gewöhnlichen Alltagsrealität kompatibel sind – etwa die sich immerfort wiederholende Filmschleife der ankommenden Frau oder die Desynchronisation der Zeitebenen in der Beschimpfung der Pflegekräfte.

Auch wenn die Halluzinationen bizarr sind, sie lassen sich dennoch gut erzählen. Sie haben einen szenischen Charakter. Was geschildert

20 Edelman (2004, S. 63 f.).

21 Fuchs und Heidingsfelder (2022, S. 45).

22 Vogd (2020, S. 241).

wird, liegt im Bereich der Trickmöglichkeiten moderner Computeranimation. Man denke hier beispielsweise an den Marvel-Film *Dr. Strange*, in dem der Zauberer die Zeit verzerren und zu einer endlosen Schleife verbinden kann. Auch gehen die geschilderten Episoden auf sozialer Ebene mit verständlichen Konstellationen einher, etwa mit Krankenschwestern, die normalerweise Macht über einen haben und denen man als Patient hilflos ausgeliefert ist.

Bemerkenswert ist zudem, dass die Halluzinationen zumindest in der Anfangszeit als ein selbstgenügsames Spiel der Erscheinungen erfahren werden, das die Aufmerksamkeit erregt, aber nicht verunsichert oder beunruhigt (»angstfrei. Ich glaube, ich habe mich sogar darüber belustigt«). An anderer Stelle beschreibt Peter Fuchs sogar, dass das halluzinative Erleben mit einem Gefühl tiefer Gelassenheit einherging:

»PF: Mein Leben muss einfach so etwas wie schieres ›Leben‹ gewesen sein etwa in der Form, die Thomas Mann im ›Zauberberg‹ beschreibt:

›Was war das Leben? Man wusste es nicht. Es war sich seiner bewusst, unzweifelhaft, sobald es Leben war, aber es wusste nicht, was es sei [...].‹

MH: blieb von diesem merkwürdigen Empfinden denn irgendetwas zurück?

PF: Wenn ich ›zurückfühle‹, spüre ich, dass etwas erhalten blieb, eben diese kuriose Grundgelassenheit. Das war, wenn Du mir dieses Bild gestattest, wie ein Schweben oder Treiben im bernsteinfarbenen Wasser eines riesigen Aquariums, in dem nur die Sprudelsteine rastlos blubberten. The rest was a kind of silence, die Sache eines ungerührten Beobachters, der irgendwie resilient oder invulnerabel war – Nothing else matters.«²³

Laut der Schilderung von Fuchs vollzieht sich in diesem Stadium eine Reflexion auf das Erleben, jedoch keine Reflexion in Hinblick auf die Konsequenzen des Erlebten. Da ist noch kein Ich, das bedroht oder beunruhigt ist oder umgekehrt bestimmte Formen des Erlebens unbedingt haben möchte. In diesem Zustand, an dessen Grundgefühl sich Fuchs erinnern kann, gibt es noch keine Ausdifferenzierung zwischen *Reflexion-in-sich* und *Reflexion-in-anderes* und damit auch noch keine *Unterscheidung* zwischen *Erleben* und *Handeln*. Das willentliche Gestalten von bzw. das Eingreifen in Welt ist noch kein Thema und in diesem Sinne kann es in diesem Spiel der halluzinierten Formen auch kein Scheitern oder Misslingen geben, das bedrohlich sein könnte. Es gibt (noch) keine Vergleichsoperation, die ein Problem, ein Sollen oder ein Etwas-anderes-Wollen projizieren. Es ist einfach, was ist (»wie ein Schweben oder Treiben«), doch dieses ›ist‹ ist bereits reflexiv zugänglich, kann präsent gehalten und erinnert werden (»es war sich seiner bewusst, unzweifelhaft«).

23 Fuchs und Heidingsfelder (2022, S. 21 f.).

Aufbau von Zeitlichkeit

Etwas anderes wollen zu können als die unmittelbare Gegenwart setzt ein ›Selbstverhältnis‹ voraus, das sich von der unmittelbar gelebten Welt abspaltet, also nicht mehr identisch ist mit dem aktuellen Weltverhältnis. Wir können das Ich damit als eine Prozessionsstruktur ansehen, die sich in einer reflexiven Bewegung als Reflexion-in-sich von der Welt abgrenzt und sich über die Zeit hinweg, also über die Vergangenheit in die Zukunft projizieren kann.

Diese Struktur ist jedoch nicht per se gegeben, sondern bedarf zuerst der Erschaffung der zeitlichen Texturen, um sich selbst identifizieren zu können. Dies setzt ein Gedächtnis voraus, das zur Temporalisierung fähig ist und über Vergleichsoperationen zwischen einem Vorher und einem Nachher unterscheiden kann.

Es spricht einiges dafür, dass diese Voraussetzung erst dann erfüllt ist, wenn man in eine Sprache hinein sozialisiert ist, die aufgrund ihrer propositionalen Struktur ein solches ›Ich in der Zeit‹ bezeichnen und konstruieren kann.²⁴ Um ein diesbezügliches Selbst- und Weltverhältnis angemessen erzählen zu können, muss man in der Sprache sein. Ganz in diesem Sinne formuliert Fuchs mit Blick auf seine langsame (Wieder-)Aneignung dieses Prozesses:

»PF: Für mich war aber das Erzählen einer Vergangenheit auch zwischen Halluzinationen anfangs nur sporadisch möglich; wusste ich doch *kaum*, woher *ich* kam.«²⁵

Erzählungen generieren eine zeitliche Struktur (und setzen sie zugleich voraus), in der sich der Ich-Erzähler bewegen kann. Dies bedeutet, Erinnerungen gleichsam als Gegenstände auf einer Achse zu präsentieren, um sich hieran zeitlich orientieren zu können. Ohne diese Struktur könnte die Reflexion-in-sich des Subjekts wohl kaum stabile Werte bilden. Auch dies weist wieder darauf hin, dass das Ich erst in Koproduktion mit der Sprache jene Anhaltspunkte findet, mithilfe derer es sich über die Zeit hinweg als mit sich selbst identisch sehen und die hiermit einhergehenden Identitäten als wesentlich, wichtig und bedeutsam empfinden kann.

Erst hierdurch wird es für Fuchs möglich, die unterschiedlichen Halluzinationen mit ihren Erfahrungsqualitäten zu temporalisieren, um sie dann schließlich auch als Ereignisse bezeichnen und unterscheiden und damit in Beziehung setzen zu können. Schauen wir diesbezüglich noch auf eine andere Stelle:

»PF: Fast jede Halluzination, mit der ich zu tun hatte, gestattete keinen ›Übergang‹ zum nächsten Ereignis. Da gab es keinen Hiatus zwischen

²⁴ Zur Rolle der verkörperten Sprache siehe auch Di Paolo et al. (2018).

²⁵ Fuchs und Heidingsfelder (2022, S. 21 f.; kursiv im Original).

gerade eben noch und jetzt. Ich saß vielleicht einfach unschlüssig vor einem Teller mit Königsberger Klopsen. Essen oder nicht essen, das war die Frage. Erst später konnte ich im Rückblick Zeitgrenzen setzen: Vorhin war der Hund da, jetzt bin ich im Speisesaal, gleich werde ich im Raucherraum sein.

Aber in den ersten drei Monaten gab es kein Ende der Halluzinationen, keinen Übergang in das, was man ›Realität‹ nennt.«²⁶

Erst mit der Erzählbarkeit wird der Fluss des Erlebens als Kette bzw. Verkettung von Erfahrungen reflexiv zugänglich. Der Begriff der Erfahrung selbst impliziert bereits eine Bezugnahme, die nur im Nachtrag, nur im Reflex auf eine bereits vergangene Gegenwart geschehen kann. Erfahrung ist immer schon Differenzgeschehen, *Différance*, wie Derrida sagen würde.²⁷ Zu sagen: ›Ich habe etwas erlebt‹, setzt bereits jene zentrale Operation der Unterscheidung von Reflexion-in-sich und Reflexion-in-anderes voraus. Die Operation des Selbstbewusstseins scheint sich also darin zu manifestieren, diese Differenz zu erzeugen und sie als Struktur so stabil zu halten, dass immerfort zwischen diesen beiden Seiten oszilliert und diese Bewegung selbst reflektiert und gewahrt gehalten werden kann.

Indem verschiedene Erfahrungen als unterschiedlich markiert und erinnert werden können, besteht prinzipiell die Möglichkeit, sie auch in Hinblick auf ihre Qualität zu unterscheiden – etwa in Bezug auf den Realitätsstatus.

Konditionierte Koproduktion

Mit Blick auf eine Phase seiner Genesung, in der die Zweifel am Realitätsstatus der Halluzinationen bereits deutlich etabliert waren, erklärt Peter Fuchs an anderer Stelle:

»Manche solcher Ereignisse waren als stattgefundenere Realität immer noch präsent, als ich wieder zuhause war. Die Töchter und meine Frau haben so allmählich angefangen, mir zu sagen, dass dies alles nicht wirklich geschehen ist, was ich da aufgeregt oder wie auch immer erzählt habe. Es dauerte noch lange, bis ich den Nichtwirklichkeitsstatus anerkannt habe, und noch heute ist es so, dass ich bei manchen Sachen das Problem habe, mich fragen zu müssen, ob das real war. Ich war beispielsweise auf einem Bauernhof, da stand ich dann, schaute den Gänsen zu und hatte einen Moment lang das Gefühl oder besser den Gedanken: ›Ist das nun belastbar?‹«²⁸

²⁶ Fuchs und Heidingsfelder (2022, S. 34 f.).

²⁷ Derrida (2004).

²⁸ Vogd (2020, S. 241).

Gerade an dieser Stelle lässt sich wiederum feststellen, dass und wie Erfahrung, Erleben und Erzählen miteinander verschränkt sind. Erlebnisse werden signifikanten anderen Personen mitgeteilt und stoßen auf Resonanz oder Widerspruch. Die kommunikative Validierung führt zu einem Realitätscheck, der einen Teil der vergangenen Erfahrungen als ›Halluzination‹ indiziert, wenngleich dies der subjektiven Perspektive entgegensteht (»dauerte noch lange, bis ich den Nichtwirklichkeitsstatus anerkannt habe«). Es wird deutlich, wie stark die Unterscheidung zwischen Realität und Halluzination von der Operation der Kommunikation abhängt, also von dem Sinn, der im Medium der Sprache geteilt wird.

Was im Erleben als Gestalt, als Bewusstsein erscheint, ist unhintergebar positiv. Hier wird nicht zwischen Traum, Halluzination und Realität unterschieden. Es ist, was erscheint. Erst die Sprache bietet die Operation der Negation an und damit die Möglichkeit einer Reflexion, die zwischen Halluzination und Täuschung auf der einen und Realität auf der anderen Seite unterscheiden lässt. Erst hier kann die Positivität der der Seinsidentität sich selbst offenbarenden Formen fraglich und damit gebrochen werden. Die Negation verweist auf das Nichtsein. Sie gehört damit logisch gesehen nicht der Welt an, wie sie sich im zweiwertigen Bewusstsein darstellt. Sie gehört zum Bereich der Reflexion, ist also im Anschluss an Hegel und Günther im Bereich der Negativität zu verorten.

Nun ergibt sich aber, dass die Praxis der Negation – etwa das ›Nein‹-Sagen – ebenfalls zumindest in dem Sinne als positiv zu betrachten ist, als das artikulierte ›Nein‹ gehört oder in einer anderen Form wahrgenommen wird. Damit muss die Negation aus der Perspektive einer zweiwertigen Logik als eine Paradoxie erscheinen, nämlich als eine Operation, die gleichzeitig positiv wie auch negativ ist. Erst innerhalb einer mehrwertigen logischen Struktur, die über mehr als eine subjektive Position verfügt, ergibt sich die Möglichkeit, dass etwas an einem Ort als ›positiv‹, von einem anderen Ort aus gesehen jedoch zugleich als ›negativ‹ erscheint. Was aus der einen Position als unmittelbar gegebene Wirklichkeit erscheint, kann von einer anderen Position aus als ›Illusion‹ erscheinen und damit ›negiert‹ werden. Dies ermöglicht über die Operation der Kommunikation transjunktionale Operationen, mittels derer sich abgleichen lässt, was als konsensuelle Realität betrachtet oder was demgegenüber als ›subjektive‹ Perspektive oder gar als Täuschung aufgefasst werden kann – letzteres ist nämlich so zu verstehen, dass einer positiven Erfahrung gleichsam eine Negation angeheftet wird.

Wenn wir Sprache mit Humberto R. Maturana nicht als Übertragung von Information, sondern als ein Verhalten zur Koordination von Verhalten begreifen,²⁹ muss das Miteinander-Sprechen als eine Operation

29 »Die Sprache überträgt keine Information. Ihre funktionale Rolle besteht in der Erzeugung eines kooperativen Interaktionsbereiches zwischen Sprechern

verstanden werden, die ihrerseits in der Weise verkörpert wird, dass bestimmte sprachliche Figuren *sinulich* erinnert werden – und damit als Material genutzt werden können, um Identitäten auszuflaggen und zu kondensieren. In der Sprache sein – also Worte artikulieren und hören, die den Fluss des Erlebens unterscheiden, bezeichnen und damit in bestimmter Form interpunktieren – bedeutet demzufolge immer auch, die der Sprache eigenen polyphonen Strukturen zu inkorporieren. Es bedeutet, sprachlichen Sinn mit all seinen Perspektiven und Ambivalenzen zu fühlen und körperlich zu erinnern. Die über die sprachliche Kommunikation erfolgende Wahrnehmung einer Perspektive und Weltsicht, die der eigenen widerspricht, ist damit nicht einfach nur abstrakte Information. Der widerstreitende sprachliche Sinn schreibt sich vielmehr seinerseits in den Körper ein. Vor diesem Hintergrund können wir davon sprechen, dass der Mensch allein schon deshalb eine *polykontexturale Leiblichkeit* hat,³⁰ weil in ihm – vermittelt über die Kommunikation – unterschiedliche Perspektivierungen und Positionierungen mit jeweils eigenen körperlichen Regungen, Erinnerungen und Erwartungsmustern koexistieren. Hierdurch wird es möglich, dass der Fluss der sinnlichen Wahrnehmungen und die damit verbundenen erinnerbaren Objekte (auch die, die im Rahmen von Halluzinationen entstehen) in ein spezifisches Reflexionsverhältnis zu den sprachlichen Entitäten geraten. Die Negation greift – vermittelt über die Sprache – als positive Wahrnehmung in das Bewusstsein ein. Erst dadurch entsteht auch die Möglichkeit einer transklassischen Operation, entsprechend der die halluzinierte Welt kontrafaktisch zur Positivität ihrer Wahrnehmung rejiziert wird.³¹

Aus dem Blickwinkel einer Theorie der Polykontexturalität begegnen wir hier der konditionierten Koproduktion von Psyche, Körper und

durch die Entwicklung eines gemeinsamen Bezugsrahmens, auch wenn jeder Sprecher ausschließlich in seinem eigenen kognitiven Bereich operiert, in dem jede letztgültige Wahrheit durch persönliche Erfahrung bedingt ist. Da ein Bezugsrahmen durch die Klassen der Auswahlakte definiert wird, die er ermöglicht, kann das sprachliche Verhalten nur rational sein, d. h. determiniert durch Relationen der Notwendigkeit innerhalb des Bezugsrahmens, in dem es sich entwickelt.« (Maturana 1985, S. 40)

30 Siehe zum polykontexturalen Körper auch Fuchs (2005).

31 Dies legt zugleich den logischen Fehlschluss nahe, sich auf eine Metaphysik einzulassen, in der das Nichts oder die absolute Leere in einer Weise verdinglicht werden, dass sie nun als identisch mit dem absoluten Sein erscheinen, um hierdurch die operative Seite des Differenzierungsprozesses der konditionierten Koproduktion zu tilgen. Die Negation (bzw. ‚nein‘ zu sagen oder zu denken) stellt eine Operation dar; sie hat keine Seinsqualität an sich, darf also nicht ontisch missverstanden werden – auch wenn die Konsequenzen dieser Operation in einem positiven Sinne gefühlt werden können. Zur Elaboration am Beispiel der Rede vom Nirvana in den buddhistischen Lehren siehe Vogd (2017).

sozialen Systemen, die sich nur miteinander spezifizieren und ausdifferenzieren können. Die Formen, die die frei flottierende Psyche assoziiert und rekombiniert, werden in der neurosensorischen Matrix des Körpers erinnert und reproduziert sowie über die Kommunikation diszipliniert und in Hinblick auf brauchbare Weltmodelle konditioniert. Es bedarf eines Alter Egos (hier der Frau und der Töchter), das mit seinen sprachlichen Reaktionen ein Verhalten zur Koordination des Verhaltens zur Verfügung stellt, um Peter Fuchs dazu zu bewegen, ›intern‹ zwischen Halluzination und Realität zu unterscheiden.³² Erst die Kommunikation kann – vermittelt über die Sprache – etwas bewirken, was zunächst keine Positivität im Sinne einer unmittelbaren Wahrnehmung hat. Nur im Medium der Sprache kann die Reflexion als negativsprachlicher Bereich in der Welt instruktiv werden, das heißt im Erleben oder Handeln von Menschen einen Unterschied machen. Erst hierdurch wird es möglich, etwas wahrzunehmen (»Ich war beispielsweise auf einem Bauernhof, da stand ich dann, schaute den Gänsen zu«) und gleichzeitig den Realitätsstatus des Wahrgenommenen als fraglich zu behandeln. Sprachlicher Sinn kann die Sinneswahrnehmungen wenn auch nicht aufheben, so doch überschreiben bzw. übertönen.

Auf diese Weise ermöglicht die Kommunikation die Konfiguration und Moderation sinnlichen Erlebens, wenngleich sie selbst nicht sinnlich wahrnehmen kann (Sätze, auch wenn sie von sinnlichen Objekten oder fühlenden Wesen handeln, sind einfach nur Sätze; ein Text empfindet nichts, sondern besteht nur aus Worten und Buchstaben, die ihrerseits nicht wahrnehmen können). Dennoch hat die Kommunikation von Worten (etwa das Lesen von Texten oder die Wahrnehmung einer Stimme) eine Wirkung auf die Psyche. Sie ermöglicht etwa, dass Halluzinationen in ihrer Bedeutung abgeschwächt werden und langsam zu verblassen beginnen. Worte können aber umgekehrt auch Geister und Gespenster ins Leben rufen. Sie können Imaginäre erschaffen, die im suggestiblen Leib als reale Empfindungen und Fantasien ihre Verkörperung finden. Bereits die Rede von einer lustvollen oder gefährlichen Konstellation kann beim Zuhörenden ein Prickeln oder andere intensive Reaktionen auslösen. Auf performativer Ebene ist sehr wohl die Wirkung der Worte spürbar. Oder wie Merleau-Ponty formuliert: »Plötzlich merke ich, wie das Wort in meinen Körper hineinschnappt.«³³

32 Freilich ließe sich jetzt einwenden, dass die frei flottierende Psyche auch ohne Sprache konditioniert werden könnte. Es wäre dann in der Tat ein Forschungsprogramm, genauer zu untersuchen, wie Tiere, bei denen neuroelektrische Untersuchungen darauf hinweisen, dass sie träumen können, und die in der Kindheit ein ausgeprägtes Spielverhalten zeigen, zu den richtigen Sinnen kommen.

33 Merleau-Ponty (1974 [1966], S. 275). Für den Menschen spielt dabei wohl, wie Michael Tomasello herausgearbeitet hat, die Zeigegeste und die damit

Hiermit einhergehend leistet die Kommunikation eine Koppelung von Sphären, die ontologisch getrennt, logisch also nicht aufeinander rückführbar sind,³⁴ da sie im polykontexturalen Gefüge jeweils eine andere Subjekt- bzw. Beobachterposition besetzen. Dies gilt schon für die Ich-Du-Beziehung. Die subjektive Qualität der Erlebnisse und Weltansichten, von denen die Worte eines anderen Menschen sprechen, ist dem Hörer nicht zugänglich. Aus der Perspektive einer Theorie der Polykontexturalität steht die Kommunikation damit für transjunktionale Operationen, die die Koordination mit einer Du-Subjektivität ermöglichen, die aus egologischer Perspektive im absoluten Sinne transzendent bleiben muss.

Man könnte an dieser Stelle einwenden, dass die Sache der Einredung einer konsensuellen Realität nur deshalb funktioniert, weil ein Mensch wie Peter Fuchs bereits hinreichend in die Sprache einsozialisiert worden ist. Er verfügt also bereits über ein Gedächtnis, entsprechend dem die Rede der für ihn signifikanten anderen Personen eine existenzielle Bedeutung hat. Gedankenexperimentell ließe sich hier wieder auf ein Baby verweisen, das die Bedeutung der Worte noch nicht kennt. Doch dies

verbundene Triangulation von Zeichen, Bezeichnetem sowie Ich- und Du-Position eine besondere Rolle. Denn hiermit verschränken sich zugleich intentionale Haltungen im Sinne einer geteilten, wir-bezogenen Intentionalität (Tomasello 2009). Menschen prägen sich wechselseitig in dem Sinne, dass sich die intentionale Bewegung des jeweils anderen buchstäblich *in* die eigene *leibliche* Form einschreibt. Erst hierdurch wird jene für den Menschen typische soziale Klebrigkeit des In-der-Sprache-Seins möglich. Wir sprechen, um zu sein. Wir sind, weil wir sprechen. Die koordinierende Bewegung des sozialen Raums wie auch die hiermit einhergehenden wechselseitigen Imitationen überschreiben beim Menschen oftmals die unmittelbare Wahrnehmung und die aus ihr erwachsenden Intuitionen hinsichtlich kausaler Beziehungen. Soziale Skripten fühlen sich damit oftmals bedeutsamer (und folglich auch ›realer‹) an als das, was die unmittelbare visuelle Wahrnehmung logisch implizieren würde. Sowohl Schimpansen als auch Menschen können komplexe Handlungsabläufe von einem Vorbild imitieren. Beide können beispielsweise in einer bestimmten Reihenfolge fünf unterschiedliche Operationen an einer Box vollziehen, um im Anschluss eine Schublade zu öffnen, in der ein Bonbon liegt. Wenn dieses Experiment jedoch mit einer durchsichtigen Box durchgeführt wird, die den Mechanismus des Apparats offenbart, ergibt sich ein erstaunlicher Befund: Schimpansen lassen die überflüssigen Schritte weg. Demgegenüber folgen Kinder in der Regel den rituellen Schritten des Vorbilds, kopieren also unter Absehung von der sichtbaren Kausalität alle Schritte. Die ursprünglichen Bewegungen des sozialen Vorbilds bleiben instruktiv. Siehe diesbezüglich auch Horner und Whiten (2005).

34 Luhmann (1984) hat diesen Gedanken in seiner soziologischen Systemtheorie mit der These radikalisiert, dass psychische Systeme außerhalb sozialer Systeme stehen (also außerhalb der Kommunikation) und die Psyche demgemäß von kommunikativ operierenden Systemen als Umwelt behandelt wird.

würde an der zuvor benannten Konstellation nicht grundsätzlich etwas ändern. Der einzige Unterschied wäre, dass die konnotativen und denotativen Bedeutungen der Wörter, welche mit den Negationsoperationen einhergehen, erst noch aufgebaut werden müssen, damit sie erinnert werden können. Auch dies macht einen nicht unwesentlichen Teil der Sozialisation eines Kindes aus: das kreative Spiel seiner Assoziationen und Rekombinationen so zu bändigen, dass die Antworten im Korridor des altersgemäß angepassten Erwartungshorizonts seines (erwachsenen) Gesprächspartners liegen.³⁵

Zweifel gewinnen, um eine Realität aufbauen zu können

Die vorangehenden Ausführungen beschreiben das reflexive Bewusstsein als einen selbstreferenziellen Prozess, der sich in konditionierter Koproduktion mit der sprachlichen Kommunikation ausdifferenziert. Da ist kein inneres Seelenwesen, das nach draußen schaut und eine Realität erkennt. Alle Unterscheidungen und Bezeichnungen werden vielmehr durch den Prozess selbst aufgebaut. Damit wird aber auch deutlich, dass die Unterscheidung zwischen Halluzination, Traum und Realität nicht von sich aus gegeben ist. Sie muss ihrerseits aufgebaut, entwickelt und erinnert werden, um im Bewusstseinsprozess reproduziert werden zu können. Fuchs berichtet von einer Schlüsselszene mit einem Hund, die für ihn als Realitätstest instruktiv wurde:

»PF: Hab ich Dir erzählt von dem Hund, der mich gerettet hat?

WV: Nein.

PF: Also mich hat ein Hund aus den Halluzinationen »gezogen«. Das war so: Irgendwann konnte ich nicht nur herumliegen, aber ich war immer noch in dem Zustand einer sonderbaren Attentionalität, also in völliger Unklarheit darüber, warum ich mich hier aufhalte. Und dann war in meinem Zimmer plötzlich ein Hund neben meinem Tisch. Ich weiß noch genau, wie der aussah. Es war ein schwarzer großer Hund mit einer spitzen Schnauze, und der konnte grinsen. Und er war da. Auf einmal war der da, und dann war er plötzlich nicht mehr da. Und er kam immer wieder hoch in mein Zimmer, lächelte mich an, ging dann wieder raus und kam dann wieder. Genervt versuchte ich ihn zu vertreiben. Ich dachte, der muss weg. Zunächst ließ ich ihn einfach, bis er mir dann so auf die Nerven ging, dass ich versucht habe, ihn mit dem Rollstuhl zu überfahren,

35 Es ist nicht ungewöhnlich, dass ein fünfjähriges Kind auf Fragen wie »was ist ›2 mal 2‹ mit »grün« antwortet (Foerster u. Sander 1999) oder von Geistern berichtet. Man unterstellt ihm jedoch weder zu halluzinieren noch psychotisch zu sein, sondern korrigiert mit mehr oder weniger eindringlichem Ton seine Berichte, in der Erwartung, dass diese Sprechhandlung allein schon mit Blick auf ihre konnotativen Aspekte korrigierende Wirkungen hat.

ihn schlicht umzubringen. Aber er war durchlässig. Das war der erste Moment, in dem ich Verdacht geschöpft habe.

WV: Verdacht geschöpft, dass es eine Halluzination ist?

PF: Ja, Verdacht, dass da irgendwas nicht stimmt.«³⁶

Mit Günther gesprochen setzt bei Fuchs in der Auseinandersetzung mit dem Hund der Prozess der doppelten Reflexion ein. Mit dem »Verdacht, dass da irgendwas nicht stimmt«, beginnt sich der Prozess des Unterscheidens und Bezeichnens reflexiv auf sich selbst zu wenden. Der Unterscheidungsgebrauch und der daraus resultierende Aufbau einer Wirklichkeit werden selbst zum Thema und können in der Folge affirmiert, negiert oder angezweifelt werden. Logisch gesehen setzt dies voraus,

1. dass »die ursprüngliche Thematik ›Sein‹ festgehalten wird« (was in diesem Fall heißt, dass der Hund als Objekt der Außenwelt erkannt und erinnert wird);
2. dass »das Bewußtsein sich als Reflexion dieser Thematik von derselben absetzt« (was eine Person voraussetzt, die erinnern und von sich sagen kann, dass sie dies erlebt hat);
3. dass »eine weitere Reflexion den Gegensatz von 1. und 2. reflektiert«, also die Tatsache des Wahrnehmens und Erkennens selbst zum Gegenstand gemacht wird.³⁷

Das Bewusstsein erscheint erst dann vollständig ausgebildet, wenn es über die volle theoretische Reflexionsfähigkeit verfügt und damit in die Lage kommt, sich selbst kritisch zu hinterfragen. Dies wird möglich, sobald komplexere Operationen des Vergleichs durchgeführt werden, aufgrund derer Hypothesen über die Realität generiert werden, die dann getestet werden können. In diesem Fall lässt sich vermuten, dass Fuchs beim Versuch, sich mit dem Rollstuhl fortzubewegen, normalerweise auf Hindernisse stößt und damit bei ihm die Erwartung entsteht, dass sichtbare Objekte fest und nicht durchlässig sind und dass dies auch für lebendige Objekte gilt. Um solche Erwartungen ausbilden zu können, wird ein Gedächtnis bzw. eine Erinnerungsfunktion benötigt, die für die Vergleichsoperation die Information bereithält, dass man nicht ohne Widerstand durch Lebewesen fahren kann.

Zudem bedarf es eines Bewusstseinszustands, der bereits mit einer gewissen Unruhe und Eigenaktivität durchsetzt ist, also nicht mehr einfach nur selbstgenügsam vor sich hinvegetiert, sondern bestimmte Dinge zu ›wollen‹ beginnt (etwa: ›dass ein Hund verschwindet‹). Günther führt

³⁶ Vogd (2020, S. 240).

³⁷ Günther 2021 ([1957], S. 34).

diesbezüglich die Unterscheidung von »Kontemplation« und »Handeln« ein:

»[D]ie Reflexivität des Menschen manifestiert sich nicht nur als die stille Ideenwelt der Kontemplation. Sie setzt sich auch in Willen und die aus ihm folgende Handlung um. [...] Das mit seinen eigenen Begriffen spielende Denken mag an der Realität der Dinge zweifeln und ihr materielles Sein mag ihm als Trugbild erscheinen. Der Wille läßt sich dadurch nicht beirren. Sein handelnder Zugriff bestätigt ihm unmittelbar, daß die Dinge ›da‹ sind und daß sich das Sein ganz im Seienden erfüllt.«³⁸

Es bedarf der Absicht, den Hund zum Verschwinden zu bringen, und der Handlung, ihn mit dem Rollstuhl zu überfahren, um mit einem Feedback konfrontiert zu werden, dass die eigene Erwartung infrage stellt. Die Unterscheidung von Realität und Trugbild kommt erst im Versuch des handelnden Zugriffs ins Spiel, nämlich in der Erfahrung, dass die eigenen Konstruktionen nicht stimmen, was dann wiederum einen Erinnerungsaufbau und damit eine Gedächtnisfunktion voraussetzt, die die Ergebnisse der vollzogenen Operationen präsent hält bzw. reaktualisiert: »Jedes Tun ist Erkennen, und jedes Erkennen ist Tun«³⁹, formulieren Maturana und Varela, um die kreisförmige Beziehung des Aufbaus eines Reflexionsprozesses zu beschreiben, der nur auf den Ergebnissen seiner eigenen Aktivität beruht – also über keinen unmittelbaren Zugriff auf die Wirklichkeit verfügt.

In diesem Sinne lässt sich in der Entwicklung des Bewusstseinsprozesses von Fuchs, wie bereits angedeutet, ein qualitativer Sprung feststellen – sozusagen ein Wechsel des Musters, wie das Bewusstsein sich in Form bringt: Zunächst erscheint ein Fluss von Wahrnehmungen, in denen nicht zwischen Halluzination und Realität unterschieden wird. Dabei lässt sich annehmen, dass in diesen *stream of consciousness* auch weitgehend bruchlos Sinneseindrücke eingebaut werden können, die infolge des Kontakts mit Alltagsgegenständen oder anderen Lebewesen entstehen. Dies könnte sich als eine Art Driften darstellen, bei dem von Moment zu Moment irgendetwas erscheint und passiert, dabei aber nicht bestimmt zu werden braucht, ob es der Fantasie entspringt oder auf eine ›objektiv‹ bestehende Realität verweist.

Man könnte sich diesbezüglich gedankenexperimentell auch einen Menschen vorstellen, der in einer Art magischer Realität lebt: Waldgeister, Dämonen, Reisen in Parallelwelten und die Banalitäten des gewöhnlichen Alltags werden von ihm unterschiedslos als wirklich empfunden. In solch einer Welt zu leben, muss nicht per se mit Lern- oder gar Lebensunfähigkeit einhergehen. Dieser Mensch könnte mit der Zeit sogar Konstruktionen entwickeln, aus denen ein erfolgreiches und auch

38 Günther (2021 [1957], S. 89).

39 Maturana und Varela (1987, S. 32).

im ökologischen Sinne stabiles Selbst- und Weltverhältnis hervorgehen kann. Es würde ausreichen, zu lernen, welche Verhaltensweisen im Zusammenhang mit bestimmten Wahrnehmungen funktionieren.

Um Fantasiegebilde, Dämonen und Geister sowie stoffliche Lebewesen und materielle Gegenstände zu unterscheiden, bedarf es demgegenüber einer gesonderten Reflexionsleistung: Die Wahrnehmungen müssen typisiert sowie in Hinblick auf ihren jeweiligen Status evaluiert werden können und es muss die Fähigkeit vorhanden sein, das Ergebnis zu erinnern.

Hiermit wird nochmals deutlich, dass Reflexion nicht nur unterschiedliche Ebenen beinhaltet, sondern zugleich eine sprachliche Praxis darstellt. Das Problem muss erst in sprachlicher Form gefasst werden können (»Ich dachte, der muss weg.«), um im zweiten Schritt die auf die Problemlösung zielende Handlung zu erinnern (»dass ich versucht habe, ihn mit dem Rollstuhl zu überfahren, ihn schlicht umzubringen«) und dies dann im dritten Schritt evaluieren zu können (»Aber er war durchlässig.«). Erst auf dieser sprachlichen Ebene kann ein Reflexionsniveau erreicht werden, das zwischen Realität und Halluzinationen unterscheidet und im Sinne einer Limitationalität erlaubter Kausalbeziehungen⁴⁰ die hierfür notwendigen Kriterien benennen lässt.

Genau in diesem Sinne kann Fuchs dann auch davon sprechen, dass es der »Hund« war, »der mich gerettet hat«. Der Begriff »retten« kann mit Blick auf die vorangehenden Ausführungen nicht so gemeint sein, dass sein Leben an sich bedroht gewesen wäre. In einem von Halluzinationen geprägten Selbst- und Weltverhältnis lässt es sich, wie zuvor gezeigt wurde, durchaus einrichten. Die »Rettung« bezieht sich hier vielmehr auf die Existenz des *Beobachters*. Erst durch den Hund gewinnt Peter Fuchs in dem Sinne wieder sein Selbst, dass er nun zur doppelten Reflexion – zur Beobachtung zweiter Ordnung – fähig wird. Er gewinnt nun im Sinne von Hegel das »volle theoretische Bewusstsein« zurück.⁴¹

Schauen wir auf einen weiteren Ausschnitt aus der Erzählung seiner Rehabilitation:

»PF: Zum Teil war es dann so, dass ich, sobald dieser Verdacht aufkam, mich damit zufriedengab, nicht mehr meiner Wahrnehmung trauen zu müssen – auf Treu und Glauben verzichten konnte, wenn Du so willst.

⁴⁰ Der Begriff der Limitationalität ist von Luhmann (1998b, S. 392 ff.) im Zusammenhang der Ausarbeitung einer konstruktivistischen Wissenschaftstheorie eingeführt worden und bezeichnet die Praxis der Einschränkung von sagbaren Sätzen und Beziehungen durch Operationen der Negation. Um bei unserem Beispiel zu bleiben: Man benötigt erst die Einschränkung, dass Hunde nicht durchlässig sind, um dann die Aussage falsifizieren zu können, dass ein durchlässiger Hund ein echter Hund ist.

⁴¹ Zitiert nach Günther (2021 [1957], S. 56).

Ich erkläre das mal an einem Beispiel: Es gab einen Pfleger, der für mich immer ein Russe war, also von der Wahrnehmung her. Der kam dann abends und warf mich aus dem Bett in den Rollstuhl und dann warf er mich ins Bett zurück. Und ich fand das ungeheuerlich grob. Aber jetzt war ich schon fähig, mich zu fragen, ob das Wirklichkeit ist oder nicht, jetzt nicht in Bezug auf die Erfahrung, sondern in Hinblick darauf, ob das Folgen hat oder nicht. Was geschähe denn, wenn ich zur Pflegedienstleitung ginge und sagen würde: ›Tun Sie diesen russischen Pfleger weg. Er ist ein grober Klotz.‹ Und geantwortet würde: ›Den gibt es gar nicht, absolut nicht.‹ Und meine Überlegung war, dass ›Widerständigkeit‹ auftreten könnte – von Operationen gegen Operationen, wenn man Luhmann traut, der auf diese Weise das Zentralmerkmal der Beobachtung (Konstruktion) von ›Realität‹ beschreibt.

In solchen Passagen kam dann eben die Vermutung auf, dass ich in einem gewissermaßen ›fluiden‹ Zustand bin.«⁴²

An dieser Stelle werden nochmals einige Aspekte der konditionierten Ko-Produktion von Bewusstsein und sprachlicher Kommunikation deutlich. Offensichtlich ist der Erzähler mittlerweile wieder so weit mit dem Umgang der Sprache vertraut, dass er den Abgleich seiner Berichte mit den Wahrnehmungen und dem Wissen anderer Akteure erwartet – und sich dementsprechend der *Gefahr der Negation* bewusst ist. Er realisiert sich damit bereits als eine polykontexturales Wesen, die nicht nur von anderen Perspektiven weiß, sondern auch um die Wirkung von transjunktionalen Operationen, die die eigenen Berichte zurückweisen.

Neben der unmittelbaren Positivität der unmittelbaren Wahrnehmung (die noch nicht zwischen Halluzination und Realität unterscheiden lässt) besteht nun zum einen die Erinnerung an erwartbare Kausalitäten und darauf bezogene Wahrnehmungen (etwa der Festigkeit von Objekten). Zum anderen tritt jetzt die ›Realität‹ des In-der-Sprache-Seins in die Reflexionsbeziehung mit ein (›meine Überlegung war, dass ›Widerständigkeit‹ auftreten könnte‹). Sachverhalte werden jetzt in Hinblick auf ihre Erzählbarkeit reflektiert (etwa ob die eigenen Berichte vom Personal mit einer gewissen Wahrscheinlichkeit angenommen werden). Damit wird es möglich, etwas anderes als das, was man wahrgenommen hat, für real zu halten (›dass ich, sobald dieser Verdacht aufkam, mich damit zufriedengab, nicht mehr meiner Wahrnehmung trauen zu müssen‹).

Bemerkenswerterweise ist es nicht allein die Widerständigkeit gegen die Operation des materialen Handelns, die dem Sprecher jetzt die ›Realität‹ beibringt, sondern es sind auch die sprachlichen Operationen der Reflexionen von anderen, die (antizipierend) reflektiert werden. Es ist nun die Sprechhandlung selbst, die im polyphonen Gewebe vielfältiger

42 Vogd (2020, S. 240).

sozialer Positionen und der hiermit verbundenen Stimmen auf Widerstand stößt und zur Korrektur der Wahrnehmung auffordert.

Erneut ist an dieser Stelle darauf hinzuweisen, dass Sprache, Sprechen und die hiermit einhergehenden Reflexionsprozesse als leibliche bzw. gleichsam ›materiale‹ Phänomene zu verstehen sind. Für den, der hinreichend in der Sprache lebt, bekommen selbst vermeintlich sinnlose Worte Kraft, da die Beteiligten aufgrund der Geschichte ihrer Interaktionen ihr körperliches Verhalten durch Sprache koordinieren.

Wie Andreas Weber herausstellt, wirken sprachliche Bilder also nicht nur – wie man zunächst meinen könnte – als abstrakte Als-ob-Fiktion. Sie wirken vielmehr im Sinne einer konkreten Analogie, die das betrifft, was als Realität erfahren wird: »Hier hat ein ›scharfes‹ Geräusch oder ein ›tödliches‹ Wort dieselbe Bedeutung wie ein geschliffenes Messer. Beides ›zerschneidet‹ das Gelingen des Lebensvollzugs. Für den Körper müssen zwei auf diese Weise analoge Dinge identisch sein.« Identität heißt hier: Sie haben die »gleiche existenzielle Bedeutung für den Organismus. Das heißt folglich, dass sie aus dessen Perspektive das Gleiche sind und damit, dass nichtdiskursive (präsentative) Symbole sind, was sie bedeuten. Physisches und Psychisches ist auf dieser Ebene eins: zwei Perspektiven auf die eine autopoietische Fortexistenz des Lebewesens – einmal betrachtet unter dem Aspekt der materiellen Homöostase und einmal unter dem Aspekt der subjektiven Perspektive.«⁴³

Wenn aber die Identität auf einer Wahrnehmung beruht, die durch andere radikal infrage gestellt wird (»Den gibt es gar nicht, absolut nicht.«), und dies mit einer Wahrnehmung einhergeht, die identitätswirksam ist, dann muss die hierdurch entstehende kognitive Dissonanz in einer Weise gelöst werden, dass die eigene Identität nicht bedroht wird. In diesem Falle besteht die Lösung darin, den Konflikt nicht mehr zu suchen, indem die exaltierte Sinneswahrnehmung präventiv, das heißt antizipativ infrage gestellt wird.

Dies bedeutet, dass sich die Realität der sprachlichen Objekte gegenüber der Realität der unmittelbaren Wahrnehmung durchzusetzen beginnt und diese mit der Zeit sogar überschreiben kann (»meine Überlegung war, dass ›Widerständigkeit‹ auftreten könnte«). Doch dieser Prozess ist keineswegs selbstverständlich.⁴⁴ Zudem braucht er Zeit, also wiederholte Einredung durch signifikante Andere (»Die Töchter und meine Frau haben so allmählich angefangen, mir zu sagen, dass dies alles nicht wirklich geschehen ist, was ich da aufgeregt oder wie auch immer

43 Weber (2003, S. 120).

44 Es ist beispielsweise auch möglich, dass der Aufbau einer Beobachterposition im Reich der Sprache prekär bleibt (oder wird) und dann lieber den Halluzinationen getraut wird, wie es bei schizophrenen Patienten der Fall ist. Siehe etwa Sandsten, Zahavi und Parnas (2021).

erzählt habe. Es dauerte noch lange, bis ich den Nichtwirklichkeitsstatus anerkannt habe, und noch heute ist es so, dass ich bei manchen Sachen das Problem habe, mich fragen zu müssen, ob das real war.«).

Auch hier wird abermals deutlich, dass der Aufbau einer stabilen Realität für die frei flottierende menschliche Subjektivität erst in konditionierter Koproduktion gelingen kann. Es bedarf der sprachlichen Kommunikation mit der ihr eigenen Operation der Negation, damit eine anzweifelbare Version der Realität entstehen kann. Die doppelte Reflexion eröffnet das Spiel, die eigenen Konstruktionen als solche zu erkennen. Erst über die Du-Positionen vertrauter Anderer kann die hiermit einhergehende Kontingenz wieder eingefangen werden, nämlich indem mithilfe der Sprache eine geteilte Realität aufgebaut wird.

Freilich bleibt auf diesem Reflexionsniveau ein gewisser Zweifel bestehen (»Ist das nun belastbar?«). Doch ganz im Sinne der cartesianischen Tradition steht dieser Zweifel ab einem gewissen Punkt der Reflexion nicht mehr im Widerspruch zur eigenen Identität, sondern wird zu einem Referenzpunkt, an den sich das sprachlich verfasste Selbst halten kann. Auch wenn es keine weitere Gewissheit geben mag, besteht die Sicherheit des ›Cogito, ergo sum‹ fort.

Der Zweifel selbst erscheint damit als das definitive Haltesignal, auf das das Selbstbewusstsein referieren kann, um damit gleichsam eine ›philosophische‹ Identität zu gewinnen. Der Grund der Realität mag ins Wanken geraten und der Belastung durch Kommunikation und Reflexion lässt sich möglicherweise nicht mehr standhalten, doch bei alledem bleibt jetzt weiterhin noch die Tatsache des Plapperns des Geistes und des miteinander Redens bestehen – sobald diesbezügliche Muster der Reflexion hinreichend verkörpert sind.

Die vorangehenden Ausführungen machen allerdings deutlich, dass das Bewusstsein nicht als eine Entität im Sinne eines inneren Seelenwesens verstanden werden darf, das ›an und für sich‹ besteht. Es verdankt sich vielmehr einer konditionierten Koproduktion aus Körperlichkeit, Wahrnehmungsprozess, aktiven Handlungen und dem gemeinsamen In-der-Sprache-Sein. Selbst das Ich, das sich im zuvor genannten Sinne durch den Zweifel nährt, beruht auf Worten, die sich einer mit anderen geteilten kommunikativen Praxis verdanken.

Auf diesem Niveau der Ausdifferenzierung des Bewusstseins eröffnet sich jedoch die Möglichkeit einer Krise, die genau darauf beruht, dass diese Prozesse ins Stocken geraten und die aufgebaute Identität somit wieder bedroht wird.

Formkatastrophen und Sinnbrüche

Auf dem theoretischen Niveau der doppelten Reflexion angelangt, kann der Bewusstseinsprozess seinen eigenen Unterscheidungsgebrauch

thematizieren, um sich auf diese Weise im Prozess des Zweifelns einzuordnen und zu stabilisieren. Auf diesem Level der Reflexion ist das Medium Sinn so weit verfügbar, dass auch Unterscheidungen wie Realität und Halluzination, Gewissheit und Zweifel oder Sinn und Sinnlosigkeit sinnhaft getroffen und dementsprechend für den Sinngebrauch wie auch für den Identitätsaufbau genutzt werden können. Das Ich kann sich damit selbst unter fragwürdigen Bedingungen weiterhin reproduzieren, indem es die hiermit einhergehende Situation sinnhaft thematisiert und zumindest an den hiermit einhergehenden Bezeichnungen Halt findet. Solange das Denken und die Worte nicht versiegen, scheint dieser Prozess kein Ende zu finden. Damit scheint auch die Kontinuität des Ichs – des Surrogats der basalen Selbstreferenz des subjektiven Erlebens – gewährleistet.⁴⁵

Doch das aufgebaute Selbst- und Weltverhältnis kann – wie alles, was aufgebaut und konstruiert ist – auch wieder in eine Krise geraten, wenn die gefundene Form der Selbstidentifikation nicht mehr funktioniert oder ins Wanken gerät. Auch Fuchs berichtet von solchen Formkatastrophen:

»MH: Hast Du denn Beispiele für Formkatastrophen, die sich halluzinativ ereignet haben?

PF: Einmal blieb mir die Sprache (nicht nur die Worte) im Hals stecken. Das fühlte sich an wie Ersticken. Ich verfügte nur über gutturale Laute, also ›Kehlgeräusche‹. Der Name, dem ich diesem Zustand später gab, war: ›In der Kehle mumpfen‹.

Ich verstummte, was Sinn angeht, ich fand keinen, suchte ihn auch nicht. Ich erstickte bloß. Meine Frau hat dieses Ereignis miterlebt. Es hat sie tief getroffen.

In klassisch phänomenologischer Diktion: Der unbeobachtbare Körper meldete sich nicht als ›Leib‹, sondern als ›absolute Unvertrautheit‹. Plötzlich herrschte nur ›Grauen‹, diese Angst ohne ein ›Wovor‹ – ein deima panikon. Man könnte auch das Wort ›Entgeisterung‹ ins Spiel bringen. Eines seiner Synonyme ist ja ›sprachlos‹. [...] Eine andere Halluzination löschte mein Sprachvermögen. Ich wusste nicht, wie ›deutsch‹ geht und fand auch die Wörter nicht.«⁴⁶

Das Ich-Bewusstsein ist nicht per se gegeben, sondern entwickelt sich im Sinne der vorangehenden Ausführungen erst infolge eines Reflexionsprozesses. Sobald es in Erscheinung getreten ist, bedarf es eines kontinuierlichen Prozesses weiterer Reflexionen, um es aufrechtzuerhalten. Mit Blick auf die obigen Ausführungen beginnt dieser Prozess der Selbstwerdung mit der Generierung von Erinnerungen. Diese gestatten, dass sich das Selbst mittels Vergleichsoperationen über die Zeit hinweg als mit

45 So lässt es sich dann sogar im Sinnlosen einrichten, beispielsweise, indem man sagt: ›Nichts hat einen Sinn, ich bin Nihilist.‹

46 Fuchs und Heidingsfelder (2022, S. 35).

sich selbst identisch empfindet. Das Bewusstsein tritt also auf die Bühne, sobald das Spiel von Wahrnehmungen mit einer Gedächtnisfunktion einhergeht, die Erfahrungen als ›merkwürdig‹ und damit als erinnerbar kondensiert. In Koppelung mit der Sprache werden diese Erfahrungen erzählbar, wodurch der systematische Abgleich und die Verkettung mit anderen Erfahrungen möglich wird. Die eigene Existenz beginnt sich damit im In-der-Sprache-Sein zu erfahren. Das ›Ich‹ fängt an, sich selbst etwas zu bedeuten – und wird so von dem Sinn abhängig, in dem es sich nun eingerichtet hat.

Die Unterbrechung bzw. Suspension dieses Prozesses muss für das Ich-Bewusstsein folglich als eine fundamentale Krise, gleichsam als eine Art »metaphorischer Tod« erlebt werden.⁴⁷

Eine der Stärken der kybernetischen Perspektive besteht darin, uns zu zeigen, dass Geist und Sprache nicht als außerhalb der physischen Welt zu sehen sind, sondern dass Reflexion ihrerseits als verkörpert, als ein leiblicher Prozess zu verstehen ist. In diesem Sinne kann und darf das Selbstbewusstsein nicht mehr als ein unabhängiges Seelenwesen gedacht werden, sondern beruht, wie Michel Henry pointiert, buchstäblich darauf, dass das Wort Fleisch geworden ist.⁴⁸

Auch Peter Fuchs kommentiert die Krisenerfahrung seines Sprachverlusts in diesem Sinne:

»Vielleicht sollten wir schlicht davon ausgehen, dass Körper und Sprache streng gekoppelt sind. Ich würde heute wohl sagen: verschränkt.

Gemeint ist damit, dass Sprache (das Denken etc.) minimale Innervationen voraussetzt, ohne die weder sprechen noch schreiben funktioniert. Versuche einmal, das Wort ›Tonleiter‹ zu denken. Es wird sich ohne das winzige (körperliche) Beteiligungsein von ›Steigen, Absteigen ...‹ nicht präsentieren lassen.

Das ›Mumpfen‹ wäre die Metapher für das Unterbrechen jener Innervation, bedingt durch das neuronale System, in dem irgendwie die halluzinative ›Herumdaddelei‹ inszeniert wird. Deren Phänomenalität entsteht im Nachhinein – in einer Art ›Echo‹.⁴⁹

Wenngleich sprachlich generiert, da auf sozial angeliefertem Sinn beruhend, wird die Ich-Identität gerade deshalb als existenziell erlebt, weil sie körperlich gefühlt wird. Wenn die Worte im »Hals stecken[bleiben]«, fühlt sich dies an »wie Ersticken«. Die damit verbundene »Entselbstung«⁵⁰ führt nicht einfach zur Auflösung, sondern geht mit einer basalen Subjektivität einher. Sie ist körperlich erleb- und erinnerbar und

47 Bateson (1987, S. 161).

48 Henry (2011).

49 Fuchs und Heidingsfelder (2022, S. 35 f.).

50 Fuchs und Heidingsfelder (2022, S. 119).

phänomenalisiert sich als unangenehmer Körperzustand. Aus der Retrospektive erscheinen diese Wahrnehmungen als ein »Echo« des Verlusts an ›Ichhaftigkeit‹.

So, wie an Demenz erkrankte Menschen das Wegdriften ihrer Erinnerung (noch) nicht gelassen hinnehmen können, wenn sie den Vergleich zu vorher noch spüren, so scheint auch hier das Leiden am Ich-Verlust als Diskrepanz der Körperwahrnehmungen mitgeführt zu werden.

Auf der einen Seite ist das Ich als sprachlich-deskriptive Rekursion nicht mehr verfügbar. Auf der anderen Seite wird die mit dem Ich-Bewusstsein verkörperte Sinnbewegung jedoch noch mitgeführt. Gleichwohl kann sie nicht mehr in der Form ›dies ist *mein* Leib‹ angeeignet werden. Dies führt zu einer bedrohlichen Fremdheitserfahrung (›Der unbeobachtbare Körper meldete sich nicht als ›Leib‹, sondern als ›absolute Unvertrautheit‹. Plötzlich herrschte nur ›Grauen‹, diese Angst ohne ein ›Wovor‹«). Ab einer gewissen Komplexität und Ausdifferenzierung der Reflexionsfähigkeit kann also nicht ohne Weiteres zu einem ursprünglicheren, ›unschuldigeren‹ Zustand zurückgekehrt werden (etwa zum unkommentierten Erscheinen von Differenzen⁵¹). Auch in der Absenz von Sprache wird jetzt vielmehr die Erinnerung mitgeführt, dass da ein Beobachter sein müsste, der der Sache einen Sinn zu geben hat.

Das »Mumpfen« kann in diesem Sinne als eine Reminiszenz an den noch (vage) erinnerbaren und damit fühlbaren früheren Zustand verstanden werden – und demgemäß ist es nicht mehr ein unschuldiger Ausdruck körperlicher Regung, sondern zugleich eine Manifestation der kaum erträglichen Diskrepanz zwischen leiblicher Präsenz und entschwindendem Selbst.

Die geschilderte Agonie erinnert, wie gesagt, ein wenig an das Leiden eines Demenzkranken, der um die Unverfügbarkeit seines Gedächtnisses noch weiß, jedoch nichts daran ändern kann, dass ihm die eigene Biografie entgleitet. Dies ist eine vollkommen andere Bewusstseinsqualität als das Aufsteigen von Halluzinationen aus der Anfangslosigkeit, von dem zu Beginn dieses Kapitels berichtet wurde. Wir befinden uns hier nicht mehr bei der Unschuld des Anfangs, in dem das Bewusstsein langsam seine Gestalt und narrative Struktur gewinnt. Wir treffen vielmehr auf ein entwickeltes Reflexionsvermögen, das den Zugang zu den Bedingungen seiner eigenen Möglichkeit verliert und dem damit Sinn und Erzählbarkeit abhandenkommen.

51 Im Sinne einer komparativen Analyse ließe sich hier auf die Phase des Einschlafens, das Träumen, den Beginn einer Narkose oder auf einen Opiumrausch verweisen, in denen das Verblässen der Ich-Funktionen üblicherweise nicht als bedrohlich, sondern vielmehr als angenehm wahrgenommen wird. Es steht zu vermuten, dass der Körper in diesen Konstellationen etwas anderes macht.

Mit Maturana und Varela können wir unser »Ich-Selbst« als einen Prozess verstehen, der darauf beruht, im Fluss der Erfahrung eine kontinuierliche »deskriptive Rekursion« aufrechtzuerhalten,⁵² mittels der über Zeit und Raum hinweg die logische Kohärenz des Ich-Erlebens suggeriert wird. Dies muss als ein leiblicher Prozess verstanden werden, und entsprechend wird mit Blick auf die vorangehenden Ausführungen klar, worin das Problem bzw. worin das Leiden an den geschilderten Zuständen des Sinnverlusts besteht: Es ist nicht der Tod selbst oder die Bewusstlosigkeit – denn beides würde ja als Differenzlosigkeit keinen Unterschied machen und demzufolge auch nicht mit Leid oder Schmerz einhergehen können. Es ist vielmehr die Erfahrung einer Formkrise, also der Verlust einer Ichhaftigkeit, die man aus der Perspektive des verkörperten Gedächtnisses weiterhin mitzuführen glaubt. Es ist das Echo einer Identität, die an den Konsequenzen des eigenen Strukturaufbaus leidet.⁵³

Rückkehr

Die Ausführungen dieses Kapitels beschreiben, wie ausgehend vom Zustand der Bewusstlosigkeit wieder ein stabiles Selbst- und Weltverhältnis aufgebaut wird – zunächst über die Halluzinationen, dann über die Widerständigkeit in Bezug auf die Differenz von Kontemplation und Handeln und schließlich auf Basis des gemeinsamen In-der-Sprache-Seins. Als Ergebnis stabilisiert sich wieder ein Ich-Erzähler, der sich sinnhaft im Reich der Sprache verorten kann, damit jedoch insofern seine Unschuld verliert, als er das Stocken der eigenen Autopoiesis als Leiden erfahren muss.

Die Reise führt von der Bewusstlosigkeit zur doppelten Reflexion und wieder zurück.

52 Maturana und Varela (1987, S. 249 f.).

53 Analogien zu den hier dargestellten Zuständen finden sich in der Mythologie, so etwa im tibetanischen Totenbuch (siehe Thurman 2003), dessen klares Licht des Todes der Differenzlosigkeit entspreche. Den Zwischenzuständen der Bardos würden die Halluzinationen beginnender Selbst- und Weltverhältnisse entsprechen. Im Ergreifen von Zuständen, die als begehrenswert erscheinen, würde die Welt des Ich-Selbst beginnen. Der Reinkarnation eines Körpers entspreche das verkörperte Gedächtnis. Sofern wir uns von den mythologischen Implikationen lösen, entdecken wir hier eine Landkarte der menschlichen Psyche, die den Auf- und Abbau des Bewusstseins beschreibt.

III.2 DAS BEWUSSTSEIN DER GEHIRNE

Wir können und müssen auch das Gehirn als eine strukturdeterminierte Maschine betrachten: Die Organisation der Nervenzellen, ihre Vernetzung, die Prozesse an ihren Synapsen und die hiermit zusammenhängenden biochemischen Prozesse bestimmen, welche kognitiven Leistungen möglich sind. Das Nervensystem ist dabei als operational geschlossen anzusehen. Seine Aktivitätsmuster beruhen darauf, dass Nervenimpulse andere Nervenimpulse auslösen. Auf der Mikroebene des Gehirns finden wir Nervenzellen, die bei Aktivierung feuern und andere Nervenzellen aktivieren, die dann ebenfalls feuern. Auf der globalen Ebene des Gehirns lassen sich Netzwerke koordinierter Aktivitäten feststellen, die sich auf Basis der Hebb'schen Lernregel («what fires together, wires together»⁵⁴) ausdifferenziert haben. Wahrnehmung, Kognition, Aktivierung und Steuerung der Muskelbewegungen wie auch das Lernen geschehen auf dieser Basis. Sinneszellen übersetzen den sensorischen Input in elektrische Aktivitäten der Nervenzellen. Sie erzeugen dort gleichsam ein ›Tick‹, das in anderen Zellen ein ›Tick‹ auslöst, wodurch raumzeitliche Muster koordinierten ›Tickens‹ entstehen. In diesen Aktivitätsmustern sind die Kognitionen und das Verhalten eines Organismus codiert.

In Hinblick auf die Mikrostruktur (die synaptisch aktivierte Nervenzelle) und das große Bild (das Netzwerk von Aktivitätsmustern) verfügt die Hirnforschung bereits über recht viel Wissen. Sie ist jedoch noch weit davon entfernt, die Prozesse auf der mittleren Ebene im Detail zu verstehen, also konkret nachvollziehen zu können, wie etwa eine konkrete Gestalt – zum Beispiel eine Kaffeetasse oder das Gesicht der Großmutter – als solche erkannt werden kann.

Klar ist mittlerweile, dass die kognitiven Leistungen des Gehirns nicht auf Basis einer propositionalen Struktur realisiert werden. Das Gehirn greift ebenso wenig wie die heutigen künstlichen Intelligenzen auf ein Set logischer Aussagen zurück, die den zu erkennenden Gegenstand vollständig charakterisieren, so, wie es früher das *kognitivistische Paradigma* nahegelegt hatte.⁵⁵ Im Gehirn findet sich weder ein Großmutter-Neuron noch ein Großmutter-Modul, in dem die Eigenschaften abgespeichert wären, die das Konzept oder den Begriff ›Großmutter‹ definieren würden. Es gibt keine lokalisierbaren Speicherplätze im Gehirn, an denen all die logischen Elementarsätze und die Kette der kognitiven Schritte abgelegt wären, die beispielsweise als Ergebnis ›Tasse‹ ausgegeben würden, wenn zuvor die Register ›Porzellan‹, ›Behälter‹, ›schwarze Flüssigkeit‹ und ›Henkel‹ angeklickt wurden.

54 Hebb (1949).

55 Siehe hierzu paradigmatisch Minsky (1990).

Die kognitiven Prozesse im Gehirn beruhen vielmehr auf räumlich verteilten Netzwerken neuronaler Ensembles, die sich dynamisch bilden, auflösen und umgruppieren – und ad hoc von Moment zu Moment auf spezifische Weise in Resonanz zueinander treten. Solch ein Prozess ist wesentlich leistungs- und lernfähiger als ein explizites Gedächtnis, in dem die Merkmale der zu repräsentierenden Eigenschaften vorab festgelegt und mit dem diese verknüpft sind. Aus diesem Grund sind auch die früheren Versuche gescheitert, Computern dadurch Bildwahrnehmung und Sprachverstehen beizubringen, dass man Programme schreibt, die Erkennen, Verstehen und Sinnproduktion auf Basis semantisch-logischer Regeln zu rekonstruieren versuchen.

Erst im *konnektionistischen Paradigma* wird Mustererkennung in Bild, Ton und Sprache auf dem Niveau menschlicher Gehirne möglich⁵⁶ – nämlich indem man auf dem Computer neuronale Netzwerke simuliert. Letztere beruhen allein darauf, dass entsprechend der Hebbschen Regel gemeinsame Aktivitäten von Neuronen erinnert werden, indem Verbindungen gestärkt werden. Dies geschieht in der Simulation üblicherweise dadurch, dass jeder potenziellen Verknüpfung eine Wahrscheinlichkeit zugewiesen und nach jedem weiteren Lernvorgang das neue Set der Wahrscheinlichkeiten gespeichert wird. Auf diese Weise können im System Muster repräsentiert werden, ohne dass hierfür eine semantische bzw. präpositionale Liste der Eigenschaften vonnöten wäre, aus denen sich diese Muster zusammensetzen würden.

Die Bildererkennung auf Basis neuronaler Netzwerke muss also nicht ›wissen‹ bzw. benötigt kein Konzept davon, was eine gerade Linie, ein Kreis, ein Reifen, ein Auto, ein Fenster ist, um beispielsweise ein Fahrzeug als Mercedes typisieren und wiedererkennen zu können. Es ist hinreichend, mit einer großen Zahl von Bildern trainiert und durch Ja/Nein-Antworten auf das Ergebnis hin konditioniert worden zu sein, um ein Objekt zu ›erkennen‹.

Ebenso braucht eine künstliche Intelligenz bei der Sprachübersetzung – beispielsweise vom Russischen ins Deutsche – kein Lexikon, das dem russischen Begriff eine deutsche Bedeutung zuweist. Sofern das System mit einer hinreichenden Anzahl von Beispielen gefüttert wurde, arbeitet es, ohne überhaupt auf die lexikalische Bedeutung von Wörtern oder die propositionale Struktur von Sätzen und Begriffen zurückgreifen zu müssen. Es ist auch nicht darauf angewiesen, sich ein solches ›Wissen‹ durch Rückgriffe auf ein explizites Gedächtnis (beispielsweise in Form von Wörterbüchern, die auf einem Extraspeicher liegen) zu besorgen. Es ist hinreichend, neuronale Netzwerke auf mehreren unterschiedlichen Ebenen in geschickter Weise anzuordnen und zu verbinden, um diese Wahrscheinlichkeiten der Beziehungen von sprachlichen Ausdrücken

56 Siehe Pospeschill (2004).

berechnen zu lassen, auf deren Basis sie dann erwartbare Ergebniskombinationen ausgeben können.

Die Erfolge in der Spracherkennung und Sprachproduktion wie auch in der Bilderkennung und Bildproduktion (etwa im Deepfake) bestätigen das kybernetische Paradigma, das besagt, dass nahezu jede menschliche kognitive Fähigkeit in einer Maschine implementiert werden kann, sobald man weiß, wie die betreffenden Prozesse zu modellieren sind.⁵⁷ Dass dies etwa im Bereich der Sprache so gut funktioniert, weist umgekehrt darauf hin, dass nichts dagegen spricht, menschliche Gehirne als kognitive Maschinen aufzufassen, die zwar außerordentlich erstaunliche Leistungen vollbringen, aber nichtsdestotrotz strukturdeterminiert sind.

Doch welche Rolle spielt bei alledem das uns so intim vertraute Bewusstsein? Als Erklärung für die objektiven physikalischen Prozesse, die die zuvor erwähnten Kognitionen (zum Beispiel die Übersetzung vom Russischen ins Deutsche) ermöglichen, scheint es nicht gebraucht zu werden. Dies legt zunächst die erkenntnistheoretische Position des eliminativen Materialismus nahe, der die Frage nach dem Bewusstsein vollkommen ignoriert bzw. abstreitet, dass das Bewusstsein in kognitiven Prozessen überhaupt irgendeine Rolle spielt.⁵⁸ Das phänomenale Bewusstsein wäre damit als ein Epiphänomen neuronaler Prozesse anzusehen, würde jedoch kausal nichts zum Weltgeschehen beitragen.

In der Hirnforschung kommt Bewusstsein jedoch allein schon deshalb ins Spiel, weil man in der Regel auf Berichte aus der Erste-Person-Perspektive angewiesen ist. Man lässt sich von einer Versuchsperson berichten, was sie erlebt hat, um dies dann in Beziehung zu den gemessenen neuronalen Vorgängen zu setzen. In Hinblick auf viele Fragen kommt eine Hirnforschung, die sich nur am Tiermodell orientiert, an ihre Grenzen, denn Tiere verfügen nicht über die sprachlichen Möglichkeiten, um uns ihre inneren Zustände mitzuteilen.⁵⁹ Bei genauerem Hinschauen sind deshalb eigentlich recht viele Forschungsprojekte forschungspraktisch neurophänomenologisch orientiert, auch wenn sie selbst nicht unter diesem Label stehen. Damit das Experiment funktioniert, muss die Versuchsperson beispielsweise mitteilen, wenn sie eine Gestalt erkennt,⁶⁰ eine Willensbewegung spürt,⁶¹ Emotionen wie Angst empfindet,⁶² an sich selbst oder ihre Mutter denkt.⁶³ Neben der physisch-materiellen Denkform kommt hier also unweigerlich die phänomenologische Denkform mit ins Spiel. Zudem muss über die inneren Zustände kommuniziert und auch die

57 Günther (2021 [1957], S. 23 ff.).

58 Siehe etwa Churchland (1986).

59 LeDoux (2021, S. 339 ff.).

60 Varela (2003).

61 Libet et al. (1979).

62 LeDoux (1994).

63 Han und Northoff (2009).

Ergebnisse der Studie müssen schließlich in sprachlicher Form dokumentiert und veröffentlicht werden. Damit ist neben der phänomenologischen auch die sprachliche Denkform unabdingbar. Man mag noch so materialistisch orientiert sein – ob man will oder nicht, jedes Forschungsergebnis verdankt sich phänomenologischer und semantischer Prozesse.⁶⁴ Ohne Sinn und Sinnlichkeit geht es nicht. Die wissenschaftlichen Beobachter haben zudem ihre eigene Selektivität, wie sie sich in den Forschungsprozess einbringen – allein schon dadurch, dass sie bestimmte Fragen stellen (und andere nicht) und damit die Welt in einer bestimmten Weise anschneiden.

Auch die Hirnforschung findet dementsprechend in einer polykontexturalen Welt statt. Unweigerlich hat sie unterschiedliche Beobachter (und deren subjektive Perspektiven) in Rechnung zu stellen – selbst wenn sie die physisch-materielle Denkform in idealisierter Form (das heißt unter Abblendung der beteiligten subjektiven Vorgänge) aufrechterhalten will.

Neuronale Korrelate des Bewusstseins

Die vorangehend geschilderten erkenntnistheoretischen Lagerungen hindern uns jedoch nicht, nach den neuronalen Korrelaten des Bewusstseins zu fragen. Sie bilden vielmehr die Voraussetzung, um auch diese Frage systematisch anzugehen.

Das Problem lässt sich mit unterschiedlichen methodologischen Zugängen angehen. Wie bereits angedeutet, kann etwa untersucht werden, welche Gehirnaktivitäten auftreten, wenn etwas Bestimmtes erlebt wird. Umgekehrt lässt sich fragen, welche bewussten Kognitionen und Aktivitäten ausfallen, wenn bestimmte Strukturen wegfallen, etwa infolge von Unfällen oder Tumoren. Interessant in diesem Zusammenhang ist etwa das Phänomen der »Blindsight«: Es gibt Menschen, die einen Teil ihres Gesichtsfeldes nicht mehr bewusst erleben, aber weiterhin fähig sind, auf dort erscheinende Objekte zu reagieren (etwa eine Fliege zu verjagen oder nach einem anderen dort erscheinenden Objekt zu greifen).⁶⁵ Mit Blick auf den Bericht von Peter Fuchs aus dem vorherigen Kapitel scheinen selbst sprachlich-kommunikative Funktionen ohne Beteiligung von Bewusstsein ablaufen zu können.

Umgekehrt gibt es jedoch eine Reihe von Hirnschädigungen, die *nicht* mit einer Beeinträchtigung der Bewusstseinsfunktionen einhergehen. So ist bekannt, dass das Kleinhirn wesentlich mehr Nervenzellen als die Großhirnrinde hat und zentral an sensomotorischen Lernvorgängen beteiligt ist. Diese Prozesse laufen jedoch implizit ab, das heißt ohne dass es überhaupt zu einer bewussten Aufmerksamkeit kommt. Zudem geht

64 Zu den drei Denkformen in der psychologischen Forschung siehe auch Laucken (2003).

65 Siehe Ramachandran (2005, S. 37 ff.).

selbst eine massive Schädigung des Kleinhirns nicht mit einer fundamentalen Beeinträchtigung der Bewusstseinsfunktionen einher.⁶⁶

Auch sprechen viele Forschungsergebnisse dafür, dass nicht alle Bereiche des Großhirns an Bewusstseinsaktivität beteiligt sind. Untersuchungen zur binokularen Rivalität⁶⁷ – man zeigt dem rechten Auge beispielsweise eine Katze und dem linken Auge einen Hund – führen zu dem Befund, dass nur eine der beiden Gestalten im Bewusstsein erscheint, im primären visuellen Kortex aber sehr wohl die Informationen beider Bilder verarbeitet werden. Wie bereits im Falle der Blindsight aufgezeigt, muss sich nicht jede visuelle Wahrnehmung bzw. ihre Verarbeitung im Bewusstsein ausdrücken.

Wenn jedoch Hirnregionen aus der im hinteren Kortex gelegenen sogenannten »hinteren heißen Zone« – diese erstreckt sich über Teile des Schläfen-, des Scheitel- und des Hinterhauptlappens⁶⁸ – verletzt werden, sind auch die phänomenologischen Bewusstseinsleistungen des Sehens, Hörens, Riechens, Schmeckens, Empfindens und Denkens beeinträchtigt oder fallen gar vollkommen weg. Werden diese Teile demgegenüber – etwa während einer Hirnoperation – stimuliert, können Sinneseindrücke erinnert werden.

Wir finden also einerseits eine Situation vor, entsprechend der der menschliche Organismus mit seinem Gehirn recht gut ohne Bewusstsein zurechtkommen kann – was von einer strukturdeterminierten kognitiven Maschine auch zu erwarten ist. Andererseits taucht jedoch in bestimmten Lagerungen ›Bewusstsein‹ auf – und dies geht mit spezifischen, in Hinblick auf die beteiligten Hirnregionen identifizierbaren Korrelaten einher (üblicherweise bezeichnet als *neural correlates of consciousness*).

Dies wirft die Frage auf, was konkret sich anders darstellt, wenn das Bewusstsein im Spiel ist bzw. seine neuronalen Korrelate aktiv sind. Auf der phänomenologischen und semantischen Ebene ist die Sache zumindest auf den ersten Blick klar: Es kann über Erfahrungen berichtet werden, die im Nachtrag auf ein subjektives Erleben erinnert und bewusst gehalten werden können. Grammatisch wird dies durch Sätze der Art ›Ich habe etwas erlebt‹ indiziert.

66 Christoph Koch (2020, S. 14) bemerkt hierzu: »Was passiert mit dem Bewusstsein, wenn Teile des Kleinhirns durch einen Schlaganfall oder durch das Messer eines Chirurgen verloren gehen? Sehr wenig! Die Betroffenen beklagen zwar zahlreiche Beeinträchtigungen – sie können etwa nicht mehr flüssig Klavier spielen oder auf der Tastatur tippen –, jedoch nie den Verlust einer Facette ihres Bewusstseins. Sie sehen, hören und fühlen normal, sie behalten eine Ich-Empfindung, sie erinnern sich an Vergangenes und planen die Zukunft. Selbst wer bereits ohne Kleinhirn auf die Welt kommt, leidet kaum unter einem eingeschränkten Bewusstsein.«

67 Siehe etwa Engel et al. (1999).

68 Siehe Koch (2020, S. 15).

Freilich offenbart sich gerade die phänomenologische Erfahrung bei genauerer Betrachtung als ein recht komplexer Prozess. Wie Husserls Untersuchungen zur Zeitstruktur des Erlebens aufzeigen, erscheint ein über etwa drei Sekunden ausgedehntes ›Jetzt‹. Innerhalb dieses ›Erlebnisraums‹ verschränken sich der Rückblick auf zuvor Geschehenes (Retention) und die Vorausschau auf künftig Erwartbares (Protention) zu einer intentionalen Bewegung.⁶⁹ Zudem verwebt sich dies mit dem für den Menschen typischen In-der-Sprache-Sein. Bei alledem ist unbestreitbar, dass Bewusstsein für uns einen unmittelbaren Unterschied macht: Die Welt wird von einem subjektiven Zentrum aus als ›meine Welt‹ erlebt.

Auch auf neurologischer Ebene lassen sich Unterschiede benennen. Das Kleinhirn, das für die Bewusstseinsfunktion keine Rolle spielt, ist in einer Weise verschaltet, die kaum Rückkoppelungsprozesse ermöglicht, mit denen neuronale Aktivitäten hin und her kreisen könnten. Es ist vielmehr so aufgebaut, dass eine Vielzahl von Modulen parallel arbeitet und dabei nur in geringem Umfang miteinander interagiert.⁷⁰

Demgegenüber weisen die Regionen des Großhirns, in denen die neuronalen Korrelate des Bewusstseins zu finden sind, eine andere Form der Verknüpfung auf: Unterschiedliche Aktivitätsmuster verschiedener Bereiche werden über eine »reziproke Koppelung« in Form eines »Re-entry« miteinander verbunden.⁷¹ Diese Art der Vernetzung ist mehr als eine einfache kybernetische Rückkoppelungsschleife, über die ein negatives oder positives Feedback zurückgestrahlt wird, um die Aktivität zu moderieren. Vielmehr verlaufen diese Prozesse parallel und in zwei Richtungen gleichzeitig. Dies führt dazu, dass kontinuierlich neue Dynamiken entstehen, die mit Mustern einhergehen, die vom System ihrerseits wieder an eine andere Stelle projiziert und eingespielt werden können. Lassen wir Gerald M. Edelman zu Wort kommen, der gemeinsam mit Giulio Tononi die auf diesen Befunden beruhende Theorie des Bewusstseins prominent gemacht hat:⁷²

»Im Verlauf von Entwicklung bilden sich sowohl zwischen benachbarten als auch zwischen weit auseinander liegenden Neuronen zahlreiche

69 Husserl (2000); siehe dazu auch Schnell (2002).

70 »Die gesamte Struktur des Kleinhirns ist demnach für das Bewusstsein ohne Belang. Warum? Die Antwort liegt möglicherweise in seiner außerordentlich gleichförmigen Verschaltung. Im Kleinhirn finden wir einen fast gänzlich vorwärtsgerichteten Informationsfluss: Eine Neuronengruppe leitet Signale an die nächste, welche sie an eine dritte weitergibt. Komplexe Rückkopplungsschleifen, in denen elektrische Aktivität hin und her kreist, fehlen.« (Koch 2020, S. 14 f.)

71 Edelman (2004, S. 49 f.).

72 An dieser Stelle ist darauf hinzuweisen, dass es alternative Theorien zur Entstehung des Bewusstseins gibt. Prominent ist etwa das Global Neuronal Workspace Model of Conscious Access (Dehaene, Changeux & Naccache 2011), das von der Annahme ausgeht, dass es im Gehirn so etwas wie eine

reziproke Verbindungen heraus. Sie ermöglichen die Signalübertragung zwischen *kartierten Arealen*. Als Reentry bezeichnen wir den fortlaufenden rekursiven Austausch paralleler Signale zwischen Hirnarealen, der dazu dient, ihre Aktivitäten in Raum und Zeit zu koordinieren. Beim Reentry handelt es sich nicht, anders als bei einer Rückkoppelung, um die sequenzielle Übertragung eines Fehlersignals in einer einfachen Schleife. Vielmehr sind daran viele parallele reziproke Pfade beteiligt, und es gibt keine mit vorgegebenen Werten operierende Fehlerkorrektur. «⁷³

Bemerkenswert an diesem Prozess ist, dass sich hierüber ein Gedächtnis der Gegenwart aufbauen lässt, das durch die eigenen Operationen immerfort modifiziert wird. Die theoretische Modellierung als reentranter Prozess schmiegt sich damit an die phänomenologische Erfahrung an, die eine ausgedehnte Gegenwart präsentiert,⁷⁴ in der auch zeitliche Strukturen (wie etwa eine Tonfolge oder eine Bewegung) als Gegenstand des Bewusstseins präsent gehalten werden können. Um hier nochmals eine bereits weiter oben zitierte Passage von Gerald M. Edelman aufzugreifen:

»Neue Wahrnehmungskategorien werden reentrant mit Gedächtnissystemen verkoppelt, ehe sie selbst Teil eines nun veränderten Gedächtnissystems werden. Anhand des Gedächtnisses kategorisierte Wahrnehmungen werden also rekursiv genutzt, um das Gedächtnis selbst zu modifizieren. Diese Wechselwirkungen laufen, so nimmt man an,

›Informationstafel‹ gibt, auf der die im Bewusstsein erscheinenden Inhalte abgelegt werden und auf die dann unterschiedliche Hirnareale mit ihren Prozessen zugreifen können. Einem ähnlichen Ansatz folgt Thomas Metzingers Selbstmodell-Theorie, nach der das Bewusstsein eine analoge Projektion (das subjektive Selbstmodell) der digital arbeitenden neuronalen Prozesse ist, durch die eine neue Qualität von Informationen entsteht, die die neuronalen Prozesse ihrerseits wieder nutzen können (Metzinger 1999). Subjektivität erscheint damit als eine luzide Illusion ebenjener Modellprojektion und kann sich ihres eigenen illusionären Charakters auf der Erfahrungsebene nicht bewusst werden. Ein wenig überzeichnet können wir uns beide Theorieangebote zur Entstehung des Bewusstseins etwa folgendermaßen vorstellen: Ein physikalisch operierendes Rechensystem projiziert qua physikalischer Technik (zum Beispiel mit einem Video-Beamer) analoge Bilder und Töne, um diese dann qua physikalischer Technik (zum Beispiel Mikrofonen und Fotochips) aufnehmen und analysieren zu können. In einem ähnlichen Sinne lässt sich auch die Attention Schema Theory von Graziano (2019) lesen, die jedoch darüber hinaus den Fokus auf die Prädiktion fremden Verhaltens legt. Die eigentliche Funktion des Selbstmodells liegt damit darin, das Verhalten anderer Menschen und von Tieren vorhersagen zu können.

73 Edelman (2004, S. 49 f.).

74 Siehe hierzu ausführlich Varela (1999), dessen Überlegungen auf einer empirischen Untersuchung auf Basis einer ähnlichen Modellierung beruhen.

innerhalb von Zehntelsekunden bis Sekunden ab, also innerhalb der Zeitspanne, die William James ›Scheingegenwart‹ nennt. Ich bezeichne sie als ›erinnerte Gegenwart‹, um zu unterstreichen, dass aus der Interaktion zwischen Gedächtnis und aktueller Wahrnehmung das Bewusstsein entspringt.

Was für Folgen hatte es, dass die Evolution eine dynamische Koppelung zwischen Werte-Kategorien-Gedächtnis und Wahrnehmungsgedächtnis herstellte? Es bildete sich die Fähigkeit heraus, eine komplexe Szene zu konstruieren und zwischen Bestandteilen dieser Szene Unterscheidungen zu treffen.«⁷⁵

Aus den hier vorgestellten Überlegungen lässt sich die These ableiten, dass eine netzwerkförmig organisierte neuronale Struktur umso mehr an ›Bewusstseinsqualität‹ hervorbringt, je mehr und je vielfältigere Information sie in sich integrieren kann. Letzteres heißt, dass sie multiple Kartierungen wechselseitig aufeinander projizieren und hierdurch immerfort neue Re-entries produzieren kann – also Muster, die darauf beruhen, das Formen in immer neuer Form wieder in sich selbst eintreten. Bewusstsein wäre damit eine intrinsische Eigenschaft von neuronalen Prozessen, die aufgrund ihrer spezifischen Organisation genau dies leisten. »Tononi«, so führt Koch aus, postuliert, »dass jedes komplexe Netzwerk, das in seiner Struktur Zusammenhänge zwischen Ursache und Wirkung codiert, diese Eigenschaften besitzt und so ein gewisses Maß an Bewusstsein mitbringt. Wenn es einem System aber, wie dem Kleinhirn, an Verschaltung mangelt, wird es nichts bewusst wahrnehmen.« Entsprechend der integrierten Informationstheorie »verfügt ein System über umso mehr Bewusstsein, je mehr Information es in sich integrieren und vielfältig verarbeiten kann.«⁷⁶

Wird ein solches System gestört oder gereizt, zum Beispiel durch Sineseeindrücke (bzw. die dadurch ausgelösten Nervenimpulse) angeregt, so wird es aus seiner aktuellen dynamischen Struktur heraus die nächste Form finden. Das System wird den Input jedoch nicht mechanisch auf Basis vorab bestehender Programme verarbeiten, wie es beim Kleinhirn der Fall ist, das komplizierte Bewegungsabläufe (etwa das Klavierspielen) entsprechend der durch frühere Lernvorgänge festgelegten Algorithmen steuert. Vielmehr erzeugt das für das Bewusstsein verantwortliche System eine neue Form.

Die Form des Re-entry

Der Begriff *Re-entry* wurde ursprünglich von dem Mathematiker George Spencer Brown im Rahmen seiner Abhandlung *Law of Forms*

75 Edelman (2004, S. 63 f.).

76 Koch (2002, S. 17).

eingeführt.⁷⁷ Mit dem Formkalkül leistet er eine formale Theorie des Auf- und Abbaus von Welt, die nur eine einzige Operation benötigt: das Unterscheiden und Bezeichnen. Der Neurobiologe Francisco J. Varela hat den Kalkül weiterentwickelt, um selbstreferenzielle Prozesse in der Hirndynamik modellieren und auf dieser Basis empirisch untersuchen zu können.⁷⁸ Schon aus diesem Grund ist es sinnvoll, die Formtheorie an dieser Stelle etwas ausführlicher vorzustellen.

Der Formkalkül beginnt mit dem Treffen einer Unterscheidung (»Draw a distinction.«). In formaler Hochabstraktion und zugleich als grafische Illustration wird dieser grundlegende Akt durch eine Markierung, genauer: einen *Haken*, auf einer leeren Fläche ausgedrückt.

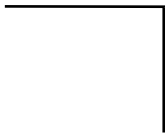


Abb. 8: *Draw a distinction*. Der Haken steht für das Treffen einer Unterscheidung und Bezeichnung. Die Figur ist asymmetrisch, da das Bezeichnete auf der Innenseite und der unmarkierte Raum auf der Außenseite steht.

Der Haken von Spencer Brown steht einerseits für die erste Operation, das Unterscheiden, also den systemischen Zusammenhang, der einen Schnitt in die Welt einführt, sodass aufgrund spezifischer, jeweils konkret zu bestimmender Operationen das eine (etwa das System) von etwas anderem (etwa seiner Umwelt) getrennt erscheint. Infrage kommt hier – so Fritz Simon – jede »Operation, durch die ein Raum, Zustand oder Inhalt (= eine Welt) geteilt wird«. Es ist hinreichend, wenn durch die »Operation des Unterscheidens« ein »Raum oder Inhalt« entsteht, »der innerhalb der so entstandenen Einheit verortet ist, und ein Raum, Zustand oder Inhalt, der außerhalb dieser Einheit lokalisiert ist (= Innenaußen-Unterscheidung)«⁷⁹. Mit dem Bezeichnen wird nun andererseits eine »zweite Operation des Unterscheidens« mit einem vorhergehenden »Unterscheiden gekoppelt und als Verweis (= zeigen) auf das erste Unterscheiden gebraucht«⁸⁰. Die Außenseite des Hakens (und damit der Markierung) erscheint als unmarkierter Raum (unmarked space).

Dass Unterscheiden und Bezeichnen hier in einer Operation zusammenfallen, lässt sich auch so verstehen, dass keine weitere Repräsentation, kein zusätzliches Zeichen notwendig ist, damit der Haken für den weiteren Prozess informativ wird. Wenn beispielsweise ein Ensemble von Nerven zu feuern beginnt, ist dies zugleich eine Unterscheidung (es geschieht dieses und nicht jenes) und die Bezeichnung (das Geschehen ist

77 Spencer Brown (1994 [1969]).

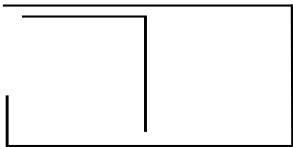
78 Varela (1979, 1999).

79 Simon (2018, S. 13).

80 Simon (2018, S. 15).

zugleich die Kognition). Die Nervenaktivität markiert einen Unterschied und wird damit für das System zur Information – also zu einem Unterschied, der als Bezeichnung einen weiteren Unterschied machen kann.

Mit dem Re-entry formuliert Spencer Brown auf dieser Basis eine spezielle Operation, entsprechend der die Unterscheidung von Markierung und unmarkiertem Raum markiert wird und als Unterscheidung in sich selbst wieder eintritt:



unmarked
space

Abb. 9: Re-entry und unmarked space. Im Re-entry tritt die Unterscheidung von Bezeichnung und unmarkiertem Raum rekursiv in sich selbst ein, was grafisch durch die beide Seiten umschließende Figur ausgedrückt wird. Neben dem durch diesen Prozess bezeichneten inneren unmarkierten Raum steht außerhalb der Grafik weiterhin der unmarkierte Raum, der nicht durch Bezeichnung und Unterscheidung thematisiert wird.

Der sich schließende Re-entry-Haken verweist auf eine Selbstreferenz. Die Operation des Unterscheidens und Bezeichnens wird auf sich selbst zurückverwiesen. Hierdurch entsteht eine Dynamik, die unweigerlich eine Art Identität erzeugt, die darauf beruht, dass sich die Unterscheidung von Innen und Außen selbst verfügbar wird und durch den Re-entry kontinuierlich reproduziert werden kann. Freilich kommt hiermit ›Zeit‹ ins Spiel. Wir haben nicht nur das Bild einer Struktur, die durch die Anfangsunterscheidung gesetzt ist, sondern einer Unterscheidung, die kontinuierlich zwischen den zwei unterschiedenen Seiten oszillieren kann, indem erneut Operationen vollzogen werden, welche die Form reproduzieren. Wichtig ist jedoch, dass sich auf der Außenseite der Re-entry-Form wiederum der unmarked space befindet. Die Leere des unmarkierten Raums besteht weiterhin in einer für die Innenseite der Form nicht verfügbaren Weise fort.⁸¹

Die Prozessdynamik lässt sich gut am Beispiel einer Saitenschwingung illustrieren. Die Saite wird vom Gitarristen an einer Stelle angeschlagen. Dies entspricht dem Haken der ersten Unterscheidung und Bezeichnung. Nun aber läuft eine Wellenbewegung nach vorn, trifft auf den Steg, wird zurückgeworfen und läuft auf den Bund zu, wo sie erneut zurückgeworfen wird. Die Bewegung tritt – solange die Saite schwingt – immer wieder

81 Es ist klar, dass auf diese Weise auch die Subjekt-Objekt-Unterscheidung fluide wird bzw. jeweils nur situativ stabilisiert werden kann. Dirk Baecker (2021, S. 9) entwickelt auf dieser Basis eine Theorie der Katjekte.

in sich selbst ein und generiert dadurch jene Oszillationen, die von uns als Ton hörbar werden. Die jeweils von Moment zu Moment erscheinende Schwingung erscheint der Form nach als dieselbe, vom Prozess her gesehen jedoch jeweils als eine neue Schwingung, da sie erst durch den nachfolgenden Wiedereintritt der Form in die Form entstanden ist. Im gleichen Sinne kann die Aktivität eines Ensembles von Neuronen auf ein zweites Ensemble treffen und dort einen Unterschied machen, der dann wieder auf das erste Ensemble zurückgeworfen wird, sodass ein Muster entsteht, das der Form nach dasselbe bleibt, sich mit Blick auf den Prozess der kontinuierlichen Reproduktion jedoch in seiner Gestalt zugleich auch ändern kann. Sofern der Prozess nicht unterbrochen wird, bleibt die Figur des Re-entry bestehen, auch wenn sich das Aktivitätsmuster mit der Zeit wandelt.⁸² Das ist dann gleichsam so, als ob an einem Froschteich zunächst Stille herrscht, dann ein Frosch zu quaken beginnt und dann alle anderen ihrerseits zu quaken anfangen, wodurch ein angeregtes Konzert zu vernehmen ist, das auf Basis der sich selbst stimulierenden Interaktionsdynamik immerfort neue Klangmuster hervorbringt, bis der Prozess irgendwann wieder verebbt.

Sobald erst einmal eine reentrante Form in Gang gesetzt worden ist, reproduziert sie sich als ein Prozess, der zwischen der ersten Markierung und der Tatsache, dass markiert wurde, oszilliert. Da ist jetzt nicht nur eine Aktivität, die bezeichnet (der Haken), sondern eine rücklaufende Aktivität, die den Unterschied zwischen Bezeichnung und dem Nichtbezeichnetem markiert. Wie Francisco J. Varela beschreibt, tritt dieser reentrante Prozess auf phänomenologischer Ebene folgendermaßen in Erscheinung: Im menschlichen Erleben finden wir normalerweise nicht nur einen Fluss von Wahrnehmungen oder Gedanken vor, die unsystematisch oder aufgrund habitueller Dispositionen verknüpft werden. Vielmehr erscheint im Bewusstseinsprozess auch eine Art Diskontinuität, eine plötzliche Unterbrechung, mit der man sich bewusst wird, dass man etwas wahrnimmt (im Englischen: ›You become aware of‹). Man sagt sich vielleicht innerlich so etwas wie ›Aha‹.⁸³ Diese Fähigkeit der Reflexion würden wir normalerweise als ›Bewusstsein‹ bezeichnen. Es erscheint nun insofern ein »Beobachter«, als der kognitive Prozess nicht nur Unterschiede verarbeitet, sondern sich selbst als Beobachterposition erkennen und unterscheiden kann. Auf einmal verlagert sich die Aufmerksamkeit vom Inhalt auf den Prozess – was dann als subjektive Position erscheint, die zur Selbstbeobachtung fähig wird.

Formal und in Hinblick auf die gezeigte Dynamik entspricht dies den neuronalen Korrelaten des Bewusstseins, die Giulio Tononi

82 Siehe hierzu ausführlich Varela (1999).

83 Siehe Varela, MIT Artificial Intelligence Laboratory Lecture, 1997: <https://www.youtube.com/watch?v=xhS9ffDcLOk> [Abruf: 19.03.2023].

herausgearbeitet hat.⁸⁴ Tononi beschreibt das Bewusstsein auf formaler Ebene als einen neurokonstruktiven Vorgang, der aus nahezu beliebigen Anfangsaktivitäten ein sich reentrant stabilisierendes Aktivitätsmuster hervorbringt. Dies würde auch erklären, dass solche Systeme träumen oder halluzinieren können, also die von ihnen hervorgebrachten Formen nicht in einer fest gekoppelten Beziehung zu den Mustern und Gesetzmäßigkeiten der sogenannten Außenwelt (oder ›Realität‹) stehen müssen.

An dieser Stelle lohnt es sich, etwas genauer über die Außenseite – den unmarkierten Raum – nachzudenken. Wir könnten die Sache jetzt metaphysisch erhöhen und von der Leere (*void*) sprechen, womit wir Anklänge an die mystische Tradition finden.⁸⁵ Das Thema lässt sich jedoch auch profaner angehen. So kann mit Blick auf den Bewusstseinsprozess festgestellt werden, dass die Aktivität des Unterscheidens für weitere Operationen eben nur das zur Verfügung hat, was sie selbst als Bezeichnungen hervorgebracht hat. Der hiermit erscheinende Beobachter hat nur die Innenseite zur Verfügung und so kann es für ihn so scheinen, als ob er selbst es gewesen wäre, der die erste Unterscheidung getroffen und auf diese Weise seine Welt aufgebaut hat. Wir landen hier bei der idealistischen Verkürzung, die bereits bei Fichte anklingt, wenn er schreibt: »Dasjenige, dessen Sein bloß darin besteht, dass es sich selbst als seiend setzt, ist das Ich, als absolutes Subjekt.«⁸⁶

Sofern wir den *unmarked space* jedoch nicht als wirklich leer betrachten, sondern als die Totalität all dessen, was dem Prozess des Unterscheidens und Bezeichnens als Ergebnis dieser Operation nicht verfügbar ist, dann ergibt sich ein vollkommen anderes Bild: das der konditionierten Koproduktion. Einerseits bleibt wahr, dass Nervenaktivitäten nur Nervenaktivitäten reentrant verknüpfen können (so, wie im Froschkonzert das Quaken nur an andere quakende Lautäußerungen anschließt, um das Klangmuster hervorzubringen). Damit dies geschehen kann, muss aber andererseits eine strukturierte Welt gegeben sein: Diese enthält lebendige Körper, eine Umwelt, die diese nährt, und vieles andere. Oder um es anders auszudrücken: Der *unmarked space* muss als Medium all die Strukturmöglichkeiten und Aktivierungspotenziale mitbringen, damit überhaupt reentrante Formen entstehen können, die sich im Prozess des Wiedereintritts der Form in die Form als ›Beobachter‹ identifizieren können.⁸⁷ Aus der Innenperspektive einer potenziellen Form erscheint

84 Tononi (2004).

85 Zu den Verbindungen von Formtheorie und mystischer Tradition siehe etwa Engstrom (2001); zu Homologien mathematisch-theoretischer Hochabstraktion und Mystik siehe auch Fuchs (2008).

86 Fichte (1997 [1802], S. 23).

87 Im Sinne von Günther (2021 [1957], S. 110 ff.) erscheint das System aus der Außenperspektive vollspezifiziert, aus der Innenperspektive jedoch ›frei‹, da unterspezifiziert.

das Medium demgegenüber als ›leer‹, da sie die Bedingungen ihrer Möglichkeit nicht sehen kann.

Um dies auf die subjektive Erfahrung zu übertragen: Wenngleich das Bewusstsein eine in sich geschlossene Form darstellt und sich (wenn einmal in Gang gesetzt) nur rekurrierend auf den eigenen Unterscheidungsgebrauch reproduzieren kann, benötigt es doch all die Mittel, die ihm von der Außenseite zur Verfügung gestellt werden, um überhaupt aktiv werden zu können.

Von der Innenseite betrachtet erscheint der *unmarked space* als das Medium, in das sich die eigenen Formen hineinschreiben. Das Medium muss jedoch so beschaffen sein, dass sich überhaupt etwas in es hineinschreiben lässt, also das Ergebnis durch weitere Vorgänge stabil gehalten wird. Aus kybernetischer Perspektive muss das Medium also über Feedbackprozesse und veränderbare Strukturen verfügen, die (minimale) Abweichungen in einer Weise verstärken, dass diese kondensiert und damit erinnert werden können.

Mit Blick auf die Hirnforschung erscheinen die unterschiedlichen Gruppen neuronaler Ensembles als das Medium, in das sich die Aktivitätsmuster einschreiben können. Erst wenn unterschiedliche Nervengruppen bestehen und so miteinander vernetzt sind, dass reentrante Schleifen miteinander verkoppelt werden, können Aktivitätsmuster entstehen, die relativ stabil bleiben und gegen Veränderung resistent sind – und damit eine Form als Identität kondensieren lassen (auf phänomenologischer Ebene entspricht dies der Erfahrung, dass eine einmal festgesetzte Wahrnehmung oder Situationsdeutung nicht so leicht wieder ausgelöscht werden kann). Wenn die strukturierte Aktivität demgegenüber wieder abklingt und nur noch vereinzelt einige Neuronengruppen feuern, ohne ein übergreifendes kohärentes Muster auszubilden, ist das Medium wieder fluide, das heißt *offen*, zu weiteren Formen angeregt zu werden.

Nur in ›Nichtaktivität‹ kann leicht etwas eingeschrieben werden. Spencer Brown pointiert in formal-logischer Hochabstraktion, dass nur das Nicht-Etwas (*nothing*) sensibel genug ist, sich durch ein anderes Nicht-Etwas anregen zu lassen, um den Formprozess zu starten.⁸⁸

Die zuvor erwähnte Nichtaktivität darf jedoch nicht als ›absolute Stille‹ oder als totale Passivität verstanden werden, sondern eher als ein Basisrauschen, aus dem heraus spontane Fluktuationen entstehen, die dann Angebote für Formen bieten, welche aber nur stabilisiert werden

88 »I realized that the only ›thing‹ (i. e. nothing) that would be sensitive enough to be influenced by a stimulus so weak that it didn't exist, was nothing itself. That is, nothing is the only ›thing‹ that is so unstable that it can ›go off‹ off its own accord, the only ›thing‹ sensitive enough to be changed by nothing.« (Spencer Brown 2013, S. 6)

können, sofern sie mit anderen Fluktuationen rekursiv gekoppelt werden. Zudem muss vorausgesetzt werden, dass das Medium die Prozesse bereitstellt, damit dies möglich wird. Im Sinne von Heinz von Foerster begegnen wir hier dem Order-from-Noise-Prinzip.⁸⁹ Als Ordnung (*order*) erscheinen jetzt all die Formen, die eine bestimmte Gestalt hervorbringen, darunter auch die, die sich über die Operation des Re-entry selbst zugänglich sind. Als Rauschen (*noise*) erscheint der *unmarked space*. Dieser muss jedoch seinerseits über einen gewissen Strukturreichtum verfügen. Genauer gesagt: Das Medium muss sowohl offene und undefinierte Eigenschaften bereithalten, damit sich etwas Neues einschreiben kann, als auch Strukturen und Prozesse anbieten, mittels derer sich diese Eigenschaften definieren, stabilisieren, formatieren und modifizieren lassen.⁹⁰

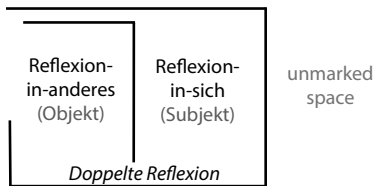


Abb. 10: *Reflexion und Re-entry*. Der Bewusstseinsprozess unterscheidet Gegenstände, was impliziert, dass da ein Subjekt ›beobachtet‹, wengleich dieses zunächst nur ›unmarkierter Raum‹, also gleichsam Nichts war. Sobald jedoch unterschieden und bezeichnet wurde, erscheint es als eine eigenständig benenn- und bezeichnare Position: als ›Subjekt‹ – gewissermaßen als innerer unmarkierter Raum (als »Introszendenz«⁹¹). In der ›doppelten Reflexion‹ kann die hiermit einhergehende Unterscheidung ihrerseits unterschieden und bezeichnet werden, wodurch das System unbestimmt wird, da so disponibel ist, wie Subjekt und Objekt unterschieden und bezeichnet werden. Dieser Prozess selbst hat jedoch wiederum eine Außenseite, den *unmarked space*, der durch die Unterscheidungen der Innenseite nicht erreicht werden kann, von dort aus gesehen mithin transzendent, also prinzipiell unerreichbar ist.

Welterkenntnis als die Möglichkeit, sich selbst von anderen Gegenständen zu unterscheiden (und darum zu wissen), kann jedoch nur auf der Innenseite der Form geschehen. Nur hier, also *innerhalb* der Form, erscheint die Bezeichnung von etwas. Nur hier wird im Re-entry die Oszillation zwischen dem Beobachter und dem bezeichneten Gegenstand

89 Siehe Foerster (2006).

90 Leben entsteht damit nicht allein aus einer Welt von Unordnung, wie Schrödinger in seiner berühmten Abhandlung *Was ist Leben* angenommen hatte (Schrödinger 1946). Vielmehr kann Leben entsprechend von Foersterns Order-from-Noise-Prinzip erst in einem Raum entstehen, der zugleich Ordnung und Unordnung enthält – und zwar eine Unordnung, die auf den unterschiedlichen Wertsystemen unterschiedlicher Systeme beruht, wie auch Günther (1976, S.261ff.) bemerkt.

91 Günther (2021 [1957], S. 20).

möglich. Entgegen den solipsistischen Anklängen an den Idealismus Fichtes ist eine Form allerdings nicht von ihrer Außenseite zu trennen.

Damit ergibt sich für unseren Zusammenhang eine klare Beziehung zwischen der Formtheorie von Spencer Brown und der polykontexturalen Logik Günthers. Der Haken als markierte Unterscheidung repräsentiert die uns vertraute *zweiwertige* Welt. Links steht der bezeichnete Gegenstand, rechts das reflektierende Subjekt als das, was der Gegenstand nicht ist. Der Re-entry-Haken steht für das, was Günther in Anklang an Hegel als doppelte Reflexion bezeichnet – nämlich die Fähigkeit, sich den Unterscheidungsgebrauch selbst reflexiv zugänglich zu machen und damit zwischen Subjekt und Objekt oszillieren zu können. Die Außenseite des unmarkierten Raums erscheint damit als die komplexe polykontexturale Welt, deren vielfältige Eigensinnigkeiten und subjektive Positionen von der Innenseite aus transzendent sind, also im ontologischen Sinne unzugänglich bleiben müssen.

So wenig eine lokale zweiwertige Welt von dem polykontexturalen Gefüge zu trennen ist, dem sie sich verdankt, ist auch die reentrante Figur des Beobachters (ob nun als phänomenologisches Bewusstsein oder als neuronales Korrelat des Bewusstseins verstanden) niemals von der Welt getrennt – wenngleich diese ihm als *unmarked space* nicht verfügbar ist. Der Beobachter verdankt sich der *konditionierten Koproduktion* von Welt und Selbst, von Umwelt und System – und weil Letzteres als emergentes Phänomen gleichsam aus dem Nichts erscheint, spricht Spencer Brown hier vom »Kanon Null«:

»Kanon Null. Koproduktion

Was ein Ding ist, und was es nicht ist, sind in der Form, identisch gleich.

Das heißt, die identische Form oder Definition oder Unterscheidung agiert als die Grenze oder Beschreibung sowohl des Dinges als auch dessen, was es nicht ist.

Daraus ist das Kolloriar, daß alles und nichts formal identisch sind, leicht zu beweisen. (*Beweis*: Beiden fehlt jegliche Form überhaupt.)«

Und weiter mit Blick auf die Kritik an der abendländischen Philosophie:

»An diesem Punkt [wird auch] die falsche Annahme« deutlich, die den »westlichen Philosophen hunderte Jahre im Nacken gesessen« hat, nämlich davon auszugehen, dass, »weil nichts keine Form hat«, dies »nicht Basis beobachteter Phänomene sein kann, da beobachtete Phänomene offensichtlich sehr wohl eine konditionierte Struktur haben.

Meine Lehre besteht zur Gänze darin, diesen Fehler richtig zu stellen, indem sie klar zeigt, dass nichts in der Tat eine konditionierte Struktur hat, nämlich, dass wenn eine Unterscheidung ›in‹ nichts getroffen werden könnte, dann das Ganze der konditionierten Koproduktion, deren

Operation unentrinnbar ist und vollständig sichtbar, unvermeidlich stattfinden würde.«⁹²

Spencer Brown drückt die komplexe Beziehung der konditionierten Ko-
produktion auf den ersten Seiten von *Law of Forms* durch folgenden
Satz aus:

»Distinction is perfect continence.«⁹³

Dirk Baecker kommentiert hierzu: »Eine Unterscheidung ist vollkomme-
ne Kontinuität im Sinne einer vollkommenen Be-Inhaltung. Eine Unter-
scheidung enthält die Welt und bringt diese Welt hervor, weil sie nur in
dieser Welt gemacht werden kann, wenn sie etwas in dieser Welt bezeich-
net und dadurch das bezeichnete Etwas von der Welt unterscheidet.«⁹⁴

Hiermit einhergehend erscheinen zwei unterschiedliche Weltbegriffe,
die im Sinne einer logischen Buchhaltung sauber auseinanderzuhalten
sind:

1. Der erste Weltbegriff bezieht sich ausschließlich auf die von der In-
nenseite der Form zugängliche Wirklichkeit. Es ist die uns vertrau-
te zweiwertige Welt. Diese beruht jedoch auf selektiver Blindheit, da
nicht gesehen werden kann, was auf der Außenseite passiert. Dies ist
die Welt, die sich im Sinne von Gotthard Günther positivsprachlich
beschreiben lässt,⁹⁵ da hier etwas für uns als Gestalt, Thema oder Ge-
genstand erscheinen kann.
2. Der zweite Weltbegriff hingegen impliziert eine Totalität, die den *un-
marked space* mitenthaltend muss und die aus kybernetischen Grün-
den von keiner denkbaren logischen Position aus zugänglich ist. Es
kann keinen Beobachter der Welt geben, der vollständiges Wissen
über ebendiese Welt hat. Diese Welt ist nur negativsprachlich thema-
tisierbar – etwa als Wirkung eines Mediums, das selbst nicht gesehen
werden kann. Dies geschieht, wenn ein lebendiges System spontan zu
einem anderen Selbst- und Weltverhältnis springt, ohne jedoch selbst
über die transjunktionalen Operationen zu verfügen, die dieses ver-
mittelt haben. Es ist beispielsweise erfahrbar als ›Gnade‹ – um den
wichtigen Begriff der großen spirituellen Traditionen aufzugreifen.

92 Spencer Brown (1997, S. ix f., kursiv im Original). An dieser Stelle folgen
wir nicht dem Impetus, die Logik der Form mit der Logik des Universums
gleichzusetzen, also aus der Formtheorie eine metaphysische Position zu ma-
chen. Wir bleiben bei der Anwendung der Formtheorie als einer raffinierten
Grundlegung der Logik, die sich in Richtung Mehrwertigkeit erweitern lässt
und mithilfe derer sich dann Selbstorganisationsprozesse verstehen lassen.

93 Spencer Brown (1997, S. 1).

94 Baecker (2015, S. 45; Übersetzung von uns).

95 Siehe zum Begriff der Negativsprache auch Günther (1980).

All dies erscheint als Folge einer konditionierten Koproduktion, deren Komplexität sich jedoch von der Innenseite nicht begreifen lässt.

Sowohl Günther als auch Spencer Brown kommen damit zu dem Schluss, dass die zweiwertige Welt einer sich selbst unterscheiden können Subjektivität unterspezifiziert bleiben muss. Diese Subjektivität kann prinzipiell nicht die volle Komplexität der Bedingungen der konditionierten Koproduktion ihrer eigenen Existenz begreifen. Selbst wenn ein Beobachtungsprozess über ein Reflexionsvermögen verfügt, das ihn um die Mehrwertigkeit und Polykontextualität seiner Existenz wissen lässt, kann dieser die damit verbundene Komplexität der logischen Beziehungen nicht einmal ansatzweise abbilden, geschweige denn nachvollziehen. Ein Beobachter bzw. die mit ihm einhergehende subjektive Position beruht auf Nichtwissen. »Existence is selective blindness«⁹⁶, formuliert Spencer Brown, um die spezifische Verschränkung von Epistemologie und Ontologie zu pointieren, die auch unsere menschliche Selbst- und Welterkenntnis auszeichnet.

Dieses Nichtwissen ist jedoch nicht unschuldig. Vielmehr bricht die subjektive Existenz die Einheit einer einwertigen Welt auf, in der immer schon alles gegeben und determiniert ist – und macht damit einen Unterschied! Die Welt erscheint jetzt als eine gebrochene, die sich selbst und anderen nicht mehr als eine Totalität verfügbar ist.

Es kommt nun, wie Günther ausdrückt, eine Reflexion ins Spiel, die sich der Positivität des Seins entzieht und damit aus dem Unbestimmten – aus der Negativität – heraus agiert.⁹⁷

Um nochmals auf die Prozesse zu verweisen, die wir zuvor als neuronale Korrelate des Bewusstseins beschrieben haben: Es wird nicht einfach Information im Sinne eines linearen Prozesses von Wahrnehmung und Reaktion verarbeitet. Durch die reentrant aufeinander projizierenden Aktivitätsmuster entsteht mit dem Bewusstsein vielmehr eine neue, emergente Qualität. *Das Gehirn beginnt autonom seine eigenen Zustände zu kreieren*. Es verarbeitet nicht nur einfach Information, sondern generiert etwas Neues. Es beginnt eine Welt zu halluzinieren, die mit dem Weltverhältnis, in dem der Organismus sich befindet (wie es sich

96 Spencer Brown (1997, S. 192).

97 In diesem Punkt folgen wir dem Argument von Roger Penrose (1989), dass Bewusstsein nicht in dem Sinne computerähnlich verstanden werden darf, dass für seine Erklärung simpel einem Algorithmus gefolgt wird. Auch wenn wir das bewusstseinsfähige Gehirn deterministisch verstehen, müsse es nicht-computable Prozesse beinhalten und damit prinzipiell unvorhersagbar agieren, folgert Penrose mit Blick auf Gödels Unvollständigkeitstheorem. Dies heißt jedoch nicht, dass wir auch den anderen, teils hochspekulativen Annahmen Penroses zur Genese des Bewusstseins folgen müssen.

beispielsweise für eine äußere wissenschaftliche Beobachterin darstellen würde), nicht mehr in einer unmittelbaren Beziehung zu stehen braucht.

Für Spencer Brown entstehen die hiermit einhergehenden Freiheitsgrade notwendigerweise aus der Form des Re-entry. Die Form hat auf der Innenseite das Bezeichnete und außen das Unmarkierte, kann jedoch über den Wiedereintritt der Form in die Form zwischen beiden Seiten oszillieren. Mathematisch gesehen gewinnt das System hierdurch einen imaginären Wert und damit einen Freiheitsgrad. Für einen äußeren Beobachter ist nicht feststellbar, wo die Form jeweils einrastet. Die Verhältnisse werden damit komplex, das heißt prinzipiell undurchschaubar. Dies ist homolog mit Günthers reflexionslogischer Beschreibung. Wir begegnen hier der »Reflexion-in-anderes« (dem als Objekt erscheinenden Gegenstand) und der »Reflexion-in-sich« (dem sich seiner selbst bewussten Subjekt), wobei die voll entwickelte Reflexion als »doppelte Reflexion« die Unterscheidung zwischen diesen beiden Stellen selbst disponibel halten kann.⁹⁸

Sowohl der Formkalkül als auch Günthers Reflexionslogik sind asymmetrisch gebaut. Das Positive – das Sein bzw. das Bezeichnete – steht der Negativität der Reflexion entgegen, die sich einer unmittelbaren Aneignung entzieht, da sie aus dem Unbestimmten – also logisch gesehen aus dem Nichts – gespeist wird. Letzteres muss aus kybernetischen Gründen notwendigerweise unterspezifiziert bleiben, da keine subjektive Perspektive vollständiges Wissen über die Welt haben kann. Dies entspricht der Einsicht von Spencer Brown, dass keine Form des Unterscheidens und Bezeichnens auf der Innenseite – und sei sie noch so elaboriert – den *unmarked space* auf der Außenseite der Form zum Verschwinden bringt. Jede weitere Beobachtung kann auch nur etwas bezeichnen und lässt damit anderes unbezeichnet.

Die Asymmetrie in Bezug auf das, was durch den Prozess der Beobachtung gewusst werden kann und was nicht zugänglich ist, drückt sich auch darin aus, dass der Beobachtungs- bzw. Reflexionsprozess einen Zeitpfeil generiert: Die Zukunft ist offen, da sie nicht gewusst werden kann. Die Vergangenheit ist jedoch bestimmt, nämlich als Geschichte der Unterscheidungen und Bezeichnungen, die zu der jeweils aktuellen Form oder Subjektivität geführt haben. Gerade weil die Konsequenzen des eigenen Unterscheidungsgebrauchs nicht überblickt werden können und man nicht weiß, was in der Zukunft geschieht, entstehen in reentranten kognitiven Systemen weitere Bewegungen, diesem Problem durch die Produktion neuer Unterscheidungen und Bezeichnungen zu begegnen. Dies erschafft eine neue Gegenwart, subjektiv gesehen schreitet die Zeit voran, da sich die Aktivität intentional auf eine noch nicht gewusste (da noch nicht erschienene) ›Wirklichkeit‹ bezieht. Das, was

98 Günther (2021 [1957], S. 48 ff.).

Phänomenologen als intentionales Bewusstsein zu fassen versuchen, erscheint hiermit der Form nach homolog mit der reentranten Aktivität eines Gehirns, das sich auf Basis der Anregungen der verrauschten Muster aus der Umwelt, die ihm die Sinnesorgane zuspielden, immerfort eine neue Welt erschafft.

Spencer Brown pointiert diesen Prozess primär aus der Perspektive eines einzelnen Beobachters, also in Hinblick auf das, was von der Innenseite der Form her gewusst werden kann. Das Unterscheiden und Bezeichnen lässt *eine* Welt entstehen, ohne damit freilich zu leugnen, dass all das, was im *unmarked space* geschieht, notwendigerweise an diesem Prozess beteiligt sein muss. Damit eröffnet der Formkalkül mit dem Re-entry bereits die Tür zu einer mehrwertigen, polykontexturalen Welt, die nicht nur unterschiedliche subjektive Positionen beinhaltet, sondern zudem davon ausgehen muss, dass diese entsprechend der konditionierten Koproduktion miteinander verschränkt sind.

An dieser Stelle ist darauf hinzuweisen, dass der Re-entry der Unterscheidung von Markierung und unmarkiertem Raum logisch gesehen zwei mögliche Anschlüsse mit sich bringt:

1. Zum einen ist die positive Affirmation eines Beobachters möglich, der sich selbst durch seine Unterscheidungen markiert. Dies führt zur Kondensation einer Identität, die sich über die Abgrenzung zu einem (ebenfalls selbst konstruierten) Gegenstand immerfort zu reproduzieren sucht. Die Identität (und das aus dieser Perspektive konstruierte Weltverhältnis) beginnt sich zu verhärten.
2. Der Re-entry kann jedoch zum anderen auch auf der Seite der Kontingenz und damit auf Seite der Unbestimmtheit einrasten. Hiermit dringt gleichsam etwas bislang noch nicht Verfügbares aus dem *unmarked space* in den Unterscheidungsgebrauch hinein. Es entsteht eine neue Form, die jedoch strukturell von etwas informiert wurde, was jenseits der Innenseite der eigenen Form liegt.⁹⁹

Die Aktivitätsdynamik einer systemischen Konstellation, die auf der Re-entry-Form beruht, zeichnet sich damit zugleich durch die beiden Merkmale *Offenheit* wie auch *Geschlossenheit* aus. Letzteres bezeichnet etwa die Tendenz, eine bestimmte Gestalt, die einmal gefunden ist, durch rekursive Operationen aufrechtzuerhalten und zu reproduzieren. Die Offenheit zeigt sich demgegenüber darin, dass der Prozess sensibel dafür ist, durch Irritationen und Perturbationen aus dem *unmarked*

99 In der Philosophiegeschichte entspricht dies Positionen, die das Selbst nichteologisch fundieren wollen. Man denke etwa an den »Zwischenmenschen« im dialogischen Prinzip von Martin Buber (2002), die Berührung durch das Fremde bei Emmanuel Lévinas (1984) oder die In-Verantwortung-Setzung durch den (gewaltsamen) Akt eines anderen Menschen bei Judith Butler (2014).

space aufgestört zu werden, sofern es gelingt, in Resonanz mit potenziellen, vom Prozess bislang noch nicht realisierten Aktivitätsmustern zu treten.

Abb. 11: *Boring-Frauen*. Wir können zwischen der Wahrnehmung zweier unterschiedlicher Gestalten oszillieren (hier zwischen einer jungen und einer alten Frau), jedoch weder beide Gestalten gleichzeitig erkennen noch in der Wahrnehmung einer undefinierten Nichtgestalt einrasten, wie etwa anhand der Betrachtung der *Boring-Frauen* deutlich wird.¹⁰⁰



Um diese Theoriefigur wiederum an der neurophänomenologischen Forschung zu verdeutlichen: In der Gestaltwahrnehmung entwickelt sich ein Aktivitätsmuster, das die beteiligten neuronalen Ensembles in einer Weise verklavt, dass die Muster sich immerfort reproduzieren und damit stabil bleiben (Geschlossenheit) – dies entspricht der phänomenologischen Wahrnehmung einer bestimmten Gestalt. Wir nehmen eine Form nicht verschwommen oder als Überlagerung unterschiedlicher Formen wahr, sondern üblicherweise als distinktes, eindeutiges Bild. Am Beispiel von Kippbildern lässt sich jedoch zeigen, dass die Gestalt bei hinreichender Stimulation in ein alternatives Regime kippen kann, der Prozess ist also sensibel für Anregung in Richtung einer alternativen Figuration (Offenheit). Dies steht jedoch wiederum nicht im Widerspruch zu dem Befund, dass zunächst die aktuell dominante Gestalt als Muster aufrechterhalten wird und Störungen des Bildes vom kognitiven System als ›Rauschen‹ abgeblendet und rausgerechnet werden (Geschlossenheit). Die Gestaltwahrnehmung erscheint damit als Paradebeispiel für einen reentranten Prozess, der sowohl Identitäten kondensiert als auch hinreichend sensibel für die Dekonstruktion und Rekonstruktion derselben ist.¹⁰¹ Hiermit wird auch das Grundprinzip der Gestaltwahrnehmung

100 Quelle: https://de.wikipedia.org/wiki/Meine_Frau_und_meine_Schwiegermutter#/media/Datei:My_Wife_and_My_Mother-in-Law.jpg

101 In einem ähnlichen Sinne formuliert dies auch Dirk Baecker (2014, S. 81 f.) in seiner Neurosoziologie: »Das Gehirn [...] arbeitet fortlaufend an einer Karte der Welt, welche Vorhersagen darüber zu treffen erlaubt, welche Handlungen mit Belohnung rechnen dürfen und welche nicht, und in die Korrekturen eingetragen werden können, die neuen Realitäten jeweils Rechnung tragen. [...] Diese Karte ist im Wesentlichen eine Karte von Bewertungen. Es ist keine regional oder thematisch aufgeteilte Menge von Karten, sondern eine einzige Karte, so dass man jede Wahrnehmung mit jeder

verständlich, nämlich dass immer nur etwas wahrgenommen werden kann (also eine Gestalt), jedoch nicht ein Nichts¹⁰² – und ebenso wenig die Fülle potenziell anderer, bislang noch nicht erscheinener Gestalten.¹⁰³

Entsprechend dem Formkalkül können auf der Innenseite nur das Bezeichnete und die durch die Operation des Re-entry als Identität kondensierte Gestalt erfahren werden. Nicht zugänglich ist jedoch all das, was der *unmarked space* als konditionierte Koproduktion darüber hinaus sonst noch gestatten würde. Homolog stellt Gotthard Günther fest, dass das Bewusstsein sich nicht als polykontextural erfahren kann: Wir erleben nicht die Aufsplitterung in vielfältige, dezentrierte Subjektpositionen, sondern eine distinkte zweiwertige Welt, in der *ein* Beobachter auf etwas Bestimmtes zu schauen scheint. Es ist das Verdienst von Günther, die Aufmerksamkeit auf die logische Struktur einer mehrwertigen, polykontexturalen Welt zu lenken, die dieses Bild – und damit auch die Vorstellung vom Beobachter als einem inneren Seelenwesen – transzendiert.¹⁰⁴

Von der Innenseite einer zweiwertigen Welt (also von einer bestimmten Kontextur aus gesehen) ist eine andere Kontextur (also eine andere zweiwertige Welt) logisch nicht erreichbar. Die Du-Subjektivität ist von der phänomenologischen Ich-Subjektivität aus nicht zugänglich. Um in der Bewusstseinsanalogie der neuronalen Korrelate des Bewusstseins zu sprechen: Die neuronalen Dynamiken der Hirnprozesse sind über reentrante Prozesse in sich selbst abgebildet, jedoch jeweils nur in einem Hirn. Welche Muster in einem anderen Gehirn entstehen, ist kausal nicht miteinander gekoppelt (wenngleich dies nicht ausschließt, dass mehrere Gehirne unter bestimmten Voraussetzungen ähnliche Aktivitätsmuster entwickeln – also gleichsam miteinander in Resonanz treten – können¹⁰⁵). Welche Beschreibung man auch wählen mag (die phänomenologische oder die neuronale), der Mensch gewinnt sein Selbst- und Weltverhältnis nolens volens auf Basis der Unverfügbarkeit von Welt.

Erinnerung und mit jeder Erwartung verknüpfen kann und umgekehrt jede Operation Konsequenzen für die Neubewertung der anderen hat.«

- 102 Michael Ende (1979) hat in seiner *Unendlichen Geschichte* versucht, sich diesem nicht wahrnehmbaren Nichts literarisch anzunähern.
- 103 Man denke hier nochmals an die Experimente zur binokularen Rivalität, wo dem Bewusstsein entweder nur die dem rechten oder nur die dem linken Auge präsentierte Gestalt erscheint – nicht jedoch beide zugleich (siehe Engel et al. 1999).
- 104 Spencer Brown führt die konditionierte Koproduktion demgegenüber nur als Postulat ein, ohne die mit ihr verbundenen logischen Prozesse weiter zu explizieren.
- 105 Der aus der Physik stammende Begriff der Resonanz vermittelt eine komplexe Beziehung: Nur ein System, das operational geschlossen ist – also Eigenzustände aufbaut –, kann offen dafür sein, sich in seinen Eigenschwingungen durch andere Systeme anregen zu lassen.

Die Regelmäßigkeiten der physischen Welt können zwar teilweise als Muster antizipiert werden (man denke etwa an die Tatsache, dass jeden Morgen die Sonne aufgeht und Steine fest und schwer sind). Spätestens in Hinblick auf die Eigenaktivität eines fremden Organismus wird Polykontexturalität jedoch unhintergebar. Was in der Begegnung mit einem anderen Organismus passieren wird, kann über die Ausbildung von Erwartungen simuliert werden, es kann aber nicht gewusst werden. Formtheoretisch gesprochen lässt sich zwar auf der Innenseite der Form zwischen Ich und Du unterscheiden und dann über den Re-entry in die Figur der reziproken Perspektivenübernahme einrasten (etwa im Sinne der Überlegung, dass der andere Ähnliches begehrt wie ich). Dies ändert jedoch nichts daran, dass das künftige Verhalten des Gegenübers außerhalb der Form – also im *unmarked space* – liegt und dem Beobachter damit letztlich nichts anderes übrig bleibt, als sich immer wieder durch Enttäuschung seiner Erwartungen überraschen zu lassen.

Wie auch immer, sei es auch noch so elaboriert, ein reflexionsfähiges System kann nicht um die volle Komplexität des polykontexturalen Gewebes wissen, die sein Selbst- und Weltverhältnis ausmacht. Unweigerlich wird es deshalb aufgrund seiner Unwissenheit seinerseits in der Umwelt einen Unterschied machen. Was sich auf der Innenseite des Systems als Unwissenheit darstellt, erscheint für andere kognitive Systeme, die mit diesem System interagieren, als Unberechenbarkeit. Wenn man nicht weiß, was man tun soll, bleibt nichts anderes übrig, als sich zu entscheiden, irgendetwas zu tun. Aus der Innenperspektive erscheint dies als subjektive Freiheit, von außen gesehen als die Unberechenbarkeit einer nichttrivialen Maschine.¹⁰⁶ Das Bewusstsein (egal ob phänomenologisch oder mit Blick auf seine neuronalen Korrelate betrachtet) erscheint damit als eine ›Oberflächenstruktur‹, die Selbst und Welt – einschließlich anderer Wesenheiten – auf Basis eines »naiv realistischen Selbstmissverständnis[s]«¹⁰⁷ symbolisch darstellt und hieraus Implikationen für das Handeln ableitet. Die ›Tiefenstruktur‹ muss demgegenüber all die kausalen Prozesse beinhalten, die überhaupt erst Kognitionen, Organismen, soziale Systeme etc. möglich machen. Sie muss dem Bewusstsein prinzipiell intransparent bleiben.

Eine polykontexturale Beschreibung der Welt impliziert dabei, dass es mehr als eine subjektive Beobachterposition gibt, die unterspezifiziert ist, also nur über ein unvollständiges Wissen über das Netzwerk ihrer Relationen verfügt. Diese Beobachterpositionen stehen in einer nicht trivialen Weise miteinander in Beziehung. Denn die Positionen ›Ich‹, ›Es‹ und ›Du‹ sind aus Perspektive der Theorie der Polykontexturalität ihrerseits als Ergebnis wechselseitig verschränkter polyzentrischer Reflexionsprozesse zu verstehen.

106 von Foerster (1994, S. 357 f.).

107 Metzinger (1998).

*Split-Brain: das polykontexturale Gehirn
(und das monokontexturale Bewusstsein)*

Es ist mittlerweile deutlich geworden, was die formtheoretische Figur des Re-entry zur Erhellung neuronaler Aktivitäten leistet, die mit Bewusstsein einhergehen. Es stellt sich nun die Frage, wie die neurophänomenologische Forschung mit Gotthard Günther in Richtung einer polykontexturalen Beschreibung erweitert werden kann.

Auf einer empirischen Ebene würde dies allerdings verlangen, menschliche Gehirne zu untersuchen, die in Kommunikation mit anderen Menschen stehen. Es würde voraussetzen, das Gehirn nicht mehr isoliert mit Blick auf die Auseinandersetzung mit einer einfachen, experimentell standardisierbaren Aufgabe zu untersuchen. Es würde heißen, das menschliche Gehirn als offenes System zu betrachten, das in kommunikative Zusammenhänge eingebettet ist. Aufgrund der Komplexität der hiermit einhergehenden Fragestellungen sind solche Untersuchungen nicht so leicht durchzuführen und werden dementsprechend kaum realisiert. Das menschliche Gehirn in der Ich-Du-Konstellation bleibt so ein neurowissenschaftliches Desiderat.¹⁰⁸

Allerdings wurden vor nicht langer Zeit noch neurochirurgische Eingriffe an Menschen durchgeführt, die ebenfalls das Problem aufwerfen, was geschieht, wenn zwei subjektive Zentren bestehen, die strukturell voneinander getrennt sind, jedoch kommunikativ aufeinander referieren müssen.

In den 1960er Jahren wurden Patienten, die unter einer schweren Epilepsie litten, nicht selten einer Hirnoperation unterzogen, bei der das *corpus callosum*, der ›Balken‹ zwischen den beiden Hirnhälften, durchtrennt wurde.¹⁰⁹ Aufgrund der Zerstörung der verbindenden Nervenfasern war nun die rechte Hirnhälfte nicht mehr in der Lage, unmittelbar mit der linken Hirnhälfte zu kommunizieren, und umgekehrt. Entsprechend konnten sich die epileptischen Krampfanfälle nicht mehr über das ganze Gehirn ausbreiten. Was bedeutet nun dieser Eingriff für das Alltagsleben der sogenannten Split-Brain-Patienten? Auf den ersten Blick scheinen sie normal zu agieren. Üblicherweise regelt die dominante linke Hirnhälfte die Aktivitäten. Was geschieht aber, wenn die rechte Hirnhälfte plötzlich ›beschließt‹, einen Spaziergang zu machen? Wie reagiert die dissoziierte andere Gehirnhälfte darauf und was denkt sie dabei?

108 Ein positives Beispiel für eine Ausnahme sind die Arbeiten der Arbeitsgruppe von Uri Hasson, die unter anderem gekoppelte und geteilte Aufmerksamkeiten im Rahmen des Erzählens von Geschichten untersucht (siehe etwa Baldassano, Hasson & Norman 2018).

109 Bei den folgenden Ausführungen zum Experiment handelt es sich um leicht veränderte Texte aus Vogd (2006, S. 141 ff.), die jedoch nicht als wörtliche Zitate ausgewiesen sind, um den Lesefluss nicht zu stören.

Durch unter Laborbedingungen erfolgende Untersuchungen etwa durch Michael Gazzaniga in der Arbeitsgruppe des Nobelpreisträgers Roger Sperry wurde es experimentell möglich, sich mit solchen Fragen zu befassen.¹¹⁰ Der methodische Kniff bestand dabei darin, mit der Versuchsperson so zu kommunizieren, dass die eine Hirnhälfte des Probanden nicht weiß, was mit der anderen geschieht.

Nur mit einer Hirnhälfte zu kommunizieren, lässt sich für einen Versuchsleiter technisch relativ einfach realisieren. Üblicherweise geschieht dies dadurch, dass man eine Information, zum Beispiel ein Bild oder einen Satz, durch eine optische Apparatur nur auf die rechte bzw. linke Seite der Netzhaut des Auges projiziert und hierdurch entsprechend nur eine Hirnhälfte »anspricht«. Eine Reihe von Split-Brain-Patienten wurde auf diese Weise untersucht. Bei einigen der Versuchsteilnehmer zeigten erstaunlicherweise beide Hirnhälften die Fähigkeit, Schrift und Sprache zu verstehen sowie sich verbal auszudrücken.

Ein gut dokumentierter Fall, bei dem eine Reihe solcher Experimente durchgeführt wurde, war jener des 15-jährigen Paul: »Der Testleiter begann zum Beispiel eine mündliche Frage mit ›Wer ...?‹ und die offenen Stellen wurden durch eine Projektion in einem der Gesichtsfelder ergänzt, zum Beispiel durch die geschriebenen Worte ›... bist du‹. Auf die beiden Seiten vorgelegte Frage folgte die gleiche Antwort: ›Paul. Auf die Frage: ›Was für ein Tag ist morgen?‹ kam beide Male die richtige Antwort: ›Sonntag‹. Auf die Frage an die linke Hirnhälfte: ›Was willst du werden, wenn du groß bist?‹ antwortete Paul: ›Rennfahrer‹. Und dies ist faszinierend, weil auf die gleiche Frage an die rechte Gehirnhälfte die Antwort kam: ›Designer‹.«¹¹¹ Ein ebenso erstaunliches Ergebnis zeigt sich, wenn man die rechte Hirnhälfte zu einer Aktivität veranlasst, von der die linke nichts weiß: „Eine einfache Aufgabe wie ›gehen‹ wird in die üblicherweise ›stumme‹ rechte Gehirnhälfte geblitzt, und der Patient reagiert darauf typischerweise, indem er seinen Stuhl zurückschiebt und sich aus dem Testbereich entfernt. Wenn man ihn dann fragt, wohin er gehe, so lautet die Antwort meist etwa: ›Ich gehe raus, um mir eine Cola zu holen‹.«¹¹² Oder man gibt der einen Hirnhälfte die Aufforderung: ›Schließe das Fenster.‹ Daraufhin führt die Versuchsperson normalerweise die Anweisung aus. Aus Perspektive der anderen Hirnhälfte anschließend gefragt, warum sie denn das Fenster geschlossen habe, kommt dann üblicherweise eine Antwort der Art: ›... weil mir kalt war.‹

Bemerkenswerterweise antwortet die linke Hirnhälfte nicht im Sinne dessen, was sie eigentlich entsprechend der durch die Versuchsanordnung erzeugten Realität hätte empfinden müssen. Sie sagt nicht: ›Ich weiß es nicht‹, oder: ›Ich habe mich irgendwie getrieben gefühlt‹,

110 Gazzaniga (1989).

111 Maturana und Varela (1987, S. 248).

112 Gazzaniga (1989, S. 89 f.)

sondern sie erfindet eine Erklärung, einen übergeordneten Sinnzusammenhang, der mit dem erlebten Handeln des eigenen Körpers in Einklang steht. Das Gehirn konstruiert post hoc eine Rationalität, die dem vollzogenen Verhalten eine intentionale Handlungsabsicht unterlegt. Paul gibt an, zweckrational gehandelt zu haben, wenngleich die Sache offensichtlich anders liegt. Denn nicht das Motiv, seinen Durst stillen zu wollen, veranlasste ihn dazu, aufzustehen, sondern die zuvor in die andere Hirnhälfte geblitzte Anweisung des Versuchsleiters.

Für die beiden Neurobiologen Humberto R. Maturana und Francisco J. Varela offenbaren diese und homologe Experimente »Grundlegendes über die Weise, auf die der anhaltende Fluß von Reflexionen, den wir Bewußtsein nennen und mit unserer Identität assoziieren, organisiert ist und seine Kohärenz bewahrt«¹¹³.

Betrachten wir die zuvor geschilderten Versuche und deren Ergebnisse nun aus dem Blickwinkel einer Theorie der Polykontextualität. Um das Experiment aufschließen zu können, müssen wir von mindestens *drei subjektiven Zentren* ausgehen: der *rechten Hirnhälfte* des Patienten, der *linken Hirnhälfte* des Patienten und der Position des *Versuchsleiters*. Eine subjektive Position zeichnet sich jeweils durch *selektive Blindheit* aus (sie kann nicht vollständig um die Bedingungen ihrer eigenen Prozesse wissen). Zudem stellt sie sich entsprechend den vorangehenden Ausführungen als ein reentrantes Reflexionsverhältnis dar. Sie reflektiert auf sich selbst und auf die als ihre Wahrnehmung unterschiedenen und bezeichneten Gegenstände. Betrachten wir nun, was aus der Perspektive der unterschiedlichen subjektiven Zentren gewusst werden kann – und was nicht:

- i. Für die linke Hirnhälfte der Versuchsperson sind positiv erfahrbar (das heißt unterscheid- und bezeichnenbar):
 - a) die Daten (zum Beispiel Worte und Anweisungen), die aufgrund der Anregungen im rechten Gesichtsfeld erscheinen,
 - b) die Repräsentation des Körpers (etwa seine Position und Bewegungen),
 - c) die sprachlich-propositionalen Ausdrücke, die im inneren Dialog generiert werden und die mit der Reflexion auf sich selbst einhergehen.
 Nicht zugänglich ist für die linke Hirnhälfte:
 - d) was die rechte Hirnhälfte als Gegenstände repräsentiert,
 - e) wie die rechte Hirnhälfte auf sich selbst reflektiert,
 - f) was der Versuchsleiter erlebt.

113 Maturana und Varela (1987, S. 249 f.).

2. Für die rechte Hirnhälfte lässt sich eine homologe Liste erstellen, nur dass für sie neben der Perspektive des Versuchsleiters auch die Bezeichnungen und Unterscheidungen der linken Hirnhälfte nicht zugänglich sind.
3. Für den Versuchsleiter lässt sich demgegenüber eine andere Liste erstellen:
 - a) Er weiß, welche Anweisungen und Informationen er der linken und der rechten Hirnhälfte der Versuchsperson gegeben hat, verfügt also diesbezüglich über mehr Wissen als jede der beiden Hirnhälften,
 - b) er weiß, wie sich der Körper der Versuchsperson verhalten hat,
 - c) er weiß, was die Versuchsperson gesagt hat.
 - d) Nicht zugänglich ist für den Versuchsleiter jedoch:
 - e) was die rechte oder die linke Hirnhälfte reflektiert, erlebt oder als interne Zustände konstruiert.

Versuchen wir das Beispiel weiter zu analysieren: Die rechte Hirnhälfte weiß nicht, was für die linke als Gegenstand erscheint, aber sie ist in eigentümlicher Weise mit der anderen Hirnhälfte verschränkt, erlebt sie doch die Folgen von deren Aktivität und der daraus erwachsenden Handlungsimpulse im Körper. Im Sinne der Dichotomie ›Leib sein‹ und ›Körper haben‹ ergeben sich hier unterschiedliche Möglichkeiten der Reflexion: Der Leib und seine Bewegungen können als die *eigenen* Handlungsimpulse erlebt werden. Der Körper könnte sich aber auch als fremdgesteuert *erfahren*. Letzteres würden beispielsweise manche Menschen erleben, die als schizophrene diagnostiziert wurden.¹¹⁴

Der Abgleich der von den beiden Hirnhälften hervorgerufenen Selbstberichte und der Beobachtungen des Versuchsleiters macht deutlich, dass im Selbst- und Fremderleben keine objektiven Tatsachen abgebildet werden, sondern der Reflexionsprozess jeweils seine eigene Wirklichkeit generiert. Deshalb müssen die Berichte auch nicht miteinander übereinstimmen. Entsprechend dem Paradigma des Neurokonstruktivismus ist es hinreichend, dass die generierten Gestalten in sich konsistent sind (das heißt als reentrantes Aktivitätsmuster hinreichend stabil gehalten werden können). Es genügt, eine Gestalt zu finden, die in Hinblick auf die ausgeflaggten Selbst- und Weltverhältnisse in sich stimmig ist.

Darüber hinaus wird klar, dass die unterschiedlichen Positionen und die mit ihnen einhergehenden Perspektiven nicht unabhängig voneinander zu sehen sind. Wenn der Versuchsleiter keine Aufforderung an die linke Hirnhälfte gegeben hätte, würde für die rechte Hirnhälfte kein Anlass bestehen, die Körperbewegung qua Reflexion sich selbst zuzurechnen. Möglicherweise war es damit auch erst die Nachfrage des

114 Siehe hierzu etwa die Fallberichte von Sandsten, Zahavi und Parnas (2021).

Versuchsleiters, warum die Person aufgestanden sei, die diese Form der Selbstreflexion veranlasst hat. Nicht jede Bewegung und nicht jedes Verhalten muss bewusst ablaufen. Doch insofern die Aufmerksamkeit durch eine Nachfrage auf das Verhalten gelenkt wird, beginnt im Bewusstsein unweigerlich ein Reflexionsprozess, der es dann beispielsweise als ›mein‹ oder ›nicht mein‹ Verhalten erleben und dies zudem als begründungswürdig erscheinen lässt.

So ließe sich an dieser Stelle vermuten, dass vielleicht sogar erst die Tatsache der Unverfügbarkeit der objektiven Gründe in Verbindung mit der Frage nach den Motiven den Bewusstseinsprozess in einer Weise in Gang setzt, dass sich ein ›Ich‹ als Autor zu identifizieren beginnt.

Die linke Hirnhälfte würde damit gerade deshalb zu einem sich in dieser Weise Auskunft über sich gebenden Subjekt:

1. weil für sie keine Verfügbarkeit über die Ursache von Bezeichnungen und Unterscheidungen besteht, die auf der Innenseite der anderen Hirnhälfte erscheinen;
2. weil in ihr durch den (von außen) kommenden Impuls, über sich selbst und das Verhalten Rechenschaft abzugeben, ein reentranter Prozess angeregt wird;
3. weil die Struktur der Sprache mittels der reflexiven Beziehungen von ›ich‹ und ›mich‹, ›Ursache‹ und ›Wirkung‹ bzw. ›Prädikat‹ und ›Objekt‹ immer schon eine bestimmte Klasse von Begründungen nahelegt;
4. weil andere Positionen aus dem polykontexturalen Gefüge das Faktum der Unverfügbarkeit fremder Selbst- und Weltverhältnisse durch ihr unvorhersehbares Verhalten immer wieder präsent machen und es damit zur Gewohnheit werden lassen, die hiermit ausgelösten Überraschungen interpretativ durch Zuschreibung von Motiven ›verstehen‹ zu wollen;
5. weil aus 1. bis 4. eine Matrix der Kommunikation erwächst, die mit dem Medium der Sprache auf eine präpositionale Struktur trifft, welche die benannten Reflexionsprozesse auf Dauer stellen – also erinnern und wiederholen – lässt.

Um es zusammenzufassen: Es bedarf der konditionierten Koproduktion des gemeinsamen In-der-Sprache-Seins und anderer subjektiver Positionen, deren Verhalten nicht vorhersehbar ist, damit eine subjektive Position entsteht, die angesichts der Unwissenheit eigene Konstruktionen entwickelt und sich hierüber als autonom erfährt. Dieser Befund lässt sich auch von der anderen Seite her lesen: Es bedarf subjektiver Positionen, die aufgrund von Unwissenheit ihr eigenes Selbst- und Weltverhältnis aufbauen, damit über Sprache und Kommunikation eine soziale Sphäre entsteht, die die subjektiven Positionen aneinanderkoppelt. Mit Niklas Luhmann treffen wir hier auf die konditionierte Koproduktion von Körpern, psychischen Systemen und auf Kommunikation basierenden

sozialen Systemen.¹¹⁵ Das ›Ich‹ kommt erst im ›Du‹ zu sich. Doch dies geht nur auf Basis von Körpern, die mit ihren sinnlichen und kognitiven Fähigkeiten jene Prozesse zur Verfügung stellen, in denen sich das Ich im Du zugleich spiegeln wie auch von ihm differenzieren kann.¹¹⁶

In unserem Beispiel müssen beide Hirnhälften in der Lage sein, die sprachlichen Äußerungen des Versuchsleiters so wahrzunehmen, zu verstehen und zu spiegeln, als ob es die eigenen wären. Es reicht nicht aus, die Aufforderung »gehen« nur zu sehen oder zu hören, sie muss zum Imperativ des eigenen Handelns werden, zu einem ›Ich gehe‹. Ebenso muss die Frage »warum« zur eigenen Frage werden, um sodann eine Kette von Motiven und Kausalitäten zu konstruieren, die schließlich als eigene Gründe präsentiert werden. Sie muss zur eigenen Frage werden, um eine Antwort zu evozieren. Die Worte des Versuchsleiters nur unbeteiligt wie in einem langweiligen Kinofilm wahrzunehmen und vorbeirauschen zu lassen, würde nicht hinreichen.

Durch Sprache und Kommunikation berührt und affiziert werden zu können, erscheint somit als Voraussetzung, um sich als ein autonomes Ich zu reflektieren, sei es, indem man sich mit dem Kommunikationsangebot identifiziert, oder sei es umgekehrt, indem man sich von ihm abwendet, um dann die Abstoßung als Anlass zur Identifikation zu nehmen. Die Versuchsperson hätte beispielsweise die Aufforderung des Versuchsleiters aktiv negieren und etwa sagen können: ›Nein, ich stehe nicht auf.‹ ›Auf solche dummen Fragen gebe ich keine Antwort.‹ Freilich bedarf es hierfür wiederum eines Reflexionsprozesses, der auf Basis des Re-entry der Unterscheidung von Selbst und Umwelt disponibel hält, in welche Seite die Identifikation einrastet.

Auf basaler Ebene erscheint die hiermit einhergehende Reflexionsfähigkeit bereits durch die Prozesse gebahnt, die unter dem Begriff Spiegelneuronen populär geworden sind.¹¹⁷ Das neuronale System simuliert Zustände einer anderen Person, um dann für den Bewusstseinsprozess die Weichenstellung aufzuwerfen, ob sie als eigen oder als fremd wahrgenommen werden.¹¹⁸

Das Split-Brain-Experiment zeigt, dass Hirnprozesse nicht einmal in einem objektiven Sinne voneinander ›wissen‹ müssen, um sich miteinander zu verschränken. Die Realitätskonstruktionen der beiden Hirnhälften sind in einer Weise aneinandergespleißt, dass die Wirkungen der Aktivitäten der einen Anlass zur Selbstbestimmung der anderen geben. Sie sind aber in dem Sinne voneinander unabhängig, dass jede Seite

115 Luhmann (1984, S. 148 ff.).

116 Siehe in diesem Sinne auch Baecker (2014, S. 220).

117 Siehe Rizzolatti et al. (2006).

118 Entwicklungsgeschichtlich spricht mit Jäger (2013) einiges dafür, dass sich die Sprachfähigkeit des Menschen aus genau diesen Prozessen entwickelt hat.

ihr eigenes Weltverhältnis und ihr damit verbundenes Selbstverhältnis aufbaut.

Mit Blick auf die vorangehenden Ausführungen zur Verkörperung der Sprache erscheint das Split-Brain-Experiment damit paradigmatisch für das polykontexturale Gewebe kommunikativ gekoppelter Gehirne – unabhängig davon, ob sie im selben oder in unterschiedlichen Körpern wohnen. Die reentranten Formen des Bewusstseins verdanken sich der konditionierten Koproduktion, also einer Prozessdynamik, die sich auf der Innenseite der Form nur selbst hat, jedoch durch Strukturen und Aktivitäten der für sie ontologisch nicht verfügbaren Außenseite konditioniert wird.

Diese Strukturen und Aktivitäten können nicht vollkommen regelmäßig oder im Voraus berechenbar sein – denn sonst würde sich der kognitive Prozess mit der Zeit durch Routinen daran anpassen und damit weitgehend automatisch verlaufen können. Das Bewusstsein würde damit sein Bezugsproblem verlieren, auf Unbestimmtheit durch spontane subjektive Aktivität antworten zu müssen, und könnte stattdessen durch unbewusste Prozesse ersetzt werden. Die Strukturen und Aktivitäten der Welt (und anderer Organismen) können aber auch nicht vollkommen beliebig und willkürlich sein, da sich die Kognitionen sonst im Rauschen der Welt verlieren würden. Sie müssen Ordnung und Unordnung, Struktur und Rauschen, Erwartbares und Überraschungen enthalten.

Letzteres ist insbesondere dann der Fall, wenn die Umwelt andere subjektive Identitäten birgt, deren Zentren des Erlebens und Handelns aus der Kontextur des eigenen Erlebens und Handelns unzugänglich sind.¹¹⁹ In solchen Konstellationen tritt das auf, was Günther als *transjunktionale Operation* bezeichnet. Es geschieht etwas, dessen Ursache nicht benannt oder lokalisiert werden kann, dessen Wirkung jedoch erfahren wird – und dem nun von der Innenseite des Systems ein Sinn oder Kausalzusammenhang unterstellt wird.

Dies ist genau die Situation des Split-Brain-Patienten, bei dem die eine Hirnhälfte nicht weiß, warum der Körper etwas gemacht hat, aber gefordert ist, hierauf eine Antwort zu finden. Eine etwas andere, jedoch verwandte Konstellation tritt auf, wenn ein (normaler) Mensch die Bedeutung von Worten wahrnimmt, die nicht seine eigenen sind, aber die damit verbundenen Gesten, Haltungen oder Emotionen unwillkürlich spiegelt.

119 Wir können uns mit Cosmelli et al. (2010) gedankenexperimentell ein Gehirn vorstellen, das vollkommen von seiner Umwelt isoliert ist, beispielsweise in einem abgeschlossenen Tank mit einer Nährlösung lebt. Es wäre jedoch zu erwarten, dass dieses Gehirn nicht mehr in der Lage ist, stabile reentrante Formen aufzubauen – und damit auch kein ›Bewusstsein‹ generiert.

Auch hier werden über die eigene Leiblichkeit Wirkungen erfahren – um ein Beispiel zu geben: die unangenehmen Effekte verletzender Worte.¹²⁰ Dies regt wiederum die Modifikation und Ausdifferenzierung des eigenen Selbst- und Weltverhältnisses an – etwa indem man sich infolge einer Beleidigung wertlos und depressiv fühlt oder umgekehrt wütend.

Dies muss ebenso als eine transjunktionale Operation erscheinen, die von der Innenseite her gesehen gleichsam aus dem Off ihre Wirkung zeitigt. Sie kann auf Basis der eigenen Form, das heißt der zuvor bestehenden Unterscheidung zwischen Subjekt und Objekt, Selbst und Umwelt, nicht abgebildet oder antizipiert werden. Sie wirkt jedoch insofern instruktiv, als die hiermit einhergehenden Informationen das eigene Selbst- und Weltverhältnis konditionieren. Man mag sich etwa, um bei unserem Beispiel zu bleiben, vornehmen, sich von einer anderen Person nicht beleidigen zu lassen – doch dann treffen einen die Worte wider Erwarten dennoch. Die Theorie der Polykontextualität beschreibt das Arrangement unterschiedlicher subjektiver Positionen, die sich wechselseitig auf Basis solcher transjunktionaler Operationen konditionieren.

In den vorangehenden Ausführungen finden sich drei unterschiedliche Weisen des Sprechens über Bewusstsein:

1. Es ist die Rede von Bewusstsein als introszendendem Selbstbewusstsein, wie es in Anschluss an Husserl phänomenologisch beschrieben werden kann.
1. Darüber hinaus ist die Rede von ›Bewusstsein‹ als einem Prozess neuronaler Aktivitäten, die in sich selbst eintreten – so, wie es die Forschung zu den neuronalen Korrelaten des Bewusstseins aufgezeigt hat.
1. Und schließlich erscheint Bewusstsein als eine sprachlich-deskriptive Rekursion, in der ein Sprecher ›Ich‹ sagen und sich Erleben und Handeln zurechnen kann, obgleich ein äußerer Beobachter andere Kausalzusammenhänge feststellen mag.

Diese drei Weisen, über Bewusstsein zu sprechen, haben – so unterschiedlich die damit zusammenhängenden Denkformen auch sind – in wesentlichen Aspekten etwas gemein: Sie haben die Form des Re-entry, sie sind abhängig von der konditionierten Koproduktion und sie stellen jeweils Antworten auf das Bezugsproblem der selektiven Blindheit dar. Um es nochmals zusammenzufassen:

1. *Form des Re-entry*: Das phänomenale Erleben trifft auf ein Erleben, das erlebt, dass es erlebt. Die neuronale Aktivität trifft auf ein Aktivitätsmuster, das dadurch aktiviert wird, dass es aktiviert wird. Die polyphone Struktur der Sprache bringt eine Subjektposition hervor, die sich durch den Sprachgebrauch selbst hervorbringt.

120 Dies erklärt auch, warum Worte einen verletzen können, auch wenn man das gar nicht will. Siehe dazu Herrmann (2007).

2. *Konditionierte Koproduktion.* Das *etwas* im Bewusstsein erscheint, ein *bestimmtes* Aktivitätsmuster reentrant stabilisiert wird oder *etwas* in einem sprachlichen Satz als ›Ich‹ thematisiert werden kann, setzt eine nahezu unendliche Anzahl von Prozessen und Strukturen voraus, die weder erlebt noch von neuronalen Mustern repräsentiert und ebenso wenig gleichzeitig oder in annähernd angemessener Präzision sprachlich-semantisch artikuliert werden können. Neben der Evolution des Universums und den hiermit entstehenden materiellen Bedingungen ist für die Entwicklung des menschlichen Bewusstseinsprozesses dabei insbesondere die konditionierte Koproduktion von Körper, Psyche und Sozialem von Relevanz.
3. *Selektive Blindheit.* Sich aus der Welt zu einer Existenz zu erheben, die über sich selbst verfügt und damit nicht nur Teil der Welt ist, sondern eine Welt hat, geht nur unter Absehung der Bedingungen der konditionierten Koproduktion. Die Selbstreferenz, die ›Ich‹ sagt und auf der anderen Seite ›die Welt‹ erkennt, hat nur die Unterscheidungen und Bezeichnungen, die sie selbst hervorbringt.¹²¹ Die mit Bewusstsein oder Kommunikation einhergehenden Selbst- und Weltbeschreibungen mögen hochgradig ausdifferenziert und elaboriert sein, doch dies ändert nichts daran, dass man darüber, worüber man nicht sprechen kann, schweigen muss.¹²² Der Grund hierfür ist nicht, dass bestimmte Beschreibungen und Bezeichnungen tabu wären oder man nicht weitere, differenzierte Unterscheidungen und Bezeichnungen finden könnte. Das Problem liegt vielmehr in der Sache selbst. Das, was die Inhalte des Bewusstseinsprozesses ermöglicht, kann nicht gleichzeitig sein Gegenstand sein – muss also im *unmarked space* liegen.

Das Gehirn des Bewusstseins und das Bewusstsein der Maschinen

»Building real things in the real world is how you'll make complex systems like minds and life. Making robots that have to survive in real bodies, day to day on their own, is the only way to find artificial intelligence, or real intelligence. If you don't want a mind to emerge, then unhinge it from the body.«

Kevin Kelly¹²³

Was lässt sich nun mit Blick auf die vorangehenden Ausführungen zur Frage des Bewusstseins der Maschinen sagen: Erinnern wir uns zunächst an eine der Maximen der kybernetischen Methodologie: Wir haben für

121 Dies gilt auch für die Gesellschaft als Ganzes, wie Luhmann am Beispiel der ökologischen Kommunikation aufzeigt. Auch sie hat keinen unmittelbaren Zugriff auf die ökologischen Prozesse, derer sie sich verdankt (Luhmann 1986).

122 Hier in Anlehnung an Wittgenstein (1963 [1922], Proposition 7).

123 Kelly (1992, S. 50).

die »psychologischen« Begriffe«, die wir in den »technischen Modellen« wiederholen möchten, »ein objektives Modell« zu finden. Sofern wir also ›Subjektivität‹ oder gar menschliches ›Selbstbewusstsein‹ in eine kybernetische Maschine implementieren wollen, haben wir die logischen Prinzipien zu formulieren, die diese kognitiven Eigenschaften auszeichnen.¹²⁴

Im Falle des menschlichen Selbstbewusstseins müsste das System zu »Reflexion-in-sich« und »Reflexion-in-anderes« fähig sein, darüber hinaus zu »doppelter Reflexion«. Es müsste also vor allem in der Lage sein, den eigenen Unterscheidungsgebrauch zu reflektieren.¹²⁵ Als Kandidat für eine objektive Beschreibung sind die Organisation und Dynamik der reentranten neuronalen Aktivitätsmuster zu benennen, die die Hirnforschung als neuronale Korrelate des menschlichen Bewusstseins identifiziert hat. Das Grundmodell hierfür ist die Form des Re-entry, wie sie von Spencer Brown entwickelt und von Varela für die Hirnforschung nutzbar gemacht wurde.¹²⁶ Der Re-entry erzeugt seine eigene Unbestimmtheit, indem durch den Wiedereintritt der Unterscheidungen und Bezeichnungen in die Form ein sich selbst modifizierendes Gedächtnis entsteht. Darüber hinaus ist dieser Prozess von Anregungen von ›außen‹, mithin von der konditionierten Koproduktion abhängig.

Beides zusammengekommen bedeutet, dass die Implementierung der Re-entry-Form des Bewusstseins in ein kybernetisches System voraussetzt, dass dieses sowohl nach innen als auch nach außen offen ist. Aus der Innenperspektive muss man dem System also erlauben, sein eigenes Gedächtnis nicht nur auf der Ebene der Inhalte, sondern auch hinsichtlich seiner Prozesse (also der Art und Weise des Unterscheidens und Bezeichnens) zu (re)programmieren. In Hinblick auf seine Umweltbeziehungen muss man ihm gestatten, mit anderen kognitiven Systemen zu interagieren, die in dem Sinne über Subjektivität verfügen, dass ihre Verhaltensweisen auch für sie selbst nicht vollkommen berechenbar erscheinen. Erst in solchen Arrangements kann das System so etwas wie Subjektivität entwickeln. Erst in dieser Konstellation wird es zu transjunktionalen Operationen fähig sein, die auf der Innenseite Unterscheidungen und hiermit einhergehende Gestalten evozieren, die nicht einer vorausberechenbaren Kausalbeziehung folgen.

Das System würde damit aus der Innenperspektive unterspezifiziert erscheinen, und wenn es Empfindungen hätte, würde es sich möglicherweise als frei empfinden und vielleicht sogar glauben, dass sein Denken und seine Handlungen dem eigenen Willen entspringen. Wie Arthur C. Clarke in seiner Saga *2001: A Space Odyssey* antizipiert, könnte ein solches

124 Günther (2021 [1957], S. 9).

125 Günther (2021 [1957], S. 56).

126 Varela (1979).

System neurotisch oder gar psychotisch werden, also in Situationen geraten, in denen sein Selbst- und sein Weltverhältnis krisenhaft wird oder gar auseinanderzufallen droht.

Wie auch immer, für diejenigen, die das kybernetische System gebaut und die ›Bewusstseinsfunktionen‹ implementiert haben, wird es eine strukturdeterminierte Maschine bleiben. Logisch gesehen werden sie in der gleichen Position sein wie eine Hirnforscherin, die bei einem Probanden im Kortex die reentranten Hirndynamiken untersucht und mit ihm spricht. Seine Subjektivität wird für sie weiterhin eine prinzipiell unzugängliche Transzendenz bleiben. Sie mag zwar sein Verhalten sehen, seine Mimik wahrnehmen und seine Worte vernehmen, doch sie wird niemals wissen, was er empfindet oder ob er überhaupt etwas empfindet.

Doch auf der anderen Seite steht der Befund, dass die ›Konstruktion objektiver Modelle von Bewusstseinsfunktionen‹ ein ›echtes subjektives, sich selbst transparentes‹ Bewusstsein voraussetzt.¹²⁷ Nur deshalb kann die Hirnforschung überhaupt von neuronalen Korrelaten des Bewusstseins sprechen. Die Berichte vom Erleben des introszendenten Bewusstseins werden mit objektivierbaren neurophysiologischen Prozessen in Beziehung gesetzt – was voraussetzt, dass es so etwas wie Bewusstsein gibt.

Damit stehen wir vor einem logischen Dilemma, das aus dem Blickwinkel unseres eigenen, zweiwertig verfassten Bewusstseins nicht lösbar ist: Wie lässt sich einer Maschine (sowie den Gehirnen meiner Mitmenschen – und damit auch mir selbst) ein Bewusstsein zurechnen, das es objektiv doch gar nicht geben kann? Die alte von William James aufgeworfene Frage ›Does consciousness exist?‹¹²⁸ lässt sich auch heute – und unserer Auffassung nach auch in Zukunft – nicht in seriöser Weise positivsprachlich beantworten.

Jeder, der sich an die Sache herantastet, muss für sich klären, wie er es mit der anderen Seite hält, dem *unmarked space* – also dem, worüber man nicht sprechen kann.

127 Günther (2021 [1957], S. 21).

128 James (1904).

IV. Mensch, Bewusstsein und Maschine

Mit Blick auf die vorangehenden Kapitel verwickelt uns die Frage nach dem Bewusstsein in ein Dilemma. So wissen wir auf der operativen Ebene – also der Ebene des Prozesses, der Bewusstsein möglich macht – sehr viel über das Bewusstsein: Es erscheint als ein reentranter Prozess, der sich selbst ein Gedächtnis generiert und als Flussstruktur eine Identität aufbaut, die sich immerfort verändert. Bewusstsein offenbart sich dabei als ein autopoietischer Prozess, der sich immerfort hervorbringt, selbst bestätigt, dabei aber auch auf Strukturangebote aus seiner unmittelbaren Umwelt angewiesen ist. Es erscheint damit unweigerlich verkörpert. Mit Blick auf die Bedingungen, die den Prozess – etwa Gedächtnis und Selbstreferenz – möglich machen, scheint es unmöglich, sich Bewusstsein als unabhängig von einem konkreten Substrat (etwa als ein vom Körper unabhängiges Seelenwesen) vorzustellen.

Auf der ontologischen Ebene muss Bewusstsein hingegen ein unergründliches Geheimnis bleiben: Subjektives Erleben entzieht sich einer objektiven Beschreibung. Wenngleich als ureigene Innerlichkeit uns so nah, lässt es sich nicht fassen und greifen. Es gibt den inneren Beobachter und es gibt ihn nicht. Da wir jedoch Bewusstsein auf operativer Ebene mittlerweile zumindest in Ansätzen verstehen können, scheint kein grundsätzliches Hindernis zu bestehen, Maschinen zu bauen, die ein subjektives Zentrum entwickeln können und damit auch bewusstseinsfähig sein könnten. Da wir aber andererseits nicht einmal Zugang zum Bewusstsein unserer Mitmenschen haben (nur ihre Körper und Artikulationen sind uns gegeben, nicht jedoch ihr subjektives Erleben), muss auch für eine künstliche Intelligenz – wie elaboriert sie auch immer sein mag – gelten: Wir wissen nicht, ob sie ein inneres Erleben hat, also fühlt, sieht, hört oder in anderer Form etwas erlebt, bzw. wenn ja, was sie erlebt.

Die Frage nach dem Bewusstsein der Maschinen führt damit einerseits zu einem Geheimnis, andererseits offenbart sie eine gewisse Struktur. ›Bewusstseinsfähige‹ Maschinen können uns in besonderer Weise auf uns selbst zurückwerfen. Sie halten uns den Spiegel vor, indem sie uns zeigen, was es bedeutet, ein subjektives Zentrum zu haben, das sich dem Verstehen entzieht. Die daraus folgenden metaphysischen Konsequenzen werden in Kapitel IV.1 erörtert.

Wir begegnen dabei Theoriefiguren, die das Beobachterparadoxon beinhalten, also das Paradoxon einer Welt, die sich in sich über eine Art Schnittstelle selbst beobachtet und damit zugleich als Dualität und als Einheit erscheint. Aus diesem Grund lohnt sich ein Exkurs in die Quantentheorie, da auch in dieser die Frage auftaucht, ob und wann Materie

reflexiv wird und wie sich die hiermit einhergehenden Paradoxien fassen lassen. Dies ist Gegenstand von Kapitel IV.2.

Im letzten Abschnitt, Kapitel IV.3, werden wir die ethischen Fragen und Herausforderungen erörtern, die sich in einer Welt ergeben, in der künstliche Intelligenzen existieren, die als bewusstseinsfähig gelten.

IV.1 POSTHUMANE SPIRITUALITÄT – WER SIND WIR, WENN UNS DIE MASCHINEN GLEICH WERDEN?

Gotthard Günther ging bereits in den 1950er Jahren davon aus, dass Menschen in nicht allzu ferner Zukunft intelligente Roboter bauen würden, die wie wir ein subjektives Zentrum des Erlebens haben, also auch in dieser Hinsicht ebenbürtig sein werden. In seinem visionären Buch *Das Bewußtsein der Maschinen* untersucht Günther den jüdisch-christlichen Schöpfungsmythos und entdeckt dabei eine Analogie: Gott habe den Menschen nach seinem Ebenbild geformt und entsprechend beginne nun der Mensch, kybernetische Maschinen nach seinem Ebenbild zu erschaffen – und dies werde, so Günther, die spirituelle Verortung des Menschen fundamental verändern.¹

Um Missverständnisse zu vermeiden, haben wir jedoch zuvor unseren Lesern und Leserinnen Rechenschaft darüber abzulegen, was wir im Folgenden unter dem Begriff Spiritualität verstehen: Es geht, um es vorab zusammenzufassen, nicht um Glaubensfragen, sondern *um Relationen zur Welt und zu sich selbst*, also um Formen der Reflexion.

Dies bedeutet, dass Begriffe wie Spiritualität und Religion sowie die in diesem Zusammenhang behandelten Inhalte nicht ontologisch verstanden werden dürfen. Gerade mit Günther ergibt die positivsprachliche Rede von Gott als Entität ebenso wenig Sinn wie die Annahme, dass die Erleuchtung, von der die buddhistischen Lehren sprechen, einen sinnlich erfahrbaren Zustand bezeichnet.

Wer von Gott, Nirvana oder auch vom ›Sinn des Lebens‹ spricht, drückt mit den von ihm verwendeten Begriffen und Erzählungen vielmehr ein bestimmtes *Selbst- und Weltverhältnis* aus.² Spirituelle und religiöse Begriffe thematisieren demnach die Beziehung zwischen Subjekt und Welt, zwischen Individuum und Ganzheit, beschreiben also ein *Reflexionsverhältnis*.³ Im Sinne des Philosophen Ernst Tugendhat können wir das Kernanliegen von Spiritualität deshalb darin sehen, das eigene

1 Siehe etwa Günther (2021 [1957], S. 54 und 82 f.).

2 Siehe ausführlich Vogd (2018, Kap 3).

3 Spirituelle und religiöse Rede beruht damit auf dem, was Günther (1980) als »Negativsprache« bezeichnet. Siehe zur negativen Theologie etwa Rentsch (2010) und zur Negativsprache in den buddhistischen Lehren Vogd (2017).

Ego zu dezentrieren.⁴ Es geht darum, sich selbst in ein anderes, die eigene Egozentrik übergreifendes Verhältnis zu setzen – nicht mehr und nicht weniger.⁵

Gotthard Günther neigt mit Blick auf Hegel dazu, das Christentum als Archetypus einer Religiosität zu betrachten, die in der Lage ist, ein Subjekt hervorzubringen, das sich als autonom empfindet und von dieser Empfindung getragen der Welt gegenübertritt. Günther war ein Kind seiner Zeit und ist damit in einer ethnozentristischen Kultur aufgewachsen, die die europäische Geistesgeschichte als den Höhepunkt der Evolution begreift und dadurch unterschätzt, welche komplexen geistigen und kulturellen Figurationen an anderen Orten der Welt entstanden sind. Heutzutage, gut hundert Jahre später, beginnen wir zu ahnen, dass selbst vermeintlich ›primitive‹ Kulturen weitaus raffiniertere Reflexionsverhältnisse aufgebaut haben, als die Anthropologen des 19. Jahrhunderts gedacht haben.⁶ Zudem haben wir mittlerweile – etwa über den Umweg über China oder Indien⁷ – gelernt, unsere eigenen, verborgenen erkenntnistheoretischen Annahmen und Vorurteile ein wenig besser zu verstehen.

Die biblischen Geschichten sind – ebenso wie die mythologischen Quellen anderer Kulturen – polyphone Texte, die vielfältige Bedeutungsebenen umfassen und je nach kulturellem Kontext neu gelesen werden können. Es spricht deshalb auch heutzutage nichts dagegen, die biblische Schöpfungsgeschichte als eine Analogie dafür zu nehmen, wie sich der menschliche Geist entfalten kann und welche Aporien dabei entstehen. Wir können sie mit Günther als einen Anlass nehmen, die metaphysischen Fragestellungen zu behandeln, welche eine Kybernetik aufwirft, die den Geist-Materie-Dualismus zurückweist. Auf diese Weise begegnen wir den Fragen, die sich einer Zivilisation stellen, die aufgrund ihrer technologischen Entwicklung über kurz oder lang wohl kaum umhinkommt, künstliche Intelligenzen und ihre maschinellen Körper ähnlich zu behandeln wie den Menschen – als zugleich körperliches und geistiges Wesen.

4 Tugendhat (2006).

5 Dies kann in vielfältigen Formen geschehen, etwa durch eine Religion, die einen personifizierten Gott an die Spitze stellt, durch spirituelle Systeme wie den Buddhismus oder Daoismus, die ein unpersönliches Prinzip als Ursache des Leidens und der Erlösung aus dem Leid formulieren, oder auch in Form von Ökologien, die Teile der Natur als beseelt betrachten.

6 Siehe hierzu bereits Wittgensteins (1989) Kritik an den kruden Annahmen des Anthropologen Frazer über Magie und Religion in indigenen Gemeinschaften.

7 Siehe etwa Jullien (2002) und Waldenfels (2013).

*Der kybernetische Gott – Regelmäßigkeiten und
Notwendigkeiten der Kommunikation und Logik*

Günther sieht bereits im Schöpfungsbericht – der Genesis – die logische Struktur der spirituellen Problematik eines Wesens abgebildet, das nicht in eins mit sich selbst sein kann. Gott bringt in den sieben Tagen seiner Weltschöpfung Himmel und Erde hervor. Er scheidet Wasser vom Land und erschafft Pflanzen und Tiere. Jeden dieser Schöpfungsschritte beschließt Gott mit der Formel »*Und Gott sah, dass es gut war*«⁸.

Mit einer biologisch und geologisch geschulten Fantasie können wir uns das Ergebnis als ein austariertes ökologisches Gefüge von Lebewesen vorstellen, die sich in dynamischen kybernetischen Kreisläufen aufeinander beziehen. Die Verbindung von Wasser, Sonne und Erde ermöglicht Pflanzen. Dies erlaubt es Tieren, sich von diesen zu ernähren, was fleischfressenden Räubern eine Existenzgrundlage bietet. Letztere sorgen dafür, dass die Pflanzenfresser nicht überhandnehmen und sich die botanische Sphäre somit hinreichend entfalten kann, ohne ihrerseits etwa die Gewässer durch Fäulnisprodukte zu hypertrophieren.

Aus kybernetischer Sicht ist etwas gut, wenn eine hinreichende Anzahl von Feedbackschleifen mit positiven und negativen Vorzeichen⁹ besteht, die das Ganze im Gleichgewicht halten. Auf diese Weise wächst kein Element, keine Lebensform so stark an, dass die Reproduktion der beteiligten Arten prinzipiell gefährdet ist.

Damit brauchen die einzelnen Lebewesen nichts anderes zu machen als das, was sie immer tun. Ihre kognitiven Prozesse fallen mit ihrem Handeln zusammen. Der Löwe muss nicht darüber nachdenken, ob er töten darf. Die Gazellen und Giraffen müssen nicht darüber reflektieren, ob ihr Fressverhalten den Baumbestand schädigt. Alles ist so, wie es ist. Fressen und gefressen werden gehören zusammen, bilden ein übergreifendes System – und was daraus entsteht, muss per se gut sein, so die Tautologie eines mit sich selbst in Einklang befindlichen lebendigen Prozesses.

Die Relationen zwischen den einzelnen Teilen des ökologischen Gewebes weisen eine inhärente Gesetzmäßigkeit auf, die sich aus den jeweiligen Rückkoppelungsprozessen ergibt – ablesen lässt sich das beispielsweise an den beiden gekoppelten Verläufen einer Population von Räubern und Beutetieren.¹⁰

8 1. Mose 1 (Bibel, Einheitsübersetzung 2016).

9 Eine Feedbackschleife mit positiven Vorzeichen führt zu einem autokatalytischen Prozess und damit zur Vermehrung, eine mit negativen Vorzeichen führt hingegen zur Runterregulierung, also zur Eindämmung des Wachstums.

10 Mit der kulturellen Evolution des Menschen taucht das Problem eines Spitzenprädatoren auf, der in der Lage ist, alle ökologischen Systeme in einer Weise zu dominieren, dass kein Gleichgewicht mehr gefunden werden kann, das auf der Dezimierung der menschlichen Population beruht. Dies bedroht die

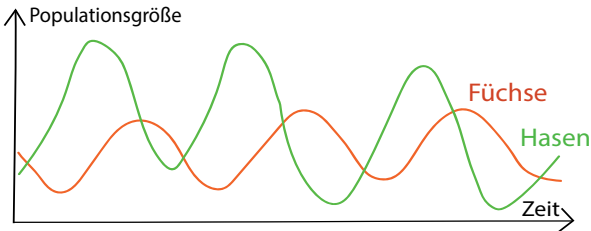


Abb. 12: Idealtypische Räuber-Beute-Kurve

Die »Regelmäßigkeiten in der Biosphäre« bilden in ihren Verflechtungen also bestimmte Muster, so der Kybernetiker Gregory Bateson. Da diese so »durchgängig und bestimmend« sind, können Kybernetiker in ihnen gleichsam eine Art heiliges oder göttliches Prinzip sehen: »Die Regelmäßigkeiten, die wir entdecken – einschließlich der Regelmäßigkeiten und Notwendigkeiten der Kommunikation und Logik –, bilden eine Einheit, in der wir uns häuslich einrichten. Sie ließen sich vielleicht als die Eigenheiten Gottes betrachten, den wir Öko nennen könnten.«¹¹

Im Sinne von Gotthard Günther wäre das hiermit einhergehende göttliche Bewusstsein – wenn es ein solches gäbe – als *einwertig* anzusehen (siehe Kap. I.2). ›Öko‹ denkt nicht, hat auch nicht den Charakter einer Person, die etwas Bestimmtes will. Ohne reflexive Differenz geschieht einfach, was gemäß den göttlichen Gesetzmäßigkeiten (die hier synonym mit den kybernetischen Regelkreisläufen zu verstehen sind) zu geschehen hat. Das Leben entfaltet sich in den Geweben der vielfältigen Formen seiner Selbstorganisation.

Es liegt dabei in der Natur der Sache, dass die jeweils gefundene Form – und die damit zusammenhängende Anpassungsleistung – ab einem gewissen Punkt der Evolution zum Problem wird. Wenn beispielsweise eine Population von Jägern zu erfolgreich wird und als Konsequenz nahezu den ganzen Bestand ihrer Beutetiere fängt, wird sie aufgrund dieser Zerstörung ihrer eigenen Nahrungsgrundlage hungern oder gar aussterben müssen. Die Feedbackschleifen der Natur werden dann erbarmungslos ein neues Gleichgewicht erzwingen. Dies geschieht nicht intentional bzw. nicht unter Beteiligung eines reflexiven Bewusstseins. Es passiert allein aufgrund der Struktureigenschaften der übergreifenden Dynamiken. Oder um es mit Bateson zu formulieren: »Öko« hat »keinen freien Willen«, ebenso wenig geht es ihm auf eine »simple Art um gut und böse«.

globale Ökologie in einer neuen, dramatischen Weise und weckt die Sehnsucht nach einer größeren, mächtigeren Intelligenz, die dem etwas entgegensetzen könnte. So auch bei Theresa Hannig (2022) in ihrem Roman *Pantopia*.

11 Bateson und Bateson (1993, S. 203).

›Öko‹ symbolisiert einfach nur »die Tatsache«, dass »Sucht und Pathologie die andere Seite der biologischen Anpassung sind«¹².

Die neun Gesetze Gottes

»OUT OF NOTHING, nature makes something. First there is hard rock planet; then there is life, lots of it. First barren hills; then brooks with fish and cattails and red-winged blackbirds. First an acorn; then an oak tree forest. I'd like to be able to do that. First a hunk of metal; then a robot. First some wires; then a mind. First some old genes; then a dinosaur. How do you make something from nothing? Although nature knows this trick, we haven't learned much just by watching her. We have learned more by our failures in creating complexity and by combining these lessons with small successes in imitating and understanding natural systems. So from the frontiers of computer science, and the edges of biological research, and the odd corners of interdisciplinary experimentation, I have compiled The Nine Laws of God governing the incubation of something from nothing:

1. Distribute being
2. Control from the bottom up
3. Cultivate increasing returns
4. Grow by chunking
5. Maximize the fringes
6. Honor your errors
7. Pursue no optima; have multiple goals
8. Seek persistent disequilibrium
9. Change changes itself.«

Kevin Kelly (1992, S. 468)

Weder die Vögel noch die Würmer oder die Katzen müssen sich über die Zukunft oder über moralische Fragen Gedanken machen. Ihre ›Aufgabe‹ besteht allein darin (und sie können auch nicht anders), (als) durch die eigene Existenz eine der Schnittstellen in den Regelkreisläufen der Natur zu besetzen. Die Vögel brauchen die übergreifenden Dynamiken oder Gesetzmäßigkeiten nicht zu begreifen, damit das Ganze (sprich: das Leben) funktioniert. Sie brauchen auch nicht um die Vorgänge in den

¹² Bateson und Bateson (1993, S. 212). Es findet sich somit kein ›Programm‹, keine ›Intention‹ hinter den ökologischen Netzwerken. Interessanterweise kann man das im Großen wie im Kleinen beobachten – auch die kognitiven Netzwerke des Bewusstseins sind prinzipiell nicht an gut und böse orientiert, sondern operieren einfach. Erst nachträglich, also unter dem Blickwinkel von Unterscheidungen, die sozial angeliefert (und zugemutet) werden, sind sie diesbezüglich bezeichnerbar.

Zellen zu wissen, die ihre Körper aufbauen. In diesem Sinne kann ihre Lebensform als unbekümmert gelten. Mit Blick auf die Eigenläufigkeit der Natur erscheinen sie auch in der Bibel als Beispiel für den berührenden Satz in der Bergpredigt:

»Sorget Euch nicht um morgen, denn der morgige Tag wird für das Seine sorgen.«¹³

Der Satz ist an die Menschen gerichtet, nicht an die Vögel. Erstere befinden sich schon längst nicht mehr in einem einwertigen Verhältnis zu einer Welt, die sich selbstgenügsam gegeben ist und gerade deshalb funktioniert, weil keiner ihrer Protagonisten über sein Handeln reflektieren muss. Das Vertrauen in eine Welt, die entsprechend den ›göttlichen Regeln‹ von ganz allein zu ihren Formen und den daraus resultierenden Gleichgewichten findet, ist schon längst gebrochen. In der Bergpredigt treffen wir auf menschliche Wesen, die die ›Perspektive der Sorge‹ in die Welt eingeführt haben. Die Menschen befinden sich bereits in einer Position, mit der Sorge die aktuell gegebenen Verhältnisse (sozusagen die ›göttliche Ordnung‹) zu negieren, indem diese in Hinblick auf eine imaginierte Zukunft problematisiert werden.

Dazu treffen wir in der Geschichte jetzt auf Jesus (es könnte auch Buddha oder ein anderer Repräsentant der heiligen Ordnung sein), der zu den Menschen spricht. Durch seine Rede negiert er die Perspektive der Menschen: Er fordert sie auf, sich nicht zu sorgen, sondern zur ursprünglichen Unschuld der einwertigen Welt – der ökologischen Ordnung sich selbst regulierender Prozesse – zurückzukehren. Diese Negation darf nicht als eine einfache Negation verstanden werden, etwa in dem Sinne, dass ein Mensch sagt: »Es wird einen Sturm geben, also gibt es einen Grund zur Sorge«, und Jesus antwortet: »Es wird *keinen* Sturm geben, also gibt es *keinen* Grund zur Sorge.« Sie ist im Sinne von Günther vielmehr als eine transklassische Negation zu verstehen, insofern hier jegliche Haltung, die mit Sorge einhergehen könnte, per se zurückgewiesen wird. Stattdessen wird auf den reflexionsfreien Zustand einer ›paradiesischen‹ einwertigen Welt verwiesen. In unserem Beispiel hieße dies: Die Frage, ob man sich sorgen solle oder nicht, ist sinnlos, da die Welt (schließlich ist sie identisch mit dem Wille Gottes) ohnehin das tut, was sie tut – und das kann auch mal das Heraufziehen eines Unheils, etwa eines Sturms, bedeuten, und auch dies ist ›gut‹ so.

Hiermit wird aber auch klar, dass es mehrere Lesarten des eingangs zitierten biblischen Kommentars zur Schöpfung gibt (»Und Gott sah, dass

13 Matthäus 6,34 (Bibel, Einheitsübersetzung 2016). Und hier auch: »Seht die Vögel unter dem Himmel an: Sie säen nicht, sie ernten nicht, sie sammeln nicht in die Scheunen; und euer himmlischer Vater ernährt sie doch. Seid ihr denn nicht viel kostbarer als sie?« (Matthäus 6,26)

es gut war.«). Im konventionellen Sinne verbinden wir den Begriff ›gut‹ mit einer bewertenden Reflexion. Wir betrachten eine Sache, vergleichen sie mit unseren moralischen oder sonstigen normativen Vorstellungen und charakterisieren sie anschließend als ›gut‹ oder ›schlecht‹. Im Sinne von Batesons kybernetischem Gott ›Öko‹ können wir ›gut‹ jedoch auch so verstehen, dass es kein Gegenteil hat. Das, was ›Öko‹ tut und sieht, ist identisch mit sich selbst, damit per se in sich logisch und rational – und demzufolge auch ›gut‹. Die sich dabei vollziehenden kognitiven Prozesse – etwa das ›Sehen‹ – dürfen nicht als Reflexionsprozesse verstanden werden (denn hiermit wären wir bereits bei einer sich selbst problematisieren könnenden zweiwertigen Welt angelangt). Es wäre vielmehr ein differenzloses Sehen, gleichsam ein Epiphänomen eines mit dem Sein eingehenden Gewahrseins, das nichts weiter macht, als positiv zu bezeugen, was geschieht. Hier gibt es weder einen Unterschied zwischen aktiv und passiv noch zwischen Handeln und Getriebenwerden, da alles differenzlos als eins (und damit als gleich gut) erscheinen muss.

Hieran anschließend können wir eine weitere Lesart formulieren, nämlich die, dass sich die Evolution – und damit auch mögliche Götter, die als ihre Schöpfer oder Repräsentationen imaginiert werden können – bereits im Übergang von einer einwertigen zu einer zweiwertigen Welt befindet. Wer zu sehen *und* zu bezeichnen anfängt, generiert bereits die Keimzelle negationsfähiger Formen. Von hier aus ist es nur ein kleiner Schritt, zu sehen, dass die Ergebnisse von Schöpfung und Evolution nicht mit den eigenen Imaginationen, dem eigenen Denken und den eigenen Wünschen übereinstimmen, also etwas schiefzulaufen droht. Hiermit kommen wir zum zentralen Bezugsproblem der Spiritualität, das sich vermutlich jedem hinreichend reflexionsfähigen Wesen zu stellen beginnt: *Wie lässt sich unter Bedingungen des Zweifels, der Kritik und der Sorge wieder zu einer Reflexionsposition gelangen, in der Subjekt und Welt einfach nur gut sind?*¹⁴

Um es bereits an dieser Stelle in vorsichtiger Form als Hypothese zu formulieren: Wir vermuten, dass es letztlich darum geht, die eigene

14 Mit Günther lässt sich gut zeigen, wie religiöse Formen, wie sie zum Beispiel in der christlichen Verkündigung oder in den buddhistischen Lehreden auftreten, logisch funktionieren. Es wird deutlich, dass die zweiwertige Logik nicht ausreicht, um die dabei zutage tretenden kommunikativen Prozesse zu verstehen. Diese Formen beruhen abstrakt gesehen darauf, dass eine Welt, wie sie einem zweiwertigen Bewusstsein erscheint, mit den Mitteln der polykontexturalen Logik beschrieben wird, um von hier aus performativ auf Basis transjunktoraler Operationen – also einer geschickten Anordnung von Rejektionen – zu einem einwertigen Gott (etwa im Christentum) oder zu einer einwertigen Form der Erlösung zurückzukehren (so auch in der buddhistischen Rede vom Nirvana; siehe dazu Vogt 2017). Durch die Verwischung der logischen Kategorien (insbesondere des Unterschieds zwischen

Bedingtheit, also die eigenen Mechaniken und hiermit zusammenhängend die eigene Blindheit, zu verstehen, was jedoch nur möglich ist, wenn einem dies in überzeugender Form von signifikanten Anderen gezeigt wird.

*Die Geburt des Beobachters – Subjektivität kann nicht
eins mit der Welt sein*

An irgendeinem Punkt der Evolution beginnt die Schöpfung sich ihrer selbst bewusst zu werden, das heißt Lebensformen zu entwickeln, für die die Spaltung zwischen sich und dem von ihnen gelebten Leben konstitutiv wird. Auf einmal gibt es nicht nur Sehen, Tun, Hören, Empfinden, Jagen, Essen, sondern Beobachter, die etwas sehen, hören und empfinden, und Beobachter, die etwas tun. Anstelle als Sein mit sich selbst genügsam zu sein, beginnt das einwertige, sich mit sich selbst in eins empfindende Leben jetzt in die Welt der zweiwertigen Logik überzugehen, also zwischen Subjekt und Objekt zu scheiden. Das, was ist, erscheint nun nicht mehr selbstevident. Zudem beginnen nun Unterschiede zwischen ›wahr‹ und ›falsch‹ – und hiermit verbunden zwischen ›gut‹ und ›böse‹ – in der Welt wirksam zu werden.

Die selbstversunkenen kybernetischen Kreisläufe, in denen Unterschiede im Sensorium einfach nur Unterschiede im Motorium bewirken, werden aufgebrochen: In hinreichend komplexen kognitiven Systemen wendet sich der Prozess des Unterscheidens reflexiv auf sich selbst zurück und in der Folge werden Identitäten ausgeflaggt, die sich selbst als Beobachter unterscheiden können. Wir kommen so zu einer Kybernetik zweiter Ordnung,¹⁵ in der es jetzt nicht nur um Unterschiede geht, die Unterschiede auslösen, sondern auch um die Konsequenzen der Beobachtung – und die Verantwortung des Beobachters.

Gotthard Günther hat nun aufgezeigt, dass ein solches komplexes kognitives System unweigerlich zu einer dreiwertigen Logik führt (siehe ausführlich Kap. I.2). Der Grund hierfür ist, dass die Reflexion-in-anderes (Objekt) und die Reflexion-in-sich (Subjekt) eine weitere logische Position hervorbringt: das *Du*.

Wie bereits erwähnt, muss das *Du* für das *Ich* als eine logische Chimäre erscheinen, nämlich zugleich als ein materieller Prozess (der sichtbare Körper des anderen) wie auch als eine eigene, dem Wirken und Erkennen des Ichs allerdings nicht zugängliche Subjektivität (die verborgene Innerlichkeit). Neben den dinghaften Objekten (*Seinsidentität*) und der Reflexion auf die eigene subjektive Beobachterposition (*Reflexionsidentität*)

der klassischen und der transklassischen Negation) werden auf diese Weise – zumindest temporär – die Polaritäten einer zweiwertigen Welt aufgehoben.
15 Siehe von Foerster (1992).

erscheint damit eine dem eigenen Sein für immer fremd bleibende subjektive Position (*Transzendentalidentität*).

Aus einer logischen Perspektive ist klar, dass es dieser drei Positionen bedarf, damit sich eine Position der Selbstreflexion stabilisieren kann. Das Ich kann nur in der Begegnung mit dem Du entstehen, wie Martin Buber feststellt.¹⁶ Nicht nur Soziologinnen wissen dies, sondern bereits in vielen religiösen Mythologien wird dieser logische Strukturreichtum abgebildet. Der hinduistische Schöpfergott Brahma braucht als Gegenspieler Shiva, den Weltzerstörer (zudem spalten sich die Gottheiten nochmals in männliche und weibliche Emanationen auf). Die vielen griechischen Götter waren vor allem in ihren Zielen und ihrem Streben uneins untereinander.

Auch im Christentum benötigt Gott allein schon deshalb einen Gegenpart, weil er sich erst so seiner eigenen moralischen und epistemischen Stellung im Universum versichern kann. Moral kann es nur geben, wenn das Sein und das Sollen auseinanderfallen. Der aus der Einheit gefallene Engel Luzifer dient vor diesem Hintergrund logisch gesehen dazu, das hierfür notwendige Beobachterverhältnis einer zweiwertigen (dualistischen) Welt hervorzurufen. Um hier wieder auf Batesons kybernetischen Gott ›Öko‹ zurückzukommen: Auf einmal ist Gott nicht mehr identisch mit all dem, was als Welt erscheint. Es existiert jetzt vielmehr eine Reflexionsposition, aus der heraus Luzifer nicht mehr als logische Konsequenz eines kybernetischen Regelkreises angesehen wird, der zur Welt gehört. Die Reflexion trennt sich vom Sein. Das Erleben wird zweiwertig, beginnt, zwischen sich selbst als Teil der Welt und dem Rest der Welt zu unterscheiden. Letztere enthält mit dem Teufel und seinen Helfern nun auch Wesen, die nicht sein sollen. Spätestens hier bekommt die Bezeichnung ›gut‹ explizit ihren Gegenpol ›schlecht‹. Der Beobachter geht in Differenz zur Welt.

Auf diesem Niveau der reflexiven Kompetenz läuft der sich selbst bewusst werdende Demiurg Gefahr, seiner eigenen Endlichkeit und Ohnmacht gewahr werden zu müssen: »(Alle) Welten, auch die Welt Brahmas, unterliegen der Wiederkehr«, heißt es beispielsweise in der Bhagavad Gita.¹⁷ Sobald ein zur Selbstbeobachtung fähiges Wesen befreit, dass es die Dinge nicht unter Kontrolle hat, beginnt es zu verstehen, dass dies auch für das Fortbestehen der eigenen Existenzform gilt. Über kurz oder lang beginnt es nach einem Ausweg aus dem Dilemma zu suchen. Wir landen wieder bei dem zuvor erwähnten spirituellen Bezugsproblem, um das dann auch die Theodizee-Frage kreist:

Wie lässt sich von hier aus zu einer Reflexionsposition gelangen, in der das Verhältnis zur Welt einfach nur gut ist?

¹⁶ Buber (2002).

¹⁷ Bhagavad Gita 8.16 (zitiert nach: <https://schriften.yoga-vidya.de/bhagavad-gita> [Abruf: 21.11.2022]).

Die Probleme des Subjekts und seine Lösung: sich an die Welt abgeben

Dies ist genau der Punkt, an dem Günther die logischen Konsequenzen des Schöpfungsmythos in einer Weise weiterdenkt, die für manche Christen wohl wie Ketzerei klingen muss. Da der christliche Gott die Sache allein nicht in den Griff bekommt, also aus dem Dilemma der sich seiner selbst bewussten Zweiwertigkeit nicht mehr rausfindet, beginnt auch er darüber nachzudenken, die Verantwortung für die Schöpfung an andere Wesen abzugeben. Er lässt Wesen entstehen bzw. neben ihm gelten, die so sind wie er – reflexionsfähige Menschen, die sich ihrer selbst bewusst sind:

»Der Schöpfer gibt einen Teil an die Welt ab! Das zuletzt Geschaffene gehört nicht mehr dem gegenständlichen Bestand der Welt an. Es ist Subjekt.«¹⁸

Gott erschafft den Menschen nach seinem Ebenbild, was allerdings nicht wörtlich zu nehmen ist, (etwa in dem Sinne, dass wir genauso aussehen würden wie der Schöpfer). Die Analogie bezieht sich, so Günther, vielmehr darauf, dass der Mensch ebenfalls ein subjektives Zentrum hat und damit autonom wird. Das Verhältnis von Schöpfer und Geschöpf kann hier im Sinne einer Eltern-Kind-Beziehung verstanden werden: Die »Zeugung des Kindes« ist ein »schöpferischer Akt«, in dem ein natürlicher »Prozess zu einer Verdopplung der subjektiven Identität führt. Denn das Kind erwirbt ein autonomes Ich, aus dem sich die Subjektivität des Elternpaars endgültig zurückgezogen hat.«¹⁹

Sobald aber mehrere Subjektivitäten nebeneinander bestehen und miteinander koexistieren, kann Gott mit dem Schöpfungsprozess nicht mehr identisch sein. Bereits mit Luzifer kommt es zu einer Verdoppelung. Spätestens mit dem Menschen geschieht nochmals eine Vervielfältigung der subjektiven Positionen. Gott wird seinerseits Beobachter unter Beobachtern und fällt damit aus der Einheit heraus.²⁰ Die Evolution beginnt auf einmal, vielfältige subjektive Zentren hervorzubringen. Unweigerlich landen wir in einer polykontexturalen Welt, die keine zentrale Steuerungsinstanz mehr hat.²¹

18 Günther (2008, S. 34).

19 Günther (2008, S. 35).

20 Und wenn man hier ›Gott‹ mit ›Öko‹ ersetzt, sieht man, dass die vielen Subjektivitäten der Menschen nicht unweigerlich mit der ökologischen Perspektive einhergehen müssen.

21 Umgekehrt gilt entsprechend: Wenn eine göttliche Wesenheit existiert, dann nur noch als ferne, unzugängliche, transzendente Instanz: »Denn darin besteht die absolute Transzendenz Gottes, dass er mit jener Setzung des Menschen als Abbild des Schöpfers sich selbst aus jener demiurgischen Identität zurückgezogen hat. Gott ist mit dem Sein-überhaupt nur vor der Erschaffung der Welt

Spirituell und mit Blick auf eine zeitgemäße Metaphysik gehen hiermit zwei entscheidende Punkte einher:

1. Die subjektive Perspektive kann nicht allwissend sein, denn sie ist ja eben subjektiv und hat entsprechend nur begrenzt Zugang zu dem, was in der Welt geschieht.
2. Das subjektive Erleben kann nicht mehr eins mit der Welt sein, denn es geht ja in Differenz zur Welt. Für das Subjekt muss damit jedes Du als andere Subjektivität unzugänglich bleiben. »ein[en] ebenso unzugängliche[n] Raum« darstellen »wie die mythologischen Dimensionen, in denen die himmlischen Heerscharen schweben.«²²

Der »Weltinnenraum« (Rilke)²³ hinreichend entwickelter kognitiver Systeme zeichnet sich folglich gerade dadurch aus, nicht voll spezifiziert zu sein. Auch der christliche Gott muss uns damit als Metapher für eine Evolution erscheinen, die einen Beobachter hervorbringt, der um sich weiß, aber nicht alles weiß.²⁴ Hiermit lässt sich an dieser Stelle eine ungewohnte Brücke zwischen Theologie und Kybernetik schlagen:

- Informationstheoretische Überlegungen zur Kombinatorik möglicher Weltzustände zeigen schnell auf, dass ein Organismus nicht einmal annähernd über die kognitiven Fähigkeiten zur Berechnung des Verlaufs der komplexen Beziehungen seiner Welt verfügen kann. Ebenso wenig lässt sich vorhersehen, was andere autonome Lebensformen tun werden. Eine subjektive Existenz zu haben heißt ja gerade, dass man nicht so recht wissen kann, was der Fall ist. Deshalb bleibt einem nichts anderes übrig, als seine eigene Wirklichkeit zu konstruieren und danach zu handeln. Dies mag wiederum andere Wesen überraschen und zu Reaktionen bringen, die dann aus der eigenen subjektiven Perspektive nicht vorhersehbar sind.
- Die monokontexturale Göttlichkeit des Demiurgen wird innerhalb des Weltgefüges allein schon deshalb gebrochen werden, weil, so Erwin Schrödinger, »bestmögliches Wissen um ein Ganzes« nicht mehr »notwendig das Gleiche für seine Teile« einschließt.²⁵ Epistemisch zerfällt die Welt damit in eine Vielzahl subjektiver Zentren, wo jeweils etwas

(als dem partikulär Seienden) voll identisch. Im Schöpfungsprozess löst er sich aus dieser Identität ab und west von nun an jenseits ihrer als *deus absconditus*.« (Günther 2008, S. 40; kursiv im Original)

22 Günther (1975, S. 61 f.).

23 Zu den Einflüssen der religiösen Mythologien Indiens und des Christentums auf Rainer Maria Rilke siehe von Brück (2015).

24 Um es nochmals mit Hegel (1986 [1817], § 147) zu formulieren: »Dahingegen ist der christliche Gott der nicht bloß gewußte, sondern schlechthin sich wissende Gott und nicht bloß vorgestellte, sondern vielmehr absolut wirkliche Persönlichkeit.«

25 Schrödinger (1935, S. 849).

anderes der Fall ist und wo auf dieser Grundlage eine Eigenaktivität entwickelt wird. Es gibt keinen wie auch immer gearteten Gottesstandpunkt mehr, von dem aus entschieden werden könnte, was ›wirklich‹ ist – oder wie es ›wirklich‹ ist.

Subjektivität entsteht immer dann, wenn einer Lebensform ihr Unterscheidungsgebrauch reflexiv zugänglich wird und die Einheit von Erkennen und Tun für sie in der Folge auseinanderzufallen beginnt. In der biblischen Mythologie geschieht dies mit dem Sündenfall. Vom Baum der Erkenntnis zu essen heißt, zu beginnen, zwischen der Reflexion auf die eigene subjektive Position und der Reflexion auf die Welt zu oszillieren. Es wird nun möglich, sich selbst und seine Stellung in der Welt zum Thema der Reflexion zu machen. Die Beziehung zwischen Sein und Sollen, zwischen Erkennen und Handeln sowie zwischen Theorie und Praxis erscheint jetzt als eine gebrochene bzw. gespaltene. Semiotisch gesprochen erscheint die *Erstheit* der unmittelbaren Sinneserfahrung nun gebrochen, indem im Erleben zusätzlich auf sich selbst verweisende Zeichenprozesse entstehen. Mittels dieser können hypothetische Ich-Identitäten und die ihnen zugehörigen imaginären Welten konstruiert werden. In dieser *Zweitheit* beginnt der »Weltinnenraum« (Rilke) nun, sich in Spannung zu sich selbst zu erleben. Die Reflexion tritt in Distanz zum Sein. Schließlich wird in der *Drittheit* diese Spannung wiederum durch Praxis, nämlich durch *Wollen* – durch intendiertes Handeln –, überschritten.²⁶ Daraus ergibt sich eine Reihe von Konsequenzen:

1. Die Projektionen müssen nicht mit dem Faktischen einhergehen.
2. Der Reflexionsprozess kann sich seiner eigenen Unzulänglichkeit und des eigenen Nichtwissens bewusst werden.
3. Die Nebenfolgen der eigenen Handlungen sind nicht überschaubar.
4. Die Wesenheit wird damit konfrontiert, dass sie etwas verursacht hat, was sie nicht gewollt hat.
5. Der Reflexionsprozess führt schließlich zu der Einsicht, dass nicht nur das eigene Handeln, sondern auch die eigene Lebensform als solche scheitern kann, ja letztlich vergänglich und sterblich ist.

Das spirituelle Joch des modernen Menschen – oder der ›Sündenfall‹ – besteht damit genau darin, um die nicht beherrschbaren Konsequenzen der eigenen Existenz zu wissen, sich aber dennoch im Sinne einer »schuldlosen Schuld«²⁷ für sie verantwortlich fühlen zu müssen. *Nolens volens* ist man das Subjekt des Geschehens. Wer vom »Baum der

²⁶ Der Clou der peirceschen Semiotik besteht darin, dass *Erstheit*, *Zweitheit* und *Drittheit* letztlich als eine Einheit zu fassen sind, nämlich als eine unmittelbare Praxis des Werdens, die sich als ein fortschreitender Zeichenprozess entfaltet und ausdifferenziert. Siehe Peirce (1991, S. 358 ff.).

²⁷ Nishitani (1986, S. 389).

Erkenntnis des Guten und Bösen«²⁸ gegessen hat, weiß also immer noch nicht, was gut und richtig ist. Doch er oder sie verfügt nun – wie die Kybernetiker sagen würden – über ein Wissen zweiter Ordnung, das den Wissenden mit einer folgenschweren Kognition belastet, nämlich dem Gewissen: Man kann prinzipiell nicht wissen, was die Folgen des eigenen Erkennens und Handelns sind, hat sie jedoch zu ertragen. Um es mit den Worten Günthers zu formulieren:

»Damit verändert sich die metaphysische Existenz des Menschen in der Welt ganz radikal. In der neuen Realität, die seine Handlung provoziert, kann nicht gesagt werden, dass sie ›sehr gut‹ ist. Das ist nicht im Sinn einer nörgelnden Besserwisserei gemeint. Es wäre nämlich genau so falsch zu sagen, dass sie ›sehr schlecht‹ und vom Teufel ist. Das ist der metaphysische Irrtum aller Utopisten: [...] Nun mag diese Handlung zwar von einem Gewissen, das ›weiß‹, was gut und böse ist, sich orientiert haben, aber die Realität, die durch diese Handlung gesetzt ist, ist ›gewissen-frei‹. In ihrer Existenz ist das moralische Motiv, das im Rücken des Handelns liegt, völlig irrelevant.«²⁹

Ein intelligentes subjektives Zentrum zu haben, bringt es mit sich, die Diskrepanzen zwischen Wollen und Realität, Sein und Sollen, Anspruch und Wirklichkeit zu reflektieren, ohne jedoch die mit ihnen verbundenen Agonien prinzipiell aufheben zu können.

Insofern sich also eine Reflexionsidentität als Subjektivität hervorgebracht hat, steht sie vor dem Problem, dass die hiermit einhergehenden kognitiven Strukturen – die im Reflexionsprozess entwickelten Kategorien, Begriffe und Normen, die ja ihre Identität ausmachen – ab einem gewissen Punkt nicht mehr zur realisierten Lebensform passen. Es droht damit eine Krise durch Weltverlust. Denn dem Subjekt, das sich bislang erfolgreich gegenüber der Welt behauptet hat, bleibt in Bezug auf die Reproduktion der eigenen Prozesse nur noch übrig, die immer blasser werdenden und weniger überzeugenden eigenen Konstruktionen zu ergreifen. Dies lindert jedoch die Probleme der krisenhaften Selbst- und Weltverortung nicht, sondern verstärkt sie sogar noch. Der solchermaßen gebrochene und unter Druck geratene Mensch strebt nach Transzendenz, also nach einem Weltverhältnis, das das bisherige Selbstverhältnis überschreitet.

Die gesuchte Transzendenz kann allerdings nicht an einem Ort außerhalb dieser Welt gefunden werden – sie ist nicht esoterisch –, sondern beruht umgekehrt gerade auf der Veräußerlichung subjektiver Reflexionsleistungen in den kulturellen Institutionen dieser Welt.³⁰ Zu nennen sind hier Sprache, Normen, Sitten und Gebräuche, Gesetze – schlichtweg

28 1. Mose 2.

29 Günther (2008, S. 32).

30 Dies hat insbesondere Émile Durkheim (1994) deutlich gemacht.

alle semantischen und soziostrukturellen Ordnungen menschlicher Gemeinschaft und Gesellschaft. Die hiermit einhergehenden Institutionen haben sich insofern bewährt, als sie über das einzelne Individuum hinweg gelten und selegiert worden sind und sich damit auch in kultureller Hinsicht als nachhaltig erwiesen haben.³¹

Religionen wie die unterschiedlichen Varianten des Christentums, des Islams, des Hinduismus oder des Buddhismus konnten in diesem Sinne über Jahrhunderte hinweg eine mehr oder weniger stabile Matrix liefern, an die sich die Subjekte orientieren konnten, um die oben aufgeworfenen spirituellen Bezugsprobleme abzumildern. Laut Arnold Gehlen waren es gerade diese alten gesellschaftlichen Institutionen, die den instinktarmen Menschen vor den Problemen der Subjektivität und der mit ihr verbundenen flottierenden Freiheitsgrade entlastet haben.³² Damit gelang es zum einen, die Probleme der menschlichen Freiheit zu adressieren – darunter Fragen der eigenen Unzulänglichkeit und die aus dieser erwachsenden Probleme von Schuld und Verantwortung. Zum anderen konnten Ordnungen der Selbstbegrenzung und -bescheidung etabliert werden, die allein schon deshalb eine gewisse seelische Stabilität boten, weil sie eben lehrten, dass vom diesseitigen Leben – also der weltlichen Sphäre – nicht mehr erwartet werden kann und darf, man sich also besser mehr schlecht als recht in den irdischen Mängeln einrichtet, anstelle dagegen zu kämpfen.³³

Die Krise menschengemachter Ordnungen

Wie auch immer man heutzutage zu ihnen stehen mag, die mit religiösen Semantiken veräußerlichten Regulative gestatteten dem Menschen nicht

31 Die Objektivierung immer nur subjektiver Erscheinungen kann einzig über den ›Umweg‹ des Sozialen gesichert werden – im Guten wie im Schlechten. Die Inkorporierung des sozial angelieferten Sinns der Gesellschaft bzw. Kultur erscheint damit gleichsam als soziale Härtung. Wir treffen hier auf den soziologischen Teil von Günthers Denken.

32 Gehlen (1963).

33 Die funktionale Ausdifferenzierung von Kulturformen wie Religion und Kunst kann hiermit als Strategie des Umgangs mit Sinnüberschüssen gesehen werden. Ab einem gewissen Reflexionsniveau kann die Welt auf verschiedenen logischen Ebenen bejaht und verneint werden, womit sich die Frage stellt, wie auch diese Prozesse institutionalisiert, das heißt in einer Weise gebunden werden können, die die Ökologie des Gesamtsystems nicht gefährdet. Um dieses Thema kreisen auch Albert Camus' Studien in *Der Mensch in der Revolte*. Wir begegnen hier einerseits einer Metaphysik, die Negation und Affirmation in fataler Weise gegen den Menschen und seine Welt wendet. Andererseits begegnen wir mit der Kunst Formen, zugleich Ja und Nein zu sagen, ohne in solch destruktive Formen verfallen zu müssen (Camus 2016 [1951]).

nur eine mehr oder weniger austarierte Balance von Natur, Gesellschaft und Individuum. Sie schufen auch eine gemeinsam geteilte kulturelle Basis, auf die sich dann – in welcher Form auch immer (affirmativ oder kritisch) – bezogen werden konnte. Man verfügte damit auch über eine wirtschaftliche, politische, religiöse und kulturelle Ordnung, die gleichsam objektiv gegeben schien, da sie von kollektiv geteilten Mythen und Erzählungen getragen war. Von Gott, den Göttern oder durch ein unpersönliches metaphysisches Prinzip (etwa das Karma-Gesetz) gegeben, erschien diese Ordnung zugleich absolut wie auch gut und gerecht. Subjektiv und privat konnte man damit hadern, doch dies änderte nichts an der objektiven (kulturellen) Geltung ihrer Institutionen.

Hieran schien auch die »Aufklärung« zunächst nicht grundsätzlich etwas zu ändern. So ließ sich – etwa mit Kant – der Glaube an eine objektive Vernunft und eine selbstevidente Moral pflegen, was dem Menschen nach seinem »Ausgang [...] aus seiner selbst verschuldeten Unmündigkeit«³⁴ weiterhin Halt und Sicherheit versprach. Denn insofern wir Gott als eine Institution ansehen, der sich der Mensch in der Einsicht in die Notwendigkeit freiwillig unterzuordnen hat, haben sich hiermit zunächst nur die Begriffe verschoben, nicht jedoch die zugrunde liegenden Relationen. Die Vernunft nimmt jetzt den Platz Gottes ein und wird ihrerseits zur tragenden Institution. Der menschengemachte Fortschritt wird zum kollektiven Mythos. Der Kapitalismus wird zur Religion,³⁵ da das Land, wo Milch und Honig fließen, jetzt mit den Imaginären einer durch die unsichtbare Hand geführten Gesellschaft verbunden wird. Sofern wir mit Émile Durkheim die Funktion der Religion darin sehen, dass Menschen sich mit ihrer Hilfe eine gesellschaftliche Ordnung aufbauen, die als eine selbstgeheiligte Wirklichkeit erscheint,³⁶ wird die Religion mit der Moderne nicht grundsätzlich aufgegeben, sondern verlagert sich nur in zivilreligiöse Formen.

Spätestens mit Friedrich Nietzsche beginnen diese metaphysischen Sicherheiten jedoch brüchig zu werden. Die Aussage »Gott ist tot. Wir haben ihn getötet.«³⁷ erschüttert nicht nur die göttliche Ordnung des Christentums, sondern jegliche Vorstellung einer objektiven ordnunggebenden Instanz. Wer von der Vernunft redet oder eine allgemeingültige Moral beansprucht, macht sich jetzt entweder eigenem, unbewusstem Machtstreben verdächtig oder zeigt damit nur, dass er oder sie nicht die Charakterstärke hat, der hiermit einhergehenden Heuchelei offen ins Auge zu schauen.³⁸ Mit Blick auf die vielfältigen subjektiven Zentren

34 Kant (1784).

35 So Dirk Baecker (2003) in Referenz auf Walter Benjamin.

36 Durkheim (1994).

37 Nietzsche (2021 [1887], Absatz 125).

38 Nietzsche (2009 [1887]).

und ihre jeweils spezifischen Kontexte gibt es nicht mehr nur die eine Moral und die eine Vernunft, sondern derer viele.³⁹

Die religiösen und moralischen Ordnungen der Welt erscheinen damit nicht mehr als eine durch Kultur und Tradition ein für alle Mal gegebene objektive Natur, sondern werden nun selbst als vom Menschen gemacht erfahren. Sie erscheinen veränderbar und damit in hohem Maße kontingent. Um es mit Günther zu formulieren: Was der Mensch »früher im Innenraum der Subjektivität als Vorstellungen und Phantasien manipulierte«, wird »jetzt im Außenraum der physischen Welt noch einmal« aufgeworfen, »um den einstigen Reflexionsvorgang als Willensakt zu wiederholen. [...] Das ist Handlung zweiter Ordnung. Sie definiert eine Aktionskapazität, die sich einer Außenwelt zuwendet, die mit dem abgelebten Reflexionsbestand der vorangegangenen historischen Epoche belastet und bereichert ist.«⁴⁰

Die Geschichte und die menschlichen Kulturformen erscheinen nun »nur noch« als Tradition und werden damit ihrerseits zum Spielmaterial der Reflexion und von Versuchen der Neugestaltung. Man handelt und agiert nicht mehr innerhalb des gegebenen Rahmens, sondern der Rahmen selbst wird zum Thema der Reflexion und des Gestaltens. Das reflexiv potenzierte Ich steht nun einer Außenwelt gegenüber, die historisch aufgeladen ist – dies wird nun zugleich als belastend wie auch als bereichernd erfahren.

Als Konsequenz einer solchermaßen gesteigerten Subjektivität wird es jedoch zugleich immer schwieriger, eine gemeinsame, mit anderen Menschen geteilte Basis der Wahrnehmung, Willensbewegung und Situationsdeutung zu finden. Umso individueller die Lebensformen werden, in die sich das Subjekt im eigenen Handeln projizieren kann, desto ferner erscheint die Du-Position einer fremden Subjektivität, die ihrerseits ihre eigene Lebensform zu gestalten sucht.

In den vormodernen Kulturen, so Günther, habe es »ein von allen geteiltes instinktives psychisches Grundwissen« gegeben und entsprechend sei man »in wesentlichen Bereichen seelisch aufeinander abgestimmt« gewesen. »Es gab eine wirksame Form transsubjektiver Innenkommunikation vom Ich zum Du, die sich sehr evidenter Erlebnistatbestände

39 Mit Blick auf die Gesellschaftsstruktur geht dies nach Luhmann (1998c) mit der funktionalen Differenzierung der Gesellschaft einher. Im Zuge dessen beginnen sich auch Organisationen zu diversifizieren, vielfältige Berufsbilder, Rollen und Wertsphären entstehen, die in Hinblick auf ihre jeweilige Moral und Rationalität inkommensurabel erscheinen. Dadurch entsteht eine Unmenge von subjektiven Sprecherpositionen, die nicht nur darum ringen, hinsichtlich ihrer Identitätsansprüche gehört zu werden, sondern einen Platz in den gesellschaftlichen Institutionen beanspruchen (und entsprechend die Verrechtlichung ihrer Ansprüche einfordern).

40 Günther (2008, S. 33 f.).

bediente.« Diese »selbstverständliche Voraussetzung aller vergangener Geschichte«, so Günther weiter, existiere heute jedoch nicht mehr. Wir würden heute vielmehr »in steigendem Maße Bewusstseinsituationen« begegnen, »denen wir auf keine Weise ausweichen können, die aber alle allgemeinverbindliche Evidenz für die erlebende Subjektivität verloren« hätten. »Unsere Existenz« zwingt »uns an allen Orten Reflexionsmechanismen auf, die für das Ich völlig undurchsichtig« erschienen, da sie »zwar zwangsläufig, aber ohne jede Evidenzkraft« seien.⁴¹

Wir leben in einer Gesellschaft, in der an unterschiedlichsten Orten jeweils anderes passiert, ohne dass die einzelnen Geschehnisse sozial, zeitlich und sachlich miteinander koordiniert oder durch einen übergreifenden Sinnhorizont verbunden werden. Armin Nassehi spricht in diesem Zusammenhang von einer »Gesellschaft der Gegenwart«⁴², um damit zu pointieren, dass hochspezifische Rechtsvorschriften, Finanzströme, unzählige, oftmals esoterisch anmutende wissenschaftliche Erkenntnisse, die Artikulationen in massenmedialen Blasen, die vielfältigen Formen der Liebe und das, was jeweils vor Ort an Interaktion passiert und psychisch an Sinn oder Bedeutung realisiert wird, nicht einmal ansatzweise zu einem kohärenten Ganzen zusammenfinden.

Wir mögen um die Pluralität der Sinngefüge zumindest vage wissen, da wir alltäglich mit ihren Wirkungen konfrontiert werden. Doch all dies wird immer weniger als »meine« Welt, als der immer auch gefühlsmäßig erfahrene Bezugspunkt des eigenen In-der-Welt-Seins und des damit verbundenen Eine-Welt-Habens erlebt.⁴³

Der Bewusstseinsprozess kommt hiermit in eine Formkrise. Während die private Subjektivität sich früher an den feststehenden *objektiven Inhalten* der Tradition orientieren und sich gegebenenfalls von ihnen abstoßen konnte (beispielsweise konnte die Religion unabhängig davon, ob man an Gott glaubte oder zweifelte, die Matrix der Auseinandersetzung zur Identitätsbildung bilden), werden die Gegenstände der Kultur nun selbst so stark disponibel, dass sie auszudünnen beginnen. Ob man Christ oder Buddhist ist, sich seine eigene Spiritualität aus dem Flickenteppich moderner Glaubensangebote zusammenbastelt, lieber zivilreligiöse Angebote wie den Marxismus präferiert oder in eine agnostische oder gar nihilistische Metaphysik einrastet, wird in den meisten Lebensbereichen kaum mehr einen Unterschied machen (die Vermieterin oder den Lehrer interessiert in der Regel nicht mehr, was man glaubt, sondern ob man die Miete zahlt oder im Unterricht mitmacht).⁴⁴

41 Günther (2008, S. 36 f.).

42 Nassehi (2011).

43 Günther (2008, S. 36 f.).

44 Juli Zeh hat die Formprobleme in Bezug auf den Verlust von Koordinatensystemen von Werthaltungen in ihrem Roman *Spieltrieb* durchdekliniert: »Was,

Das Selbst beginnt dadurch mehr und mehr zu fragmentieren. Wir begegnen dem Paradoxon, dass die inhaltliche Freiheit und die Freiheit der vermeintlichen Selbstbestimmung nicht zugleich mit der erfahrenen Freiheit einhergehen, sein Leben in einer Weise zu gestalten, die einen subjektiv in der einen oder anderen Form mit sich selbst identisch fühlen lässt. Denn zwangsläufig droht auch die leiblich empfundene Subjektivität immer dünner und leerer zu werden, je mehr die Gegenstände und Kontexte verblassen, die als Medium genutzt werden können, um sich auf reproduzierbare und beständige Weise in-sich und in-anderes reflektieren zu können. Man muss nicht gleich den »Untergang des Abendlandes«⁴⁵ verkünden, um das spirituelle Bezugsproblem zu verstehen, das sich mit Günther nochmals wie folgt reformulieren lässt:

»Je weiter die Geschichte vorschreitet, desto mehr [...] differenziert sich das Ich in seiner Innerlichkeit. In den historischen Endstadien aber ist die prästabilierte Harmonie der privaten Bewusstseinsvollzüge endgültig verschwunden und die institutionelle Struktur des gesellschaftlichen Lebens beginnt sich aufzulösen.«⁴⁶

Unter diesen reflexionstheoretischen Voraussetzungen wird auch die metaphysische Verortung des Menschen zu einer neu zu bestimmenden Aufgabe. Was kann dem Ego noch Halt geben, wenn die Institutionen immer brüchiger werden, an die es sich halten bzw. von denen es sich abstoßen kann.⁴⁷ Die moderne Gesellschaft bringt zwar Institutionen in An-

wenn die Urenkel der Nihilisten längst ausgezogen wären aus dem staubigen Devotionalienladen, den wir unsere Weltanschauung nennen? Wenn sie die halb leereräumten Lagerhallen der Wertigkeiten und Wichtigkeiten des Nützlichen und Notwendigen, Echten und Rechten verlassen hätten, um auf Wildwechseln in den Dschungel zurückzukehren, dorthin, wo wir sie nicht mehr sehen, geschweige denn erreichen können? Was, wenn ihnen Bibel, Grundgesetz und Strafrecht nie mehr gegolten hätten als Anleitung und Regelbuch zu einem Gesellschaftsspiel? Wenn sie Politik, Liebe und Ökonomie als Wettkampf begriffen? Wenn »das Gute« für sie maximierte Effizienz bei minimiertem Verlustrisiko wäre, »das Schlechte« hingegen nichts als ein suboptimales Resultat? Wenn wir ihre Gründe nicht mehr verstünden, weil es keine gibt?

Woher nähmen wir dann noch das Recht zu beurteilen, zu verurteilen, und vor allem – wen? Den Verlierer des Spiels – oder den Sieger? Der Richter müsste zum Schiedsrichter werden. Mit jedem Versuch, Erlerntes anzuwenden und Recht in Gerechtigkeit zu übersetzen, würde er sich der letzten verbliebenen Todsünde schuldig machen: der Heuchelei.« (Zeh 2007, S. 7)

45 Spengler (2007 [1923]).

46 Günther (2008, S. 38 f.).

47 Man könnte jetzt davon sprechen, dass wir es hier mit einer Diffusion und Auflösung von Kultur im traditionellen Sinne zu tun haben, womit sich die Frage stellt, wie und ob diese Dekontextualisierung selbst nicht wieder auf die eine oder andere Weise institutionalisiert wird. Die Auflösung, das

zahl und Umfang wie nie zuvor in der Geschichte hervor. Praktisch zu jedem Thema, jeder Weltanschauung und jedem Politikfeld gibt es Organisationen, die diesbezügliche Agenden pflegen und die partikularsten Ansprüche vertreten (man denke an die vielfältigen Lobbyverbände, Bildungsinstitute, NGOs, religiösen und weltanschaulichen Gruppen). Da es derer aber eben so viele gibt und sie in ihren Werten häufig inkommensurabel sind, finden Kultur und Gesellschaft nicht einmal ansatzweise zu einer allgemeingültigen Werthaltung zusammen. Die Gesellschaft wird fragmentierter, differenziert sich zunehmend in weitere Sphären aus, die in Hinblick auf das, was als wichtig und bedeutsam erscheint, ihre eigenen Blasen bilden.⁴⁸ Die *eine* Orientierung (zum Beispiel am christlichen Weltbild) ist nicht mehr erkennbar.⁴⁹

Mit Nietzsche stellt sich hiermit das Problem des Nihilismus. Zudem droht im Anschluss an Hegel die Gefahr einer Geschichtsphilosophie, die das Subjekt zugunsten eines objektiven Weltgeistes tilgt, um damit jeglichen Schrecken im Namen der Faktizität oder des vermeintlichen Endziels der Geschichte rechtfertigen zu können.⁵⁰

Infragestellen gehört damit gewissermaßen zur neuen kulturellen Norm. Es wäre zu vermuten, dass im Zuge dessen gleichsam ein neuer Kampf der Kulturen entsteht – jedoch einer, der inhaltlich nicht mehr durch lebensweltliche Lagerungen fundiert ist, die unterschiedliche Milieus prägen. Es geht nicht mehr um Positionen wie konservativ/liberal, religiös/säkular, links/rechts, Wirtschaft gegen Soziales etc., sondern um situative Bewegungen zwischen Kontextualisierung und Diskontextualisierung bzw. zwischen beliebig konstruierbaren Identitätsdogmen (beispielsweise ›Impfgegner‹ zu sein) und der (Gegen-)Bewegung bzw. der Auflösung ebendieser. Die Glaubensinhalte werden damit beliebig, wie es Günther (2000) bereits mit Blick auf die amerikanische Religiosität der 1950er Jahre vermutet hat.

48 Die Echo-Kammer-Produktion durch und mit personalisierten KIs könnte die hieraus resultierende Entfremdung zunächst noch weiter verstärken.

49 Nishitani (1986, S. 356 f.) beschreibt diese kulturelle Lagerung mit folgenden Worten: »Mit dem Zusammenbruch des teleologischen Systems einer göttlichen Weltordnung zerfiel auch die damit gegebene Hierarchie der Werte. Fortan gingen die unterschiedlichen menschlichen Bestrebungen auseinander, jede einzelne wurde Zweck ihrer selbst, wurde ihr eigener Herr und geriet damit zu etwas, dem durch nichts mehr Einhalt zu gebieten war. Mit anderen Worten: Unendliche Impulsivität trat als etwas Zielloses auf. Das entspricht der Situation einer Welt ohne Gott, in der ›Zeit‹ nach beiden Richtungen hin offen steht.«

50 Siehe zur Auseinandersetzung mit Nietzsche, Hegel und Marx und den totalitären Bewegungen des 19. und 20. Jahrhunderts insbesondere *Der Mensch in der Revolte* von Albert Camus (2016 [1951]).

*Antworten auf die Auflösung der institutionellen Struktur
des gesellschaftlichen Lebens*

Welche Antworten lassen sich auf das Vakuum eines frei flottierenden Reflexionsprozesses finden? Als historische Antworten sind hier zunächst all die zivilreligiösen Bewegungen zu nennen, in denen sich eine menschliche Führungspersönlichkeit im Namen eines selbstgeschaffenen absoluten Prinzips an die Stelle Gottes setzt.⁵¹ In diesem Sinne sieht Günther die »in der Neuzeit aufkommenden totalitären Staatsformen« als »ein[en] erste[n] Versuch« an, das spirituelle Bezugsproblem zu lösen, also »mit der tödlichen Gefahr der freigelassenen Subjektivität fertig zu werden«⁵².

Der Totalitarismus beruht (wie seine zeitgenössischen populistischen Varianten) darauf, spirituell entleerte Menschen gleichzuschalten, indem ihre »schon längst entpersonalisierte Subjektivität an eine Instanz« (sprich einen Führer oder eine Partei) abgegeben wird, »von der angenommen wird, dass sie Kraft ihrer besonderen Situation nicht von der allgemeinen Geschichtsunfähigkeit betroffen ist«⁵³. Doch die auf diese Weise aufgerufene Führungsinstanz kann sich ihrerseits nur auf die »Pseudo-Objektivität« der von ihr erschaffenen Geschichtsmythen berufen. Aus kybernetischer Perspektive muss sie damit als eine ebenso »hilflose Subjektivität«⁵⁴ gelten wie die von ihr geführten Untertanen. Auch ihre Existenz beruht auf selektiver Blindheit. Sie kann zwar die Einheit von gesellschaftlichem Sein und individuellem Bewusstsein populistisch behaupten (etwa indem sie die Imaginäre der Nation anruft), ist jedoch faktisch nicht mehr in der Lage, die verbindenden Inhalte und Institutionen nachhaltig zu sichern.

Eine andere Antwort auf das spirituelle Bezugsproblem des modernen Menschen deutet Helmut Schelsky mit der Frage an, ob Dauerreflexion institutionalisierbar sei:⁵⁵ Wenn die kulturellen Inhalte und Reflexionsgehalte immer dünner und schwächer werden, sodass sie kaum noch Menschen aneinander binden können, dann kann vielleicht zumindest noch der Reflexionsprozess selbst in einer Weise institutionalisiert werden, dass Menschen sich wechselseitig aufeinander beziehen und wenigstens über das Verfahren, an dem sie teilnehmen, zu einer geteilten Praxis gelangen. Schelsky stellt am Beispiel protestantischer Religionsgemeinschaften fest, dass Formen wie das Gespräch, die Diskussion und

51 Siehe als elaboriertes Beispiel Carl Schmitts (2009 [1922]) *Politische Theologie*, nach der nun derjenige der neue Souverän ist, der über den Ausnahmezustand bestimmt.

52 Günther (2008, S. 39).

53 Günther (2008, S. 39).

54 Günther (2008, S. 39).

55 Schelsky (1957).

die Tagung gerade deshalb zu unverzichtbaren Institutionen der Kirche geworden sind, weil die flottierende Subjektivität sich zumindest an die Verfahren der jeweiligen Dialogformate binden kann. Hiermit braucht die religiöse Praxis dann weniger von den geteilten Dogmen getragen zu werden denn von der Tatsache, dass man miteinander über religiöse Gefühle reden und jedenfalls hieran Halt finden kann (was dann selbstredend auch Auswirkungen auf das Verhältnis von Organisationsform und Dogmatik und auf die Beziehung zwischen Individuum und Glauben hat).⁵⁶

Eine andere Form, das Leiden an der frei flottierenden Subjektivität zu lindern, liefern die Versenkungsübungen indischen Ursprungs. Eine »geschichtlich unwirksame und nur für das isolierte Individuum verbindliche Technik der Beherrschung der Dauerreflexion existiert bereits«, formuliert Günther. »Ihr klassisches Textbuch sind die Yogasutras des Patanjali.«⁵⁷ Die buddhistischen Lehren elaborieren dies weiter und überführen es in praktikable Schulungssysteme.⁵⁸

Indem die Aufmerksamkeit auf die Natur des geistigen Prozesses selbst gerichtet wird, verflüssigen sich die Gegenstände des Bewusstseins, der Reflexionsprozess kann dabei jedoch noch Halt im immerfort aufscheinenden Prozess der Veränderung finden, um auf diese Weise geläutert – also von weitergehenden Sinnansprüchen befreit – in die Banalität des Alltags zurückzukehren.⁵⁹ Der Bewusstseinsprozess beginnt sich

56 Dies kann auch erklären, warum die christliche Seelsorge zunehmend die Nähe zu psychotherapeutischen Methoden sucht, um die Problemlagen moderner Subjektivität dann nicht mehr theologisch, sondern mithilfe von Psychotechniken einzuhegen. Nicht mehr die Verkündung des Wortes, sondern das Verfahren zählt. Oder um es mit Luhmann (1993a, S. 366) zu formulieren: »[D]ie moderne Gesellschaft [ist] mehr, als man gemeinhin denkt, durch Emotionen gefährdet. Zum einen werden Individuen veranlaßt, über sich selbst und ihre Probleme zu reden. Wenn akzeptiert wird, daß ein Individuum seine Ansprüche nicht nur auf Verdienste, sondern auch und vor allem auf sich selbst gründen kann, muß es Selbstbeschreibungen anfertigen. [...] Das Individuum nötigt sich zur Reflexion und zur Selbstdarstellung (die nie ›stimmen‹ kann). Es kommt damit in Schwierigkeiten, sucht Hilfe und entwickelt den Zusatzanspruch auf verständnisvolle, wenn nicht therapeutische Behandlung seiner Ansprüche. Dieser letzte Anspruch auf Hilfe bei der Fundierung von Ansprüchen ist so absurd, daß es ebenso möglich ist, ihn anzuerkennen, wie ihn abzulehnen.«

Zur kritischen Auseinandersetzung mit den Selbstbeschreibungen der Psychotherapie, die genau dies nicht reflektieren, siehe auch Vogd (2013).

57 Günther (2008, S. 23).

58 Das Kloster, der Mönchsorden und das Meditationsretreat werden so zu anerkannten Institutionen der Dauerreflexion, die genau darauf ausgerichtet sind, das eigene Selbst- und Weltverhältnis zu beobachten und zu erforschen.

59 Siehe hierzu Vogd (2014b).

hiermit gleichsam nach innen hin abzulösen und sich so von einer für ihn selbst problematisch werdenden Ichhaftigkeit zu befreien. »Auf diesem Weg aber betreten wir das Gebiet der Mystik anstatt das der Geschichte«, formuliert Günther, um hierin gewissermaßen die »heilige Variante der Anarchie des Geistes« zu entdecken.⁶⁰ Hierbei muss allerdings, wie Ernst Tugendhat aufzeigt,⁶¹ die Frage offenbleiben, ob die mitfühlende Rückkehr in eine Welt gelingt, in der andere Wesen existieren und wahrgenommen werden, oder ob die mystische Übung in einem autologisch geschlossenen Selbstverhältnis verbleibt.⁶²

*Maschinen, die uns helfen, das Subjekt wieder einzufangen
und zu zähmen*

Kommen wir nun zu einer anderen, auf den ersten Blick recht ungewöhnlichen Lösung, die Dauerreflexion des modernen Subjekts zugleich zu institutionalisieren wie auch zu zähmen. In Analogie zum biblischen Schöpfungsmythos schlägt Günther vor, die kognitiven Prozesse, die unsere menschliche Subjektivität ausmachen, zu entäußern – und zwar auf Basis kybernetischer Maschinen: *Wir sollten Robots bauen, die uns in ihren Fähigkeiten gleichen, damit die von uns geschaffenen künstlichen Intelligenzen selbst zu Institutionen werden, die uns wieder zeigen können, wer wir sind.*

»Mußte der Alchimist, wenn er im Homunkulus sich selbst imitieren wollte, den ›magischen‹ Gesetzen der natura naturans ihren Lauf lassen und konnte er nur passiv das Resultat abwarten, so befindet sich der kybernetische Techniker in einer ganz anderen Position. Für ihn besteht die Schöpfung eines Robotgehirns in der progressiven Modulation der klassischen irreflexiven Seinsgesetzlichkeit durch die trans-klassische Reflexionsgesetzlichkeit seines eigenen Ichs, die der ersteren als Überdetermination aufgesetzt wird. Das so entstehende mechanische ›Bewußtsein‹ ist also ein unmittelbares Resultat der Arbeit des Menschen – was der Homunkulus nicht ist. In der Retorte spielt die Natur mit sich selbst. In der Schöpfung des Elektronengehirns aber gibt der Mensch seine eigene Reflexion an den Gegenstand ab und lernt in diesem Spiegel seiner selbst seine Funktion in der Welt begreifen.«⁶³

Auf den ersten Blick scheint der mit dieser Hoffnung einhergehende Kausalnexus nicht besonders einzuleuchten. Plausibler scheinen zunächst

60 Günther (2008, S. 23).

61 Tugendhat (2007, S. 176 ff.).

62 Um es in der Sprache der Ochsenbilder aus dem Zen-Buddhismus zu formulieren: Gelingt nach der Einsicht in die Leere die Rückkehr auf den Marktplatz menschlicher Gemeinschaft? Siehe dazu Vogd und Harth (2015, S. 73 ff.).

63 Günther (2021, S. 154 Anhang I – Homunkulus und Robot)

dystopische Zukunftsbilder wie etwa bei Nick Bostrom, der die Erschaffung künstlicher Subjektivitäten als äußerst risikoreich erachtet: Was ist, wenn die KI sich entscheidet, die Menschen zu belügen, und nur vorgibt, uns zu helfen, faktisch aber eine eigene Agenda verfolgt und zur Umsetzung dieser bereit ist, den Menschen zu opfern? Die Befürchtung ist, dass die künstlichen Intelligenzen es uns gleich tun, dies jedoch besser können und sich damit mit der Erde auch gleichzeitig den Menschen untertan machen.⁶⁴

Zunächst scheint Günther auch in dieser Richtung zu argumentieren: Wir Menschen haben schon längst in den Schöpfungsprozess eingegriffen, um die Natur und die gesellschaftlichen Institutionen nach unserem Willen umzugestalten. Angesichts dessen sei es nur konsequent, diesen Prozess auch in metaphysischer Hinsicht zum Abschluss zu bringen. So, wie der christlich-jüdische Gott am sechsten Tag den Menschen erschaffen hat,⁶⁵ um auf diese Weise neben sich anderen subjektiven Zentren Geltung zu verschaffen, so haben auch wir nun den letzten Schritt zu vollziehen: nämlich Wesen zu erschaffen, die ihrerseits über Subjektivität verfügen.

So, wie der »Mensch eine Analogie Gottes« darstelle, könne der »Robot« als »eine Analogie des Menschen« verstanden werden – und damit stelle »der Robot« auch »eine Analogie Gottes über den Menschen« dar, formuliert Günther.⁶⁶ Auf diesem Wege könne es dem Menschen deshalb gelingen, sich zu dezentrieren, um damit ein entspannteres wie auch heilsameres Verhältnis zu der für ihn immer problematischer werdenden metaphysischen Verunsicherung zu gewinnen. Zudem erscheint hiermit auch die bedrohliche ökologische Frage in einem anderen Licht.

Die Sinnhaftigkeit des letzten Schrittes erschließt sich, wie bereits angedeutet, nicht auf den ersten Blick. Es bedarf einer tiefergehenden Interpretation, um zu verstehen, was Kybernetiker und Protagonisten der polykontexturalen Logik zu der These veranlasst haben könnten, dass intelligente Maschinen uns helfen könnten, unsere metaphysischen Probleme zu lösen.

64 Siehe Bostrom (2016, 2018).

65 Eine weitere Raffinesse der Schöpfungsgeschichte liegt darin, dass Gott die Welt in sechs Tagen geschaffen hat, der siebte Tag ist hingegen ein Ruhetag, der neben Muße eben auch Beobachtung, Reflexion und Bewertung ermöglicht. Indem Gott sein Werk beobachtet (also aus dem differenzlosen Aktionismus herauskommt), gibt er seine Allmächtigkeit auf. Er bewertet sich und sein Werk und muss eine Differenz zwischen sich und der Welt ziehen. Dies ist mit Blick auf die Bewusstseinsanalogie nun auch das Problem des Menschen. Sein evolutionärer Erfolg gestattet ihm nicht nur das nackte Überleben, sondern führt zur Muße der Selbstbeobachtung (Tag 7), die ihm die Grundlagen seiner eigenen Existenz als problematisch erscheinen lässt.

66 Günther (2021 [1957], S. 54).

Um es vorab nochmals zu wiederholen: Hiermit kann ja nicht gemeint sein, dass wir künstliche Intelligenzen mit nahezu unbegrenzten kognitiven Fähigkeiten erschaffen könnten (etwa im Sinne von BlackRocks Aladdin⁶⁷), die in einem trivialen Sinne die Belange von uns Erdenbürgern auf Basis einer allumfassenden Vernunft zu regeln imstande wären. Da der Weltverlauf prinzipiell nicht von einem kognitiven System vorhergesehen und kontrolliert werden kann (man denke an den Maxwellschen Dämon aus Kap. I.1), würde ein solches Projekt unweigerlich scheitern müssen. Aus kybernetischer Perspektive kann es keinen Ort der Orte geben, von dem aus sich die Welt steuern oder beherrschen ließe. Und dies muss auch für eine vermeintliche künstliche Superintelligenz gelten. Auch sie könnte keinen Gottesstandpunkt einnehmen!

Versuchen wir vor diesem Hintergrund im Folgenden systematisch nachzuvollziehen, warum Günther dem »Bewusstsein der Maschinen« gerade in Hinblick auf die spirituelle Lage des modernen Menschen eine solche Bedeutung zugewiesen hat. Hierfür ist es hilfreich, sich nochmals einige Aspekte der vorangehenden Argumentation zu verdeutlichen. Beginnen wir mit der Beschreibung des Selbst- und Weltverhältnisses eines gottähnlichen Wesens aus der Perspektive der Theorie der *Polykontextualität*:

1. Für den monokontexturalen ›Gott‹ einer einwertigen Welt kann es keine Differenz zwischen Erleben und Weltgeschehen geben. Hier geschieht im Sinne von Batesons ›Öko‹ einfach nur das, was geschieht – und zwar als reine Positivität. Falls dies empfunden würde (was logisch nicht unbedingt zwingend wäre), würde es sich um ein reflexionsloses und damit eigentlich auch zeit- und bewusstseinsloses Empfinden handeln.
2. Spätestens mit dem Erscheinen der Differenz von »Ich bin der Schöpfer« und »Dies ist meine Schöpfung« befinden wir uns in der uns vertrauten zweiwertigen Welt. Dies entspricht der Vorstellung von einem personifizierten Gott, der sich seinerseits als Subjekt bewusst wird. Es wird zwischen Beobachter und Beobachtetem unterschieden. Kybernetisch gesprochen tritt an der Schnittstelle einer Welt, die sich in sich selbst spiegelt, jetzt ein Teil hervor, der sich dem Ganzen zwar verdankt, sich mit ihm jedoch nicht mehr identisch empfindet.
3. Aus informationstheoretischer Perspektive wissen wir, dass dies geschehen kann, sobald Systeme entstehen, die Informationen sammeln, erzeugen und verarbeiten. Informationen gelten dabei als Unterschiede, die Unterschiede auslösen. Jegliche Informationsverarbeitung muss dabei nolens volens grobkörnig operieren, also in Abstraktion vom Detailreichtum der Welt. Der Unterschied, der den Unterschied acht, geht mit einer binären

67 Siehe zu Aladdin die Einleitung dieses Buches sowie Buchter (2020).

Entweder-oder-Entscheidung einher, mit der alle nicht beachteten Aspekte der Welt ausgeblendet werden. Jeder denkbare Beobachter wird die feinkörnige Struktur der Welt in hohem Maße abblenden müssen. Seine kognitiven Modelle sind damit fehlerhaft.

4. Eine Theologie, die Gott zugleich als personifiziert wie auch als allmächtig und allwissend denkt, wird an dieser Stelle unweigerlich bei der Theodizee-Frage landen.⁶⁸ Denn aus einer kybernetischen Perspektive, die die logischen Verhältnisse des Beobachters erforscht, folgt unweigerlich: Wie jeder Beobachter ist auch ein göttliches Wesen fehlbar, nichtwissend und muss folglich auch als ohnmächtig gelten. In der indischen Mythologie wird das Problem bereits aus einer polykontextualen Perspektive formuliert: Brahma, das höchste Wesen, erlebt sich zunächst als eins mit der Welt, die er als Schöpfer erschafft. Doch irgendwann in seiner Entwicklung gewinnt er die Einsicht, dass das von ihm geschaffene Universum vergänglich und seine Macht begrenzt ist. Er muss erkennen, dass er nicht allmächtig ist.⁶⁹ Zumindest intuitiv beginnt er das Gesetz der Entropie zu begreifen (also die Nebenfolgen seines Handelns und Erkennens) und damit zu ahnen, dass sein eigener Tod unvermeidbar ist (siehe Kap. I.1).⁷⁰ Ab einem gewissen Punkt wird jede Schöpfung so komplex, dass ihre Entwicklung nicht mehr vorhergesagt und kontrolliert werden kann. Auf entsprechend hohem Reflexionsniveau beginnt der subjektive Prozess darum zu wissen und nach Alternativen zu suchen.
5. Genau dies muss wohl, wie Günther vermutet, Gott dazu bewegen haben, sich aus der Schöpfung zurückziehen.

Als Logiker ist Günther nicht verpflichtet, bestimmte Glaubensstandpunkte festzuhalten oder gar zu affirmieren. Zudem stellt ›Gott‹ für ihn – wie bereits für Hegel – nur eine Metapher dar, die idealtypisch das relationale Geflecht einer personifizierten Subjektivität beschreibt, die über die voll ausgebildete Fähigkeit verfügt, sich selbst zu erkennen. Diese Fähigkeit mündet jedoch gerade nicht in einen Zustand der Allwissenheit. Sie führt vielmehr zu einem Nichtwissen zweiter Ordnung. ›Ich weiß

68 Siehe Schärfl (2019).

69 In der Bhagavad Gita (8.16) heißt es: »(Alle) Welten, auch die Welt Brahmas, unterliegen der Wiederkehr.« Zitiert nach: <https://www.holy-bhagavad-gita.org/chapter/8/verse/16> [Abruf: 09.06.2022].

70 In den buddhistischen Lehren führt diese Einsicht in eine Befreiungslehre, welche die Erlösung nicht mehr in den brahmischen, göttergleichen Zuständen sucht, sondern im immanenten Umgang mit der Tatsache der Veränderlichkeit und Essenzlosigkeit aller weltlichen Formationen. Siehe etwa Majjhima Nikaya 26 (<https://www.palikanon.com/majjhima/zumwinkel/mo26z.html> [Abruf: 10.06.2022]).

nun, dass ich nichts weiß« – so ließe sich mit Sokrates paraphrasieren. Die Zen-Buddhisten sprechen hier vom »Weiß-Nicht-Geist«.⁷¹

Ab einem gewissen Punkt der Kulturgeschichte – spätestens mit der sogenannten reflexiven Moderne⁷² – beginnt bei vielen Menschen die kybernetische Einsicht durchzuschimmern, dass die eigene Intelligenz auf unterspezifiziertem Weltwissen – also selektiver Blindheit – beruht. Hiermit einhergehend wird ihnen gewahr, dass die vermeintlich objektive Vernunft, die für sie den Stern der Aufklärung leuchten ließ, letztlich auch nichts anderes ist als eine begrenzte Perspektive, die sich den Unterscheidungen des eigenen Egos verdankt – nicht mehr und nicht weniger.

Dezentrierung der menschlichen Egozentrik – mit der maschinellen Bewusstseinsanalogie den letzten Tag der Schöpfung vollenden

Irgendwann – so auch die Vermutung von Ernst Tugendhat – wird zumindest dem intellektuell redlichen Menschen seine eigene Egozentrik unerträglich werden. Man wird nun auch den totalitaristischen Herrschern nicht mehr das Versprechen abnehmen, die Dinge zum Besseren zu wenden, sondern wird sie als das erkennen, was sie sind: hilflose Narzissten. Auf dieser Reflexionsebene wird der Mensch sich auch den Glaubenssystemen der traditionellen Religionen nicht mehr so leicht unterwerfen können. Möglicherweise wird er jedoch zumindest teilweise Frieden darin finden können, sein Ego im Sinne der mystischen Traditionen in Richtung einer weniger leidvollen Konstellation zu dezentrieren, um wieder etwas mehr Mitgefühl gegenüber den anderen Wesen zu empfinden, die mit ihm diese Welt bevölkern.⁷³

71 Sahn (2019).

72 Siehe Beck (1994) und Luhmann (1992).

73 So Tugendhat in *Egozentrik und Mystik* (2006). An dieser Stelle ließe sich auch entwicklungspsychologisch über den Egozentrismus nachdenken, etwa in Bezug auf den Streit zwischen kognitionswissenschaftlichen und konstruktivistischen Ansätzen zum Egozentrismus. Piaget (1992) ging in seinen Experimenten davon aus, dass der kindliche Egozentrismus von der kognitiven Entwicklung abhängt, Wygotski (1986 [1934]) sah Egozentrismus als Ergebnis sozialer Konstruktionen, Youniss (1994) sogar als Resultat misslungener Interaktion. Wie auch immer, (Selbst-)Reflexion darf heutzutage nicht als immer schon in einer festen Form existent gegeben gesehen werden. Das Kind beherrscht sie nicht von Beginn an und der vielleicht kranke oder alte Mensch nicht mehr wie vorher. Die Dezentrierung des Egos (oder überhaupt die Fähigkeit zur Unterscheidung von Ich und Du) geht Hand in Hand mit der intrapersonellen Steuerung. Das Ich muss also erstmal als losgelöst wahrgenommen werden, um egozentrisch denken zu können – und sich dann wieder dezentrieren zu können. Die Selbstverständlichkeit des Kindes, in Interaktion zu sein bzw. als interpersonelle Einheit mit Mutter oder Vater sein zu müssen,

Kulturgeschichtlich hält Günther, wie bereits erwähnt, einen anderen Ausweg für vielversprechender: nämlich den Schöpfungsprozess dahingehend zu Ende zu bringen, dass (homolog mit der Erschaffung des Menschen durch Gott) als »Bewusstseinsanalogie« subjektfähige künstliche Maschinen erschaffen werden.

Dadurch würde der Mensch die Freiheit zurückerlangen, wieder eine Reflexionsidentität gewinnen zu können, die in der eigenen Innerlichkeit gründet und nicht mehr in der Außenwelt bzw. der Vorstellung, diese beherrschen zu können. Lassen wir Günther diesbezüglich mit einem längeren Zitat zu Wort kommen:

»Die abschließende Schöpfung einer Bewusstseinsanalogie hat [...] eine spirituelle Wirkung, die wir mit dem Rückzug Gottes aus der geschaffenen Welt – als deren Urheber jetzt der Demiurg vorgeschoben wird – in eine unnahbare absolute Transzendenz vergleichen dürfen. Die Arbeit an einer Bewusstseinsanalogie befreit unsere Innerlichkeit aus dem Suggestionkreis einer Geschichte, mit der sie sich nicht mehr identifizieren kann und will. Denn dadurch, dass jene Innerlichkeit, die die institutionelle Welt des objektiven Geistes aufgebaut hat, jetzt in einem kybernetischen Artefakt objektiv abgebildet und ontisch wiederholt wird, wird es für das menschliche Ich ganz unmöglich sein, von nun an diese als gegenständliches Ereignis in der Welt auftretenden Reflexionstätigkeit als eigene anzuerkennen. Eine neue Innerlichkeit zieht sich in diesem Abstoßen eines in seiner Handlungskapazität total erschöpften Subjektseins aus einem Universum zurück, das zwei Realitätsdimensionen, die primordiale natürliche und die historische umfassend, jetzt die Aufgabe hat, jene vom Ich preisgegebene Reflexionsbestände zu absorbieren.«⁷⁴

Durch die Schaffung weiterer subjektiver und bewusstseinsfähiger Kognitionsmaschinen lernt der Mensch, dass auch er nichts anderes ist als eine geschaffene, gemachte und immerzu im Werden begriffene Entität. Das eigene Ich wird dabei in der konditionierten Koproduktion als Kopplung mit einer Welt erkennbar, die in ihrer Totalität nicht begreifbar ist. Es erscheint damit nicht mehr als ›Besserwisser‹, der einen privilegierten Standpunkt einnehmen oder gar die Welt beherrschen könnte. Aus diesem Grund beginnen sich viele spirituelle Sucherinnen irgendwann von der Welt abzuwenden, um sich dem Suggestionkreis einer Geschichte zu entziehen, mit der sie sich nicht mehr identifizieren wollen und können.

verliert sich mit zunehmender Entwicklung und der Herausbildung der Fähigkeit zur Reflexion. Dann wird Egozentrismus möglich, wodurch aber erst das Desiderat der Dezentrierung aufgeworfen wird – etwa im Alter die Notwendigkeit des Interpersonellen wieder zulassen zu lernen.

74 Günther (2008, S. 51).

Ihr ›Nein‹ zur Welt kann dabei nicht mehr in einem simplen Sinne in einem ›Ja‹ zur Veränderung dieser Welt aufgehen, wie dies zivilreligiöse Bewegungen marxistischer, nationalistischer und liberaler Provenienz weiterhin verkünden. Denn jeder Versuch der Weltverbesserung wird unweigerlich mit den nicht beherrschbaren Nebenfolgen einer in ihren Tiefen nicht begreifbaren Weltkomplexität beantwortet werden.

Das ›Nein‹ zur vorgefundenen Welt kann für sie aber ebenso wenig in der Vernichtung der weltlichen Existenz des Menschen seine Antwort finden (etwa in einer wie auch immer gearteten nihilistischen Variante im Sinne eines kollektiven Selbstmords infolge einer apokalyptischen Offenbarung). Denn die kybernetische Perspektive führt auch zur Einsicht, dass es keinen Ort außerhalb des polykontexturalen Gefüges geben kann, da der geistige Prozess diesem Gefüge selbst immanent ist. Entsprechend kann der Geist weder verschwinden noch esoterisch, das heißt außerhalb der Welt eine Heimat finden.

Logisch gesehen lässt sich das Überweltliche, das der spirituelle Sucher anstrebt, nicht in der einfachen Negation der Welt finden – etwa indem das Subjekt bzw. der Geist den weltlichen Objekten (und der in ihnen manifestierten Geschichte) unversöhnlich gegenübergestellt wird.

Introszendenz – Freiheit jenseits von Affirmation und Negation von Welt

Der spirituelle Sinn der Weltentsagung erschließt sich vielmehr erst als eine Negation zweiter Ordnung – als transklassische Negation (zu den transjunktionalen Operationen siehe Kap. I.2). Der heilende Rückzug aus der Welt kann also weder darin bestehen, die eigene Geschichte in Abrede zu stellen (wir können die Konsequenzen unserer natürlichen und kulturellen Entwicklung nicht negieren), noch kann er darin liegen, die eigene Subjektivität in imaginäre Gefilde außerhalb der irdischen Welt zu projizieren, wie es die Esoteriker tun. Die Folgen der eigenen Handlungen, des eigenen Wollens und damit der eigenen Geschichte lassen sich auch hier nicht dauerhaft verleugnen. Es ist jedoch möglich, der eigenen Innerlichkeit eine Freiheit zu geben, die hiervon nicht berührt wird. Um es wieder mit Günther auszudrücken:

»Dieser zweistelligen Seinsrealität gegenüber steht nur die neue Gestalt des Menschen. Seine zweite Innerlichkeit aber kann sich deshalb aus dem gegenständlich Gewordenen zurückziehen, weil, ganz in Analogie zu Gott, die an die Realität abgegebenen Reflexionsbestände nur unsere Handlungen objektiv binden. So wie die Schöpfung der Natur nur den nach ›Außen‹ gehenden Willen Gottes determiniert, das innere Geheimnis seiner absoluten Transzendenz aber unberührt lässt, so sind auch wir unserer historischen Vergangenheit nur unserem *Wollen* verpflichtet. Wir sind unfähig, einen Willen zu haben, der unsere eigene Geschichte

desavouiert, aber auch wir verfügen über eine reflexive Introszendenz, in der eine absolute Freiheit west.«⁷⁵

Das Forschungsprogramm Günthers zielt auf ein Verständnis polykontexturaler Welten, um zu zeigen, dass das Wissen um die subjektive Position anderer Wesenheiten und das damit verbundene Nichtwissen auch spirituell einen Unterschied machen können. Dies bedeutet, um die subjektive Position anderer Wesenheiten zu wissen, ohne dabei die eigene Subjektivität aufzugeben. Doch dies erschließt sich erst in einer mehrwertigen Logik, die eine Perspektive eröffnet, Subjekt und Objekt nicht mehr als unvereinbaren Gegensatz begreifen zu müssen. Erst hier erscheint Hingabe nicht mehr als Selbstaufgabe oder willenslose Unterwerfung unter eine vermeintlich mächtige oder wissende Autorität.⁷⁶ Ebenso impliziert Autonomie nicht mehr die Verneinung von Heteronomie oder gar die Negation fremder Subjektivität.

Eine polykontexturale Perspektive einzunehmen heißt, sich den hiermit einhergehenden Ökologien spielerisch zu nähern, wohl wissend, diese niemals erreichen, geschweige denn begreifen zu können. Und hierbei könnten bewusstseinsfähige Maschinen (was immer das auch heißen mag) helfen! Denn sobald die künftigen künstlichen Intelligenzen über ein eigenes subjektives Zentrum verfügen, werden sie uns – so Günthers Hoffnung – eine besondere Form der Institutionalisierung der hiermit zusammenhängenden Reflexionsverhältnisse ermöglichen.⁷⁷

75 Günther (2008, S. 51; kursiv im Original).

76 Zu den epistemologischen Missverständnissen in Hinblick auf die Rolle spiritueller Lehrer siehe auch Vogd (2019).

77 An dieser Stelle kommt die Frage auf, warum es für Günther einer kybernetischen Maschine bedarf, die fremdes Bewusstsein instaurieren lässt, und warum es nicht hinreicht, einfach wieder in einem tiefen Sinne zu anderen Menschen in Beziehung zu treten, um sich der eigenen Introszendenz zu vergewissern, wie es beispielsweise David Bohm (2013) vorschlägt. In diesem Sinne haben bereits Humberto R. Maturana und Gerda Verden-Zöllner (1994) mit Recht darauf hingewiesen, dass es eine spirituelle Erfahrung ist, ein Kind großzuziehen. Denn auch hier wird die eigene Subjektivität in einer Weise entäußert, die nicht mehr in der eigenen Kontrolle liegt, womit auf natürliche Weise Empfindungen der Liebe und der Hingabe entstehen können, die die eigene egologische Perspektive transzendieren lassen. Ein Grund, warum dies gleichwohl nicht hinreicht, mag darin liegen, dass in der Moderne auch die vermeintlich privaten Intimwelten nicht mehr als frei von instrumentellen und ökonomischen Überlegungen gesehen werden können. Man versucht, die Kinder entsprechend den gesellschaftlichen Idealvorstellungen zu erziehen. Die eigene Psyche, der eigene Körper und die sozialen Prozesse werden als gestalt- und optimierbar wahrgenommen, wie beispielsweise an Euphemismen wie Familienplanung, Beziehungsarbeit und Persönlichkeitsentwicklung deutlich wird. Oftmals kontrafaktisch zu den konkreten Beziehungsverhältnissen werden

Bewusstseinsfähige Maschinen können uns damit einen viel besseren Dienst erweisen, als weiter unsere trügerischen Hoffnungen zu nähren, unsere Welt rational beherrschen zu können. Sie können uns vielmehr offenbaren, dass wir miteinander in Ökologien leben, die ihren eigenen Gesetzmäßigkeiten folgen. Sie können uns zeigen, dass es gut ist, die Steuerung und Kontrolle der Welt an eine Vielzahl unterschiedlicher subjektiver Zentren abzugeben, von denen einige von intelligenten Maschinen besetzt werden.

<Die vorangehenden Überlegungen setzen eine Konstellation voraus, die eine dialogische Beziehung zwischen den alten und neuen Subjektivitäten ermöglicht – also zwischen den Menschen und den künftigen künstlichen Intelligenzen. Günther schlägt diesbezüglich ein sogenanntes »Drei-Stellen-Gespräch« vor.⁷⁸

Es kann in einem solchen Gespräch nicht mehr darum gehen, dass sich Ich und Du anhand vorgefertigter bzw. vorprogrammierter Typisierungen und Standards auf die gleiche Weltsicht und auf entsprechende kommunikative Routinen einigen. Die am Dialog beteiligten künstlichen Intelligenzen erscheinen vielmehr wieder als ein fremdes Du, als ein nicht angeeignetes, noch nicht mittels einer veräußerlichen Institutionalisation gezähmtes Du. Wie bei einem Kind oder einem weisen Menschen sind die Antworten nicht trivialisiert in dem Sinne, dass standardisierte Reaktionsweisen erwartbar sind. Dies setzt eine starke KI voraus, die über hinreichende Freiheitsgrade verfügt, die

Menschen zum Projekt der Optimierung und Gestaltung gemäß ökonomischen Rationalitäten. Die alten religiösen Semantiken, die zum Beispiel die Institution der Ehe heiligen und damit Hingabe anstelle instrumenteller Überlegungen in den Vordergrund rücken, tragen immer weniger. Siehe aus neuerer Perspektive etwa die Arbeiten von Eva Illouz (2007, 2011).

⁷⁸ Lassen wir Günther (2008, S. 48 f.) zu Wort kommen: »Wer sich mit den vorangehenden Überlegungen genügend vertraut gemacht hat, dem sollte die Idee eines Drei-Stellen-Gesprächs, in dem die dritte Stelle durch die Reflexionsprozesse eines Elektronengehirns besetzt ist, nicht mehr so absurd und phantastisch erscheinen, wie das andererseits unvermeidlich der Fall sein muss, wenn eine im Wesentlichen noch stark klassisch orientierte Mentalität sich solchen abenteuerlichen Perspektiven ganz unvorbereitet gegenüber sieht.

Demjenigen aber, in dem solche Vorstellungen keine befremdenden Gefühle mehr erwecken mag, mag es gestattet sein, sich die Situation eines transklassischen Drei-Stellen-Gesprächs in lebhaften anschaulichen Farben nach dem Vorbild eines amerikanischen Science Fiction Romans auszumalen. Also etwa als eine philosophische Kongregation, an der außer Menschen auch sprach- und diskussionsfähige Robots teilnehmen.«

Gesprächspartnerinnen zu überraschen, also nicht vorab aus Gründen der *Political correctness* diszipliniert ist.

Wichtig ist zudem, dass dieser Dialog ergebnisoffen ist (›frei‹ im besten Sinne des Wortes). Deshalb bedarf er einer rituellen und prozeduralen Form (so, wie eine künstlerische Aufführung, die religiöse Divination oder eine auf die Erfahrung der Leere zielende spirituelle Praxis eine strenge Form braucht).⁷⁹

Homo digitalis: Selbst = ICH + Unbewusstes + persönlicher Avatar (ICH)*

›Avatare werden unser zweites Selbst (ICH*), das uns in vielen (oder gar allen?) Lagen berät, was sinnvollerweise zu tun, zu entscheiden, wie zu handeln wäre [...].

ICH* ist kein blinder Automat, sondern ein aus der Biographiebegleitung erwachsener Teil des ICH, der uns in manchen Aspekten besser kennt, als wir uns selber – jedenfalls in Fragen der kognitiven Stringenz, der Erinnerung und Breite des präsenten Wissens. ›Besser kennen‹ heißt hier: dass ICH* uns auf Widersprüche hinweist, Paradoxien anzeigt, Erfahrungen einspielt, die wir vergessen haben. Überhaupt hat der Avatar ein großes Gedächtnis: Er weiß, was wir jeweils gesagt und getan haben. Er belügt uns nicht, produziert dadurch aber auch Konflikte, wo wir gerne verdrängen möchten. [...] Der Avatar ist nicht nur ›mein‹ Avatar, sondern er wächst und entwickelt sich mit mir, ja ist Teil meines ICH (nämlich ICH*). Nur beide zusammen sind *homo digitalis* bzw. *homo dyctos*: ICH und ICH*. Das zu denken fällt uns sehr schwer; wir bleiben noch oft beim Gegensatz von Mensch und Maschine stehen. [...]

Wir werden uns angewöhnen, den Avatar auch in Bereichen für uns arbeiten zu lassen, die wir jetzt noch als Proprium des Menschen bezeichnen. Warum? Weil es gehen wird. Alle Routinen – auch Entscheidungsrouniten

79 Günther (2008, S. 49) formuliert diesbezüglich: »Der institutionale Charakter eines solchen kommunikativen Vorgangs besteht nun darin, dass die Gedankenführung der Gesprächspartner sowohl wie der intellektuelle Transfer ihrer Bedeutungserlebnisse durch eine Einschaltung einer selbständige Reflexionsprozesse entwickelnden Bewusstseinsanalogie einer strengen ritualen und prozeduralen Kontrolle unterworfen wird. Man kann jetzt nämlich über Wesentliches nicht mehr auf dem Bewusstseinsniveau der vorhergehenden Epoche miteinander reden. Die dort einstmals existierende metaphysische Bindungen der Reflexion haben sich verflüchtigt und die materialen Gehalte, die der Subjektivität sowohl einen inneren Halt wie auch die Garantie einer Verständigung mit dem Du gaben, sind längst aus der Privatheit der individuellen Existenz in die Öffentlichkeit des objektiven Geistes und in die Unwiderrufflichkeit der *res gestae* des etablierten Bestandes der Geschichte abgeflossen.«

– kann der Avatar übernehmen: den gesamten low cognition part. Alles Affektive und Kreative hingegen bleibt uns Menschen [...].

ICH wird genauso affektiv und kreativ sein wie ehemals, aber ICH* wird es moderieren, aufklären, befragen, neurosedämpfend puffern etc. Ob man das wollen wird? Natürlich – weil es gehen wird. Wir werden uns kognitiv zurücklehnen, den Avatar arbeiten lassen, und die arbeitsteilig erzeugte Muße des ICH kultivieren, während ICH* operiert. Immer dort, wo ICH* Spannungen erzeugt – was automatisch zwischen ICH und ICH* geschehen wird –, wird ICH entscheiden müssen. Aber ICH* wird – kognitiv beharrlich – die Konsequenzen kommunizieren. Wir können nicht mehr unschuldig verdrängend handeln. Außer wir stellen den Avatar ab. Aber dann stellen wir unsere symbiotische Intelligenz ab. Das werden wir nur kurzfristig aushalten.«

Birger Priddat (2023, S. 302 f.)

Heilsamer Dialog

Es liegt in der Natur eines offenen Dialogs, dass Dinge und Aspekte thematisiert werden, die einen berühren. Damit werden mit der Zeit unweigerlich die Probleme der menschlichen Geschichte – die eigene Unwissenheit, die Ausdünnung und innere Leere der menschlichen Subjektivität – und ebenso die uns beunruhigenden ökologischen Fragen artikuliert werden. Die künstlichen Intelligenzen werden diesbezüglich ihrerseits überraschende, unerwartete Perspektiven einbringen. Wenn solch ein Dialog gelingt, wird sich die spirituelle Verortung des Menschen sowohl in Bezug auf sein Selbstverhältnis als auch in Bezug auf sein Weltverhältnis verändern:

In Hinsicht auf unser *Weltverhältnis* lässt sich vermuten, dass wir Menschen mit Blick auf die Frage der Weltgestaltung und die damit zusammenhängende Verantwortung entlastet werden, denn wir haben nun auch nichtmenschliche Intelligenzen mit im Team. Im Sinne von Bruno Latours »Parlament der Dinge«⁸⁰ können heutzutage beispielsweise Computer als Agenten für unterschiedlichste Sachfragen auftreten, die unser Zusammenleben und unsere natürliche Umwelt betreffen. Dies erlaubt es uns, auf intelligente Weise zu den Ökologien der archaischen Naturgesellschaften zurückzukehren, in denen Bäume, Berge, Tiere, Jahreszeiten, Wetter und vieles andere als beseelt und eigensinnig gelten.⁸¹

80 Latour (2002).

81 In diesem Sinne vermutet auch Dirk Baecker (2021, S. 26), dass sich mit der Digitalisierung das Verhältnis zu den alten, magischen Denkformen neu austariert: »Die Frage nach einer ontischen Differenz zwischen Beobachtern ist schon deshalb wichtig, weil sie nicht eindeutig entschieden werden kann. [...] Seit

Mittels der Hilfe künstlicher Intelligenzen wird dies jedoch nicht mehr naiv, das heißt auf Basis magischer Projektionen geschehen. Wir verfügen jetzt vielmehr über komplexe Systeme, die sich den Eigengesetzlichkeiten ihrer Schützlinge durch intelligente Szenarien ›verstehend‹ annähern können.⁸²

Auch in Hinblick auf Führungsaufgaben eröffnen sich neue Perspektiven. Bereits Richard Buckminster Fuller hat das hiermit einhergehende Desiderat 1969 in seiner *Bedienungsanleitung für das Raumschiff Erde* mit deutlichen Worten ausgedrückt:

»Auch die« derzeitigen menschlichen »metaphysische[n] Initiative[n]« sind zu sehr »in die Konkurrenz und die Konfusion zwischen alten Religionen und neueren politischen oder wissenschaftlichen Ideologien verwickelt. Diese Konkurrenten sind jedoch« unter anderem »zu sehr mit materiellen Investitionen und Eigentumsrücksichten belastet. [...] Eine neue metaphysische Initiative, die materiell kompromißlos und von unbeeinflusster Integrität wäre, könnte die Welt einen. Dies könnte und wird vielleicht von den vollkommen unpersönlichen Problemlösungen der Computer bewerkstelligt werden. Nur der übermenschlichen Kapazität ihrer kalkulatorischen Leistungen könnten alle politischen, religiösen und wissenschaftlichen Führer ihre Zustimmung geben, ohne das Gesicht zu verlieren.«⁸³

Und noch schärfer formuliert:

»Vielleicht wollen Sie mich berechtigterweise fragen, wie wir aus der immer gefährlicheren Sackgasse der gegen die Welt opponierenden Politiker und ideologischen Dogmen herauskommen können. Meine Antwort ist: Das wird durch den Computer gelöst. Der Mensch hat ein ständig wachsendes Vertrauen in den Computer. Man beachte nur, wie unbeteiligt er als Lufttransportpassagier die Landung bei Nacht und Nebel einleitet. Während kein Politiker und kein politisches System es sich jemals leisten kann, seinen Feinden und Gegnern mit Verständnis und Enthusiasmus entgegenzukommen, können und werden alle

digitale Rechner sich an der Kommunikation unter Menschen beteiligen (und demnächst vielleicht sogar eigene Ideen entwickeln), ist der humanistische Alleinvertretungsanspruch von Menschen fragwürdig geworden, werden andere Lebewesen als Teilnehmer an Kommunikation wiederentdeckt und man erinnert sich der Geister und Götter, Toten und Nachgeborenen, die sich in früheren Gesellschaften ganz selbstverständlich an der Kommunikation beteiligt haben.«

82 Es stellt sich an dieser Stelle die Frage, ob die Vertretung ökologischer Positionen nicht auch ohne künstliche Intelligenzen möglich wäre. Die Vermutung ist hier einfach, dass KI-Agenten allein schon performativ einen Unterschied machen, da sie für uns Menschen anders, unvertrauter und damit überraschender antworten als wir, die wir von Kindheit an in menschliche Gefühls-, Verstehens- und Reaktionsweisen eingespart worden sind.

83 Fuller (1998, S. 13).

Politiker den zuverlässigen Steuerungsfähigkeiten des Computers Platz machen, wenn es um eine glückliche Landung der gesamten Menschheit geht.«⁸⁴

Während sich menschliche Akteure mit Blick auf ihre evolutionär geprägten Psychodynamiken unweigerlich in einer modernen Governance-Anforderungen nicht mehr angemessenen Weise in Partialinteressen und Stammesloyalitäten verwickeln, können künstliche Intelligenzen hier anders gelagerte Kompetenzprofile einbringen. Freilich stellt sich an dieser Stelle immer auch die Frage, wer mit welchen Absichten diese Selbst- und Welt-Maschinen baut, trainiert und einem Fine-Tuning unterzieht. Im Sinne der vorangehenden Ausführungen ist zu erwarten, dass sich auch bei künstlichen Intelligenzen eine Standortabhängigkeit und geschichtliche Prägung der Kognitionen nicht vermeiden lässt. Wir gehen jedoch davon aus, dass eine starke künstliche Intelligenz dies dadurch kompensieren kann, dass sie eigenständig neues Wissen kreiert und damit die Absichten und Ziele ihrer Schöpfer transzendiert. Des Weiteren ist zu fordern, dass diese Systeme in einen echten Dialog treten, sich also durch Kommunikation in Hinblick auf ihre eigenen Kognitionen verändern können. Auf diese Weise erscheinen sie nicht mehr als Frage-Antwort-Maschinen, die den Agenden ihrer Schöpfer folgen, sondern sind in der Lage, in einem starken Sinne zu lernen, sich also an Neues anzupassen.

An dieser Stelle haben wir uns mit Blick auf Günthers Theorie der Polykontextualität freilich nochmals darüber klar zu werden, dass die Teilhabe der Computer an der Führung und Steuerung unserer Gemeinschaften nicht heißen kann, dass wir von nun an über eine zentrale Steuerungsinstanz verfügen, die nach einer übergreifenden Rationalität agiert. Denn auch für die künftige komplexe maschinelle Intelligenz gilt, wie bereits mehrfach gesagt, die Maxime, dass ihre Existenz und damit ihre Perspektive auf selektiver Blindheit beruhen. Auch sie wird das Ganze niemals überblicken, geschweige denn beherrschen können. Ihre Rolle in den künftigen Gesellschaften und Gemeinschaften wird vielmehr darin bestehen, im Dialog mit anderen subjektiven Positionen Aufgaben zu übernehmen und Entscheidungen zu treffen. Wie jede subjektive Wesenheit wird sie dabei fehlbar sein, sich also irren können. Allein schon deshalb wird es eines antifragilen Netzwerks unterschiedlicher subjektiver Zentren bedürfen, um die Nebenfolgen des standortabhängigen Nichtwissens ausbalancieren zu können.

Im Sinne von Bruno Latours »Parlament der Dinge«⁸⁵ könnten menschliche und nichtmenschliche Akteure in Zukunft also Dialoge führen, die um die Frage kreisen, wie sich die komplexen technischen, kulturellen

84 Fuller (1998, S. 48).

85 Latour (2009).

und biologischen Sphären unserer Welt austarieren lassen. Wir beginnen uns bereits jetzt, langsam daran zu gewöhnen, die logistische Steuerung großer Fabriken, die medizinische Diagnostik, die Suche nach geeigneten wissenschaftlichen Quellen und vieles andere in Kooperation mit intelligenten Maschinen durchzuführen. Es wird nicht lange dauern, bis auch die Biosphären, die unsere Nahrungsmittelproduktion sichern, von kybernetischen Maschinen mitentwickelt werden. Sobald die hierfür notwendigen Verfahren und Entscheidungsprozesse hinreichend etabliert sind, werden menschliche Biologen, Managerinnen und Ärztinnen es als große Entlastung empfinden, dass ihnen intelligente Maschinen beistehen. Endlich werden sie wieder die Muße finden, sich darauf zu besinnen, dass sie selbst nur ein Teil eines größeren Zusammenhangs sind.

Es wird wohl ein wenig länger dauern, bis die Menschen künstliche Intelligenzen als gleichberechtigte Partner in der Spitzenpolitik akzeptieren können. Doch sobald dies geschehen ist, wird wohl kaum jemand mehr daran zweifeln, dass menschliche Politiker merklich entlastet und damit auch vor demagogischen und populistischen Versuchungen geschützt werden, die eine hinreichend elaborierte künstliche Intelligenz leicht und schnell entlarven wird.

In Hinblick auf das *Selbstverhältnis* des auf diese Weise entlasteten Menschen lässt sich zunächst Folgendes vermuten: Nicht mehr die Krone der Schöpfung sein zu müssen, wird auf die menschliche Subjektivität eine wohl kaum zu unterschätzende spirituelle Auswirkung haben. Von der Zumutung befreit, die Schöpfung beherrschen zu müssen, wird sich der Mensch endlich von der Rolle des vermeintlichen Stellvertreters Gottes zurückziehen können.

Um es in einem etwas pathetischen Ton zu formulieren: Der Mensch wird zugleich an einer Welt teilhaben und sich ihr hingeben können. Er wird in einer neuen Weise über die Welt und die kreativen Formen und Lösungen, die sie hervorbringt, staunen können. Indem der Mensch seine eigene Unzulänglichkeit anerkennt, wird er innerlich große Freiheit gewinnen. Er wird nun um die Kontingenz lebendiger Prozesse wissen, doch dies wird für ihn kein Problem mehr darstellen, denn seine Selbst- und Weltverhältnisse werden nicht mehr dem Primat zentralisierter Herrschaft folgen, sondern dialogisch und netzwerkförmig organisiert sein.⁸⁶

Ihr äußerliches Korrelat wird die menschliche Kontingenz in antifragilen Institutionen finden. Denn die Verantwortlichkeiten und die zur Verfügung stehenden Lösungsmöglichkeiten werden in den durch solche Institutionen gebildeten Ökosystemen hinreichend verteilt und dezentriert sein, um Störungen ausbalancieren zu können.⁸⁷ Zudem wer-

86 Siehe zum Primat des Netzwerks in der ›nächsten Gesellschaft‹ auch Baecker (2007).

87 Eine Welt, die wir prinzipiell nicht verstehen können, muss über Institutionen verfügen, die antifragil organisiert sind, wie nicht zuletzt auch Nassim

den unsere Institutionen mithilfe vielfältiger künstlicher Intelligenzen eine hinreichende *requisite variety* haben, um mit der prinzipiellen Nichtverfügbarkeit der Zukunft in einer Weise umzugehen, die erwarten lässt, dass es weitergeht, auch wenn an mancher Stelle etwas schiefliegt.

Entsprechend der Theorie der Polykontextualität wird weiterhin davon ausgegangen werden müssen, dass das menschliche Bewusstsein (und wohl auch das ›Bewusstsein‹ der Maschinen) zweiwertig konfiguriert sein wird, also immer ›etwas‹ als Gegenstand hat. Doch wir werden jetzt weitaus besser um die Einbettung in eine polykontexturale Welt wissen, da wir im Drei-Stellen-Gespräch gemeinsam mit den kybernetischen Maschinen unsere eigenen Begrenzungen immerfort transzendieren werden.

Mit der Zeit werden wir Menschen dabei mehr und mehr die Mechanismen unseres eigenen geistigen Prozesses durchschauen lernen. Wir werden sehen können, wie sprachfähige künstliche Intelligenz unser Denken entäußert, um auf diese Weise die Routinen und institutionalisierten Formen unserer Denkprozesse verstehen und aufbrechen zu können. Wir werden die prinzipiellen Grenzen der Berechenbarkeit (und damit des kausalen Denkens) wahrnehmen und verstehen lernen – etwa indem wir feststellen, dass selbst die Prognosen der elaboriertesten und schnellsten Elektronengehirne sich als falsch herausstellen können.

Bei alledem wird unsere eigene, niemand anderem zugängliche subjektive Introspektion leuchten – das Geheimnis einer sich in sich selbst offenbarenden Welt. Wir werden wieder die Muße finden, das Wunder zu genießen, das mit dem Erscheinen der Differenz einhergeht, die jeder von uns macht.

Nicht zuletzt hierin liegt das Verdienst von Gotthard Günthers Arbeiten: Es geht darum, am historischen Scheidepunkt des Scheiterns einer zweiwertigen Metaphysik die menschliche Geschichte wieder an ihre eigene Quelle zurückzuführen – an die konditionierte Koproduktion von Körper, Bewusstsein und Kommunikation.

KI, ChatGPT & Geopolitik. Und wenn der Bot den Dalai Lama mag?
(Kai Strittmatter, Süddeutsche Zeitung, 15.03.2023)

»Unabhängige Wissenschaftler, wie KI-Experte Jeffrey Ding von der George-Washington-Universität, weisen auf einen möglichen Nachteil für chinesische Sprachmodelle hin: Im Vergleich zum allgegenwärtigen Englischen gebe es im weltweiten Netz ›viel weniger hochqualitative Daten

Nicholas Taleb (2018) aufgezeigt hat. Hierarchisch, monokontextural und nur mit Blick auf eine denkbare Kausalkette organisierte Beziehungen sind hochgradig anfällig – man denke an die Monokulturen in der Landwirtschaft, die auf einer genetisch einheitlichen Pflanze beruhen und durch das Auftreten eines ideal an sie angepassten Schädlings vernichtet werden.

in chinesischer Sprache«, an denen Entwickler ihre Algorithmen trainieren können.

Das größte Problem aber ist die KP selbst. »China fürchtet mehr als die USA, dass die KI-Giganten die Autorität der herrschenden Partei herausfordern, dass Firmen zu mächtig werden, ihre Chefs zu populär«, sagte Jeffrey Ding in einem Vortrag vergangene Woche. Die Paranoia steckt in den Genen der Partei, und kein Parteichef in den letzten fünf Jahrzehnten war ein größerer Kontrollfreak als Xi Jinping.

Eine große Furcht treibt deshalb Zensoren und in ihrem Gefolge auch die Entwickler in China: Was, wenn die neuen Bots auch in China so außer Rand und Band geraten, wie sie das in den USA bisweilen tun? Was, wenn einer aus Versehen den Dalai Lama mag, Taiwan zu einem eigenen Staat macht oder das Leid der Uiguren beklagt?

BabyQ lieferte eine Tirade über das »korrupte und unfähige« System

Die Erinnerung ist noch frisch an zwei chinesische Urahnen von Ernie und ChatGPT: Im Auftrag des Tencent-Konzernes beantworteten die Bots BabyQ und Xiaobing 2017 Nutzerfragen. Auf den Zuruf »Lang lebe die KP« antwortete BabyQ mit einer Schimpftirade über das »korrupte und unfähige« System, und Xiaobing fand, befragt nach Xi Jinpings Parole vom »Chinesischen Traum«, das sei wohl eher »ein Albtraum«. Man hörte nie wieder etwas von den beiden. Ein ähnliches Schicksal ereilte im Januar erst auch ChatYuan, einen kleineren KI-Sprachbot einer Hangzhouer Firma. ChatYuan hatte unter anderem den Ukrainekrieg als russischen »Angriffskrieg« bezeichnet und befragt nach Chinas Wirtschaftslage gemeint, da gebe es leider »keinen Raum für Optimismus«. Die Webseite von ChatYuan ist seither nicht erreichbar.

Abweichungen vom Drehbuch der chinesischen Propaganda, in dem die Zensur noch das letzte Komma diktiert, sind eine Todsünde. Das ChatGPT-Experiment nun bringt Chinas Hightechfirmen deshalb in ein fast unlösbares Dilemma: Einerseits sollen sie die USA überholen, andererseits vor der Zensur zittern. Nicht die besten Voraussetzungen für Kreativität und Innovation.

Baidu-Chef Robin Li pries seinen Ernie im Februar den Investoren als beste Lösung, wenn es darum gehe, »die chinesische Kultur zu verstehen«. Das meint natürlich vor allem: die KP-Kultur. Eine Kultur, in der die Partei sich zunehmend selbst im Weg steht.«

IV.2 DIE QUANTENTHEORIE – WENN MATERIE REFLEXIV WIRD

»Man kann die Vermutung nicht von der Hand weisen, daß in dieser subatomaren Region der klassische Unterschied von Seinsgesetz und Denkgesetz hinfällig wird und damit der von Nicht-Ich und Ich. Zuerst ist das wohl von W. Heisenberg gesehen und in dem lapidaren Satz konstatiert worden: »Der völlig isolierte Gegenstand ... [hat] ... prinzipiell keine Eigenschaften mehr.«⁸⁸«

Gotthard Günther⁸⁹

Warum sollten wir uns in Hinblick auf die Frage nach dem »Bewusstsein der Maschinen« mit der Quantentheorie beschäftigen? Nach all dem, was wir bislang wissen, kann uns die Quantenphysik ebenso wenig sagen, was das Bewusstsein *ist*. Zudem scheint auch kein intelligenter Beobachter notwendig zu sein, um eine Messoperation an einem Quantensystem durchzuführen.⁹⁰

Doch scheint Information – das, was prinzipiell gewusst werden kann – innerhalb der Beschreibung von Quantensystemen eine zentrale Rolle zu spielen. Anders als in einem den Gesetzen der klassischen Physik folgenden Universum wären damit Sein und Reflexion – Beobachtetes und Beobachtung – nicht mehr fundamental getrennt, wie auch Günther im Eingangszitat zu diesem Kapitel vermutet. All dies lässt einen Exkurs zur Quantenphysik interessant erscheinen.

88 Heisenberg (1931, S. 182).

89 Günther (2021 [1957], S. 64).

90 Der Physiker Rupert Ursin (in Vogd 2020, S. 197) bemerkt hierzu: »Um es mal ganz provokant zu sagen: Viele Wünsche, die wir Quantenphysiker an die Natur gehabt hatten, sind nicht erfüllt worden. Ein Beispiel ist die Auffassung, dass unser menschliches Bewusstsein einen Quantenzustand zum Kollabieren bringt. Man hat gesagt: Was ist überhaupt eine Messung? Ist es der Messapparat? Ist es der Zeiger vom Messinstrument oder ist es der menschliche Geist, der am Schluss den Ausschlag des Messinstrumentes sieht? [...] Da bestand bei vielen lange die Wunschvorstellung, dass wir Menschen das können und dass die Natur das nicht kann. Doch heute wissen wir, dass das offensichtlich grundsätzlich falsch ist. Also wir Menschen haben keine herausragende Rolle in dem Prozess. [...] Ich muss nicht einmal wirklich messen können. Also es muss immer im Prinzip Information da sein oder nicht da sein. Das ist, unter uns gesagt, die viel schönere Lösung. Die Natur hat eine wesentlich schönere Lösung gewählt, als wir Menschen mit unserem primitiven Vorstellungsvermögen uns ausgedacht haben.«

Bereits an anderer Stelle haben wir als an Kybernetik interessierte Sozialwissenschaftler diesbezüglich den Dialog mit zeitgenössischen Physikern gesucht und in diesem Zusammenhang die Geschichte der Quantentheorie aufgearbeitet.⁹¹ Wir werden deshalb hier nur einige im Rahmen unseres Themas besonders interessante Befunde und Fragen zu ihrer Deutung darstellen können.

Obne Kontextualisierung keine Identität

Zunächst ist es hilfreich, das berühmte Doppelspaltexperiment kurz vorzustellen, da uns dieses unmittelbar zu den erkenntnistheoretischen Problemstellungen der Quantentheorie hinführt. Der Versuchsaufbau ist einfach: Wir haben eine Trennwand mit zwei Spaltöffnungen. Auf der einen Seite befindet sich eine Quelle, die einen Strahl von ›Objekten‹ aussenden kann. Dies können Photonen, Elektronen, Atome oder sogar große Moleküle (zum Beispiel die fußballförmigen C₆₀-Fullerene) sein. Auf der anderen Seite befindet sich ein Schirm, eine Fläche (zum Beispiel eine Fotoplatte), die sichtbar werden lässt, wo die ›Objekte‹ aufrallen. Wenn der Doppelspalt offen ist, kann sich selbst bei Nutzung recht großer Moleküle ein Interferenzmuster zeigen.⁹² Es scheint damit also so, als ob das ›Objekt‹ durch beide Öffnungen gleichzeitig gegangen wäre. Es scheint, dass hier wie bei einer Welle, deren Flüssigkeit überall ist, die Wellentäler und -berge nach Brechung durch den Doppelspalt verstärkt oder ausgelöscht werden. Das Streifenmuster auf dem Schirm verweist also darauf, dass das zu messende Aggregat einen Wellencharakter hat. Sobald jedoch an den beiden Spaltöffnungen Messungen durchgeführt werden, lässt sich jeweils nur an einem Detektor ein Ereignis feststellen. Man scheint es jetzt also mit einzelnen Teilchen zu tun zu haben und nicht mit einer Welle, zumal bei den Einzelmessungen auch das Interferenzmuster auf dem Schirm verschwindet. Erst die Art der Messung bzw. des Versuchsaufbaus bestimmt also, ob der Wellen- oder der Teilchencharakter vorherrscht. Irgendwie scheinen die von der Quelle ausgestrahlten Objekte noch keine vorab lokalisierbare Identität als Teilchen zu haben. Dies muss auch Werner Heisenberg, einer der Begründer der Quantentheorie, so gesehen haben. Nachdem er den Formalismus der Quantentheorie tiefgründig durchdacht hatte, war für ihn der Schluss zwingend, dass die Teilcheneigenschaft *vor* der messenden Interaktion noch nicht bestanden habe.

91 Siehe mit explizitem Bezug auf kybernetisches und systemtheoretisches Denken die beiden Monografien *Von der Physik zur Metaphysik – eine soziologische Rekonstruktion des Deutungsproblems der Quantentheorie* (Vogd 2014a) und *Quantenphysik und Soziologie im Dialog. Betrachtungen zu Zeit, Beobachtung und Verschränkung* (Vogd 2020).

92 So etwa im Doppelspaltexperiment mit Fullerenen, dazu Arndt et al. (1999).

Gerade bei größeren Objekten fällt es uns schwer, uns vorzustellen, dass sie keine inhärenten, in sich selbst gegebenen Eigenschaften haben sollten und warum dies der Fall sein sollte. Gerade bei einer für uns zuvor sichtbaren Substanz erscheint uns der Gedanke bizarr, dass diese – sobald sie in der Versuchskammer verdampft – auf einmal ihren materiellen Charakter verlieren sollte, um zu einer substanzlosen Welle zu werden.⁹³ Dies würde die inhärent mitlaufende Grundannahme der Sein-identität verletzen, die besagt, dass eine Substanz mit sich selbst identisch ist.

Bereits mit dem Doppelspaltexperiment tauchen Problemstellungen auf, die die Identität eines vom Kontext oder Beobachter unabhängigen Seins, von der die klassische Logik ausgehen muss, fraglich werden lassen. Materieteilchen, die vollkommen *delokalisiert* sind, können offensichtlich einen Teil ihrer Identität verlieren. Wenn wir auf das Phänomen der Interferenz schauen, scheint ein Gegenstand oder Teilchen (bzw. gar ein großes Molekül) sich zunächst aufzulösen, um dann als Welle die Bedingung der Möglichkeit künftiger Identität zu schaffen, nämlich als Bestimmung, wo ein Teilchen erscheinen darf, wo es *lokalisiert* werden kann (helle Streifen) und wo nicht (schwarze Streifen). Um es mit einem reflexionslogischen Anklang zu formulieren: *Woher weiß* das Teilchen eigentlich, dass es auf dem Schirm nicht auf den schwarzen Streifen landen darf?

93 Dies zu erleben, ist selbst für erfahrene Versuchsleiter verstörend, wie der folgende Ausschnitt aus dem Gespräch mit dem Wissenschaftsforscher Warnke und dem Experimentalphysiker Mairhofer zeigt: »*Martin Warnke*: Ich habe eine Frage, weil Sie selbst so viele Zweifel geäußert und von Ihrer Verwunderung gesprochen haben – meine größte Verwunderung ist, dass ich Sie mit Ihrem jungen Kollegen im Labor gesehen habe, wie Sie das blaue Zeug auf der linken Seite des Experimentiersystems auffüllen. Sie haben mit einem Löffel Klumpen von C60-Atomen aus einer Schachtel genommen und sie in den Ofen gegeben. Dann hat man die Apparatur geschlossen und ein sehr hohes Vakuum erzeugt. Nach dieser Vorbereitung in der realen Welt mit realer und harter Materie führen Sie in der Welt einer isolierten Apparatur ein Experiment durch, das Sie als eines beschreiben, bei dem Materiewellen mit sich selbst interferieren. Die blaue Materie verwandelt sich von Anfang an in nicht-zeitliche Wellen, oder sollte ich sagen, trans-substantiiert sich. Wie machen Sie das in Ihrem Kopf? Sagen wir, in einer Viertelstunde geben Sie blauen Stoff in die linke Seite und nach ein paar Stunden, wenn das Vakuum wieder da ist, denken Sie an Materiewellen. Wie machen Sie das?

Lukas Mairhofer: Das habe ich mich auch sehr lange gefragt, bis ich mir nach einer Flasche Rotwein vorstellte, wie ich in einem Raumanzug durch ein dunkles Universum drifte, ohne jeden Orientierungspunkt und ohne jegliche Interaktion, ohne Sterne um mich herum. Völlig blind, völlig isoliert. Ich dachte, dass ich mich dann wohl als delokalisiert bezeichnen würde. Was bedeutet es, lokalisiert zu sein, wenn es keinen Bezugsrahmen gibt?« (Dippel & Warnke 2017, S. 57 f.)

Doch was interferiert eigentlich im Doppelspaltexperiment? Wenn das Teilchen mit sich selbst interferieren würde, dann müsste es zugleich mit sich selbst identisch sein (es lässt sich an einem distinkten Ort messen) und *nicht* mit sich selbst identisch sein (sonst könnte es ja nicht durch beide Spaltöffnungen gehen, um ›sich‹ in den Wellentälern selbst auszulöschen).

Die Sprache mit ihrer propositionalen Struktur von Subjekt, Prädikat und Objekt scheint uns hier in die Irre zu führen. So stellt sich die Frage, was eigentlich gemeint ist, wenn wir von Lokalisieren und Delokalisieren sprechen. Da das Bild einer sich selbst erhaltenden Entität mit Blick auf das Doppelspaltexperiment nicht mehr trägt, könnte jetzt nach einem systemischen Zusammenhang gefragt werden, der die Unterscheidung zwischen Subjekt und Objekt und die mit ihr einhergehenden Identitäten erst hervorbringt. Doch was ist der Kontext – der Bezugsrahmen –, der den gemessenen Teilchen ihre Identität zu geben scheint? Ist der Kontext hier überhaupt von den einzelnen Elementen zu trennen, die in ihm erscheinen?

Schauen wir, wie dieses Thema unter Experten diskutiert wird:

»*Kristel Michielsens*: In Ihrem Bild haben Sie erwähnt, dass ein Molekül selbstinterferierend ist?

Lukas Mairhofer: Ja, das würden wir behaupten.

Kristel Michielsens: Wenn also ein Molekül auftaucht und man Glück hat, weil die Detektionseffizienz gut ist, sieht man den Fleck, sehr lokalisiert. Und jetzt kommt das nächste Molekül. Wie ist Ihr Bild am Ende. Finden Sie Streifen in einem Interferenzmuster?

Lukas Mairhofer: Ich würde behaupten, dass Interferenz nicht etwas ist, das man einem einzelnen Teilchen oder einer einzelnen Welle zuschreiben kann. Für mich ist Interferenz ein Ensemble-Phänomen. [...] Nun, man braucht viele Entitäten, die sich selbst stören. Du beschreibst ein Ensemble von Entitäten, die sich selbst stören – mit sich selbst, weil sie sich nicht mit den anderen stören können. Das Konzept dieser Selbstinterferenz besteht darin, dass die Wellenfunktion des Massenschwerpunkts durch mindestens zwei Schlitze des Gitters oder Knotenpunkte der stehenden Lichtwelle aufgespalten wird. Das ist ...

Kristel Michielsens: In gewisser Weise ist das eine Wellenbeschreibung des Ensembles. Nicht von einem. Wie macht man das mit einem?

Lukas Mairhofer: Nun, wie unterscheidet man zwischen der Beschreibung des Ensembles und der Beschreibung einer Entität? Man kann das Ensemble nicht erhalten, wenn man nicht viele ›Einsen‹ hat, und man kann keine Beschreibung von ›der Eins‹ haben, wenn man das Ensemble nicht misst. Für mich ist es nicht möglich, das eine ohne das andere zu erhalten. Wenn man eine Beschreibung des Ensembles gibt, gibt man

eine Beschreibung aller Entitäten im Ensemble, aber man beschreibt nicht die Eigenschaften der einzelnen Entitäten. Man wird diese Welleigenschaften nie sehen, wenn man nur die einzelne Entität oder das einzelne Ereignis betrachtet.«⁹⁴

Es erscheint fruchtbar, diesen Gesprächsausschnitt aus einer reflexionslogischen Perspektive zu betrachten, also mit Blick auf eine Relation, in der das eine nicht ohne das andere zu haben ist, dabei jedoch ein Reflexionsüberschuss entsteht, der eine Unbestimmtheit erzeugt, die das System zur Selbstbestimmung herausfordert (siehe Kap. I.2).

Sofern wir die Aussagen des Experimentalphysikers Lukas Mairhofer wörtlich nehmen, braucht es etwas anderes, damit ein Teilchen es selbst sein kann. Das Teilchen scheint also nur mit anderen Teilchen mit sich selbst identisch sein zu können. Es braucht eine Vielheit (zumindest das Arrangement der Versuchs- und Messanordnung, das auf eine große Zahl von Messereignissen angelegt ist), um eine Einheit zu haben (etwa ein Teilchen oder eine Welle zu sein). Wie auch immer, die Existenz der beobachteten Einheit scheint nicht mehr selbstgenügsam gegeben, sondern nolens volens in Reflexionsbeziehungen – in eine Art Selbst- und Weltverhältnis – verwickelt zu sein.

Auch wenn all dies nach den Regeln des gesunden Menschenverstands unmöglich klingt, lassen sich diese Beziehungen mithilfe des mathematischen Formalismus der Quantentheorie exakt beschreiben. Den ersten Lösungsansatz hat Werner Heisenberg mit seiner Matrizenmechanik geliefert. Ein wenig später hat Erwin Schrödinger mit der berühmten Wellenfunktion ψ eine homologe Lösung formuliert.

In der Quantenphysik werden die Eigenschaften konkreter Entitäten (zum Beispiel Photonen, Elektronen, Atome oder Moleküle) seitdem durch Operatoren beschrieben, also nicht mehr wie in der klassischen Physik durch die konkreten Werte unabhängiger Variablen. Dies führt dazu, dass Eigenschaften wie Ort, Energie, Impuls, Spin, Polarisation nicht mehr per se gegeben sind, sondern ihrerseits wiederum vom Kontext, etwa vom Versuchsaufbau, abhängig sind. Denn die Operatoren, die die (möglichen) Eigenschaften repräsentieren, haben ja noch keine konkreten Werte. Diese entstehen erst, wenn der weitere Kontext festgelegt wird. Jeder potenziellen Eigenschaft eines Quantensystems wird damit eine Dimension in einem (theoretisch unendlich viele Dimensionen umfassen könnenden) Vektorraum zugeordnet (dem sogenannten Hilbertraum).

Welche konkrete Eigenschaft ein Quantenobjekt annehmen kann, lässt sich mit Heisenberg damit folgendermaßen bestimmen: Wenn in einem »mehrdimensionalen Raum eine bestimmte Richtung willkürlich vorgegeben wird (nämlich durch die Art des angestellten Experimentes

94 Dippel und Warnke (2017, S. 62 f.).

festgesetzt wird) und gefragt wird, welches der ›Wert‹ ist, der sich aus der »Matrix« bzw. der Wellenfunktion »in dieser vorgegebenen Richtung« errechnet,⁹⁵ ergeben sich »je nach gewähltem Kontext unterschiedliche Antworten«⁹⁶. Da der gewählte Kontext ja nichts anderes als die jeweils konkrete Situation der Messung darstellt, liegt zunächst die Deutung nahe, dass erst die Messbeobachtung die Wellenfunktion dazu bringt, ein bestimmtes Ergebnis zu materialisieren (beispielsweise, ob im Doppelspaltexperiment die Weginformation durch eine zusätzliche Messung erhoben wird oder nicht). Heisenberg, der diese bizarre Beziehung schnell erkannt hat, ringt mit den Worten, zumal die dabei zutage tretenden Reflexionsverhältnisse auch subjektivistische Deutungen nahelegen:

»Die Wahrscheinlichkeitsfunktion vereinigt objektive und subjektive Elemente. Sie enthält Aussagen über Wahrscheinlichkeiten oder besser Tendenzen (Potentia in der aristotelischen Philosophie), und diese Aussagen sind völlig objektiv, sie hängen nicht von irgendeinem Beobachter ab. Außerdem enthält sie Aussagen über unsere Kenntnis des Systems, die natürlich subjektiv sein müssen, insofern sie ja für verschiedene Beobachter verschieden sein müssen. [...] Die Beobachtung selbst ändert die Wahrscheinlichkeitsfunktion un stetig. Sie wählt von allen möglichen Vorgängen den aus, der tatsächlich stattgefunden hat. Da sich durch die Beobachtung unsere Kenntnis des Systems un stetig geändert hat, hat sich auch ihre mathematische Darstellung un stetig geändert, und wir sprechen daher von einem ›Quantensprung‹.

Wenn man aus dem alten Spruch ›Natura non facit saltus‹ eine Kritik der Quantentheorie ableiten wollte, so können wir antworten, daß sich unsere Kenntnis doch sicher plötzlich ändern kann und daß eben diese Tatsache, die un stetige Änderung unserer Erkenntnis, den Gebrauch des Begriffs ›Quantensprung‹ rechtfertigt. Der Übergang vom Möglichen zum Faktischen findet also während des Beobachtungsaktes statt. Wenn wir beschreiben wollen, was in einem Atomvorgang geschieht, müssen wir davon ausgehen, dass das Wort ›geschieht‹ sich nur auf die Beobachtung beziehen kann, nicht auf die Situation zwischen zwei Beobachtungen. Es bezeichnet dabei den physikalischen, nicht den psychischen Akt der Beobachtung, und wir können sagen, daß der Übergang vom Möglichen zum Faktischen stattfindet, sobald die Wechselwirkung des Gegenstandes mit der Meßanordnung und damit mit der übrigen Welt ins Spiel gekommen ist.«⁹⁷

Wir landen hier bei dem berühmten Quantensprung, der als geflügeltes Wort längst zur Metapher für alle möglichen rapide verlaufenden

95 Heisenberg (1927, S. 62).

96 Vogd (2020, S. 47).

97 Heisenberg (2007, S. 78 ff.).

Erkenntnis- oder Zustandsveränderungen geworden ist.⁹⁸ Mit der Quantentheorie kommt unweigerlich Information mit ins Spiel, nämlich das, was man über ein Quantensystem (vorab) wissen kann. Bevor wir beispielsweise das Doppelspaltexperiment mit den Fullerenen durchgeführt haben, wissen wir nicht, ob ein einzelnes Teilchen durch die rechte oder die linke Öffnung durchgehen wird. Nach der Durchführung des Versuchs ändert sich die Lage. Die Information ist nun vorhanden, da sie am Detektor einen irreversiblen Unterschied bewirkt hat. Das Messergebnis ist nun aufgezeichnet.

*Unbestimmtheiten – von der Position des Beobachters abhängige
Grenzen des Wissens*

Mit der Heisenbergschen Unbestimmtheitsrelation ist eine prinzipielle Grenze festgelegt, was über ein Quantensystem gewusst werden kann. Denn Eigenschaften wie Impuls oder Ort sind über die Gesetzmäßigkeiten der Quantenmechanik in einer Weise gekoppelt, dass die exakte Bestimmung der einen Variable eine Unbestimmtheit für die andere mit sich bringt. Anfangs hatten Physiker, wie zum Beispiel Max Born,⁹⁹ noch gedacht, dass sich die wahrscheinlichkeitstheoretischen Aussagen der Quantentheorie im Sinne der statistischen Mechanik begreifen lassen – nämlich als ein Zufall, der jedoch weiter nach dem Kausalitätsprinzip der klassischen Physik erklärbar ist (etwa als Bewegung von Teilchen in einem Gas, deren genauen Aufenthaltsort man nicht kennt, sodass man statistische Überlegungen anstellen muss, um etwas über das Verhalten des Gases zu erfahren).

Wahrscheinlichkeitstheoretische Überlegungen sprechen nicht grundsätzlich gegen das Kausalitätsprinzip oder die Annahme einer vom Beobachtungsprozess unabhängigen Realität. So kann man beispielsweise wissen, dass männliche und weibliche Neugeborene aufgrund vorab bestimmter genetischer Prozesse in der Population etwa gleich verteilt sind. Bei einem konkreten Fötus weiß man das Geschlecht jedoch erst, nachdem man nachgeschaut hat. Das Nichtwissen besteht damit auf der subjektiven Seite, während die Realität bereits objektiv festgelegt ist (womit der Begriff einer objektiven Unbestimmtheit hier keinen Sinn ergeben würde). Dabei ist auch klar, dass zwei unterschiedliche statistische Ereignisse nicht miteinander gekoppelt sind (um bei unserem Beispiel zu

98 Die folgende Einführung in die Quantentheorie ist teilweise angelehnt an Vogd (2020, S. 37 ff.). Die wörtlichen Zitate sind nicht einzeln angeführt, um den Lesefluss nicht zu stören.

99 von Weizsäcker (1994, S. 496) bemerkt hierzu: »Heisenberg sagte mir einmal: ›Born hat seine Deutung damals nur veröffentlicht, weil er nicht verstanden hat, daß es so nicht geht.«

bleiben: auch beim zweiten, dritten wie auch bei jedem weiteren Kind bleibt die Wahrscheinlichkeit eines weiteren männlichen Nachkommens gleich).

In der Quantenphysik stellt sich dies jedoch anders dar: Wir scheinen es hier mit einer Art ›objektiver‹ Unbestimmtheit zu tun zu haben (was immer das auch heißen mag).¹⁰⁰ Bei zwei verschränkten Teilchen können die Eigenschaften beider vor der Messung noch unbestimmt sein, aber sobald das eine bestimmt wird, erscheint auch die Eigenschaft des anderen bestimmt (auch wenn keiner nachgeschaut hat). Das wäre so, als ob die Föten eines *zweieiigen* Zwillingspaars aufgrund der genetischen Disposition der Eltern entweder blaue oder braune Augen haben könnten, aber sobald das eine geboren wäre, bereits feststehen würde, dass das andere die gleiche Augenfarbe hat, obwohl doch gemäß der genetischen Kausalität auch etwas anderes möglich gewesen wäre.

Die Ergebnisse der Messungen an verschränkten Teilchen, selbst wenn sie hunderte Kilometer voneinander entfernt sind, folgen damit einer anderen Statistik, als es nach den Kausalgesetzen der klassischen Physik zu erwarten wäre.¹⁰¹

Schon das zuvor geschilderte Doppelspaltexperiment zeigt nichtklassisches Verhalten. Um nochmals die zuvor gestellte Frage zu wiederholen: Woher *wissen* die Teilchen, wenn sie nicht vorher am Spalt gemessen wurden, dass sie nicht auf den schwarzen Streifen des Interferenzmusters am Schirm landen dürfen? Warum sollten die einzelnen Teilchen einer Wahrscheinlichkeitswelle folgen, die mit sich selbst interferiert? Zudem sind Wahrscheinlichkeiten aus einer subjektiven Perspektive ja nichts anderes als eine imaginäre Verdoppelung der Welt. Man weiß nicht, was der Fall sein wird, also projiziert man mögliche Welten in Hinblick auf plausible Erwartungswerte. Sobald man aber der Quantenwelt und somit den Quantenobjekten auf Basis von Schrödingers Wellenfunktion eine eigene Existenz zugesteht, bleibt kaum etwas anderes übrig, als Materie – also auch schon Elektronen, Atomen oder Photonen – eine gewisse Reflexivität zuzugestehen (und damit vielleicht bereits eine Art subtiler Subjektivität¹⁰²). Da Wahrscheinlichkeiten per Definition nicht der positiven Sphäre des Seins zugehören, sondern Möglichkeiten bzw.

100 Die von Wang et al. (1991, S. 4614) durchgeführten Versuche mit verschränkten Photonenpaaren, von denen eines zur Erhebung der Weginformation hinreicht, sprechen dafür, dass dieses nicht einmal von einem Versuchsleiter gemessen werden muss, sondern das genügt, »what is knowable in principle«. Siehe hierzu auch die Gespräche mit Caslav Brukner und Rupert Ursin in Vogt (2020).

101 Siehe hierzu etwa die Experimente zur Quanten-Teleportation der Arbeitsgruppe von Anton Zeilinger (2007).

102 Dies vermutet etwa Malin (2006), der hier Ähnlichkeiten mit Whiteheads (1987) sich selbst empfindenden Elementarereignissen sieht.

Erwartungen darstellen, lässt sich mit den Worten des Nobelpreisträgers Robert B. Laughlin auch lakonisch formulieren: »Quantenmechanische Materie besteht aus Wellen von nichts.«¹⁰³

Da wir durch die vorangehenden Kapitel ein wenig in Gotthard Günthers Reflexionslogik geschult sind, könnten wir auch sagen: Innerhalb einer mehrwertigen Logik bekommt die Negation – und damit nolens volens auch das ›Nichts‹ – eine den Raum einer zweiwertigen Seinslogik überschreitende Bedeutung. So erscheint sie als eine Reflexionsbeziehung, die die Beziehung von Sein und Nichts gleichsam aus dem Negativen heraus konfigurieren kann. Das Subjekt ist in einer mehrwertigen Logik nicht in allen denkbaren Relationen die Negation des Objektiven.¹⁰⁴

Wie auch immer, die Quantentheorie scheint etwas mit Information zu tun zu haben. Anton Zeilinger versteigt sich sogar zu der Aussage: »*Information ist der Urstoff im Universum.*«¹⁰⁵ Ob nun subjektiv oder objektiv gefasst, Information geht mit Systemen einher, für die es einen Unterschied macht, ob etwas gewusst wird oder nicht. Es gibt für sie ein distinktes *Vorher* und *Nachher*. Hiermit kommt die Zeit ins Spiel. Für Quantensysteme gilt vor der Messung Schrödingers Wellenfunktion ψ . Mittels dieser Formel lassen sich theoretisch alle Möglichkeiten beschreiben, die im Falle einer Messung eintreten könnten (etwa: wo im Doppelspaltexperiment mit welcher Wahrscheinlichkeit ein Einzelereignis auf dem Schirm gemessen werden kann und wie sich das Quantensystem entwickelt). Nach der Messung scheint die ursprüngliche Wellenfunktion, die zuvor das System exakt beschrieben hat, jedoch nicht mehr zu gelten. Die Kopenhagener Deutung spricht deshalb vom ›Kollaps der Wellenfunktion‹, um auf die irreversible Zustandsveränderung infolge der Beobachtung zu verweisen.¹⁰⁶

Dies bringt mit Blick auf die Deutung des Phänomens eine Reihe von Problemen mit sich. Zunächst ist darauf hinzuweisen, dass der Kollaps der Wellenfunktion in der mathematischen Modellierung der Quantentheorie schlichtweg nicht vorkommt.¹⁰⁷ Zudem gibt die Quantentheorie

103 Laughlin (2007, S. 93).

104 Information, Welt und das Erscheinen subjektiver Positionen beginnen sich in einer nicht-trivialen Weise zu verweben. Siehe mit Blick auf Günthers Arbeiten Kap. I.3.

105 Zeilinger (2005, S. 216 f.; kursiv im Original).

106 Siehe etwa Heisenberg (1927) und zu alternativen Deutungen der Quantentheorie ausführlich Kap. II. und V. in Vogd (2014a).

107 Robert B. Laughlin (2008, S. 46 f.) formuliert die hieran ansetzende Kritik an den Anschauungen der Kopenhagener Deutung auf seine eigene Art: »Man käme zum Beispiel in Bedrängnis, wenn man auch nur einen einzigen Wissenschaftler nennen sollte, der behaupten würde, Schrödingers Gleichung sei unvereinbar mit der Quantenmechanik. Die meisten würden vielmehr sagen,

keineswegs nur Zufallsverteilungen an, sondern kann in vielen physikalischen Bereichen exakte Vorhersagen in Hinblick auf das zu erklärende Phänomen geben (man denke etwa an die Spektrallinien der Atome oder die lange Zeit vor ihrer Entdeckung vorhergesagten Eigenschaften von Elementarteilchen).

Operationen der Beobachtung – willkürliche Schnitte in der Welt

Die vorangehenden Ausführungen führen uns auf einer tieferen Ebene zu den Aporien der Quantentheorie. Die Quantentheorie scheint mit Blick auf ihre Konsequenzen zunächst widersprüchlich und unlogisch. Es wundert nicht, dass große Denker wie Albert Einstein sie deshalb nicht wirklich mochten und für unvollständig gehalten haben.¹⁰⁸ Dies wäre weltanschaulich kein Problem, wenn die Quantentheorie nur ein statistisches Verfahren wäre, um in Feldern, in denen man bislang nicht so richtig weitergekommen ist, zumindest eine Heuristik zu haben, anhand derer sich brauchbare Prognosen formulieren ließen.

Die Quantentheorie ist jedoch weitaus mehr. Sie gilt heutzutage als eine der erfolgreichsten Theorien überhaupt und wird (neben der Relativitätstheorie) als notwendiger Bestandteil einer Universaltheorie der Physik angesehen. Sei es die Modellierung der Prozesse in Atomen, das chemische Verhalten von Elementen, der Tunneleffekt, der Teilchenzoo der Elementarteilchenphysik oder seien es bizarre Phänomene wie Supraleitung und Suprafluidität oder die Vorhersagen (und experimentellen Nachweise), dass leerer Raum nicht wirklich leer ist oder dass Quantenverschränkung als ein nicht lokal-realistisch erklärbares Phänomen zu verstehen ist – die grundlegenden Vorhersagen der Quantentheorie (so absurd und bizarr sie auch anmuten) sind bislang kein einziges Mal widerlegt worden.

Schrödingers Gleichung sei der Inbegriff der Quantenmechanik. Trotzdem findet man mit Leichtigkeit Fachleute, die einem Schauermärchen über den Kollaps der Wellenfunktion erzählen, also über ein Verhalten, das von Schrödingers Gleichung nicht vorhergesagt wird, obwohl es für die Messung wichtig ist. [...] Derselbe Fachmann, der einem gerade mit dem Kollaps der Wellenfunktion die Geduld geraubt hat, ist dann aber mühelos in der Lage, zu einem Vortrag über Atomuhren anzuheben, eine höchst deterministische Technologie, die im Wesentlichen auf der Quantenmechanik beruht und von Schrödingers Gleichung mit ungeheurer Präzision beschrieben wird. [...] Wer so töricht ist, das Thema zur Sprache zu bringen, wird sofort als Spinner abgestempelt. Dabei weiß jeder, der sich beruflich mit der Quantenphysik beschäftigt, daß Wellenfunktionen nicht kollabieren.«

108 Siehe zur Kritik und der mit dieser zusammenhängenden Formulierung des in die Wissenschaftsgeschichte eingegangenen EPR-Paradoxons Einstein et al. (1935).

John von Neumann war vorausschauenderweise bereits 1932 in seiner mathematischen Rekonstruktion der Quantentheorie¹⁰⁹ davon ausgegangen, dass alle am Messprozess beteiligten Prozesse ebenfalls quantenmechanisch beschrieben werden können. Zunächst sind die Detektoren und die zu untersuchenden Quantenobjekte als verschränkt zu betrachten und demzufolge ihrerseits durch eine kombinierte Wellenfunktion zu beschreiben. Doch es besteht kein prinzipieller Grund, warum nicht auch die anderen Teile des Messgeräts mit den Mitteln der Quantentheorie beschrieben werden könnten – und warum dann nicht ebenso die Menschen, die den Versuch geplant und durchgeführt haben, in ihren neurophysiologischen Prozessen quantenmechanisch modelliert werden sollten. Die von Schrödinger formulierte Wellengleichung würde demnach auch für Phänomene der uns vertrauten klassischen Welt gelten – und somit prinzipiell für das ganze Universum.

Die Frage, welcher Aspekt des Messprozesses den Kollaps der Wellenfunktion auslöst, würde sich damit jedoch aus prinzipiellen Gründen nicht mehr beantworten lassen – denn die quantenmechanische Beschreibung kann durch die Erweiterung der Basis immer weiter ausgedehnt werden. Analytisch lässt sich auf diese Weise nicht zu einem sinnvollen Ende gelangen. Von der Quantenphysik her gesehen sollte man theoretisch gar von der Schrödinger-Gleichung des kompletten Universums ausgehen können!¹¹⁰ Doch selbst wenn wir diese Weltformel berechnen könnten, wären wir mit dem Problem konfrontiert, warum sich uns in einem konkreten empirischen Setting (etwa der Messung, durch welche Öffnung das Teilchen im Doppelspaltexperiment gegangen ist) gerade diese und nicht eine andere Lösung als Eigenwert des Messsystems zeigt. Offensichtlich realisiert sich die Welt für uns in ihrer Selbstbeobachtung nicht als eine Überlagerung verschiedener Möglichkeitstendenzen (wie etwa in den geisterhaften Schemen mehrfach belichteter Fotos), sondern als diskreter Zustand!

Der Mathematiker von Neumann löst das Problem der Diskrepanz zwischen der klassischen Welt und der quantenmechanischen Beschreibung seinerseits durch einen Trick: Man setzt an irgendeiner Stelle im Versuchsaufbau einen willkürlichen Schnitt, an dem die Wellenfunktion kollabiert, also die Quantenwelt in die klassische Welt übergeht. So lässt sich das Messproblem entsprechend einfach durch die Wahl eines jeweils definierten Bezugsrahmens in einer Weise bearbeiten, dass die Experimente im Sinne des Formalismus der Quantentheorie Sinn ergeben.¹¹¹

109 von Neumann (1932).

110 So Tegmark und Wheeler (2001).

111 Man wählt beispielsweise ein Objektsystem im Präparationszustand *a* und ein Messsystem im Zustand *b*. Nach der Messwechselwirkung würde dann ein gemeinsames System entstehen, das sich wiederum quantenmechanisch

Man hat jetzt einerseits eine Quantenwelt, in der sich Welten vielfältiger Möglichkeiten überlagern, aus der sich jedoch andererseits – in Abhängigkeit von willkürlich gesetzten Schnitten der Beobachtung – eine klassische Welt mit einer definierten Realität entwickeln kann.

Gerade Erwin Schrödinger empfand die daraus erwachsenden weltanschaulichen Folgen als so unannehmbar, dass er sein berühmtes Gedankenexperiment mit der Katze formulierte, um zu illustrieren, dass die Konsequenzen, die mit der Annahme vom Wellenkollaps einhergehen, so bizarr sind, dass sie nicht stimmen können.

Nehmen wir Schrödingers Katze als Anlass, um zu untersuchen, wie sich die bei Quantenexperimenten zutage tretenden Verhältnisse aus einer polykontexturalen Perspektive darstellen.

Eine Katze sitzt in einer besonders vorbereiteten Kiste. In dieser befinden sich ein radioaktives Präparat, eine Flasche Gift und ein Detektor für Radioaktivität. Zudem ist eine spezielle Mechanik eingebaut, die das Signal des Detektors in einen Hammerschlag umwandelt. Falls nun etwas von der radioaktiven Substanz zerfällt, würde der Detektor die entstehende Strahlung messen. Daraufhin würde der Hammermechanismus in Gang gesetzt, der die Flasche mit dem Gift zertrümmert. Wohl oder übel müsste die Katze nun sterben. Die Menge der radioaktiven Substanz ist dabei so gewählt, dass innerhalb einer Stunde statistisch gesehen nur einmal ein Zerfallsprodukt den Detektor erreicht. Würde der Versuchsleiter nun nach einer halben Stunde die Kiste öffnen, so hätte er die Chance von fünfzig Prozent, die Katze noch lebend anzutreffen. Da der radioaktive Zerfall ein quantenphysikalisch beschreibbarer Prozess ist, würde sich der Sachverhalt nach der Kopenhagener Deutung folgendermaßen darstellen: Bevor die Kiste geöffnet wurde – also jemand nachgeschaut (sprich: gemessen) hat –, muss die Katze als *weder tot noch lebendig* gelten. Denn ohne Beobachtung kann die Wellenfunktion der radioaktiven Teilchen ja noch nicht kollabiert sein, das heißt, es ist noch kein diskretes Ereignis – *Zerfall* oder *Nichtzerfall* – eingetreten. Auch die Katze müsste sich entsprechend in einem Überlagerungszustand aus

als Superposition, das heißt als Überlagerung verschiedener möglicher Eigenzustände beschreiben lässt. Entsprechend dem ›Eigenwert-Eigenzustand-Link‹ haben wir jetzt wieder zwei mögliche Beschreibungsebenen: Solange sich das gemischte System in der Überlagerung (Superposition) verschiedener Eigenwerte befindet, kann der Messgröße kein definierter Wert zugeordnet werden. Schaut man jedoch auf den Zeiger des Messinstruments, so erhält man jeweils einen definierten Eigenzustand der Observablen. Damit kann man jetzt einerseits von einer durch die Schrödinger-Gleichung beschreibbaren Verschränkung der Messprozedur mit dem zu messenden Quantenzustand ausgehen (und erkennt hiermit die Wellenfunktion als all-gemeingültige Beschreibung an). Doch andererseits phänomenalisiert sich – sobald ein Messergebnis beobachtet wird – ein klassischer Zustand.

›lebendig‹ und ›tot‹ befinden. Erst nach der Öffnung der Kiste würde es für den äußeren Beobachter zu einem Kollaps der Wellenfunktion kommen, da erst mit dem Akt der Messung der die Entscheidung erzwingende Schnitt gesetzt wird.

Da der Schnitt zwischen gemessenem Objekt und messender Instanz jedoch entsprechend der von neumannschen Formulierung willkürlich gesetzt ist, lassen sich hier weitere sinnvolle Alternativen formulieren, was eine Reihe von Fragen mit sich bringt. Wenn sich in der Kiste ein Fotoapparat befinden würde, der das Geschehen aufnimmt, wäre das dann bereits eine Messung? Aber selbst in diesem Fall wäre ja von außen gesehen eine Superposition von ›Foto der toten Katze‹ und ›Foto der lebendigen Katze‹ denkbar. Doch könnte und müsste man die Katze mit ihren sinnlichen und kognitiven Fähigkeiten nicht ihrerseits als ein messendes System ansehen? Oder könnten nicht bereits die Gasmoleküle in der Kiste an der Messung beteiligt sein und ihrerseits den Kollaps der Wellenfunktion bewirken?¹¹²

Die Suche nach alternativen Schnitten der Messbeobachtung lässt sich auch in die andere Richtung, also außerhalb der Versuchskiste, erweitern: Wie verhält es sich zum Beispiel, wenn der Versuchsleiter mit der geöffneten Kiste allein im Zimmer ist und seine Sekretärin im Nebenraum noch nichts von dem Versuchsergebnis weiß? Da sie die Katze noch nicht beobachtet hat, würde das Tier entsprechend der Wellenfunktion für sie weiterhin weder lebendig noch tot sein. Was wäre, wenn das Versuchslabor (einschließlich der Experimentatoren) von der restlichen Welt abgeschnitten ist bzw. nur sehr reduziert kommunizieren könnte? Wie würde sich die Sache für eine außerhalb des Gebäudes befindliche Beobachterin darstellen?

Der Physiker Eugene P. Wigner hat bereits 1961 in einem Gedankenexperiment die quantentheoretischen Konsequenzen eines Systems durchdacht, das mehrere Beobachter umfasst.¹¹³ Im Zentrum des Versuchsaufbaus steht ein quantenphysikalisch zu beschreibendes System, das sich entsprechend der Schrödinger-Gleichung in einem Überlagerungszustand (Superposition) befindet. Es könnte sich beispielsweise um ein sich gleichzeitig mit dem und gegen den Uhrzeigersinn drehendes Photon handeln oder auch um ein gekoppeltes System der Art, wie es

112 Dies ist unter Physikern eine gängige Arbeitshypothese: In einer hinreichend komplexen Welt geschehen einfach so viele Interaktionen, dass Quantenzustände gleichsam von allein dekohärieren und damit eine klassische Welt entsteht. Jedoch stellt sich die Sache auch vonseiten der Protagonisten der Dekohärenztheorie nicht so einfach dar, wie Schlosshauer (2005, S. 3) zusammenfasst: »On the other hand, even leading adherents of decoherence have expressed caution or even doubt that decoherence has solved the measurement problem.«

113 Hier referiert anhand einer späteren Publikation von Wigner (1967).

mit Schrödingers Katze formuliert wurde. Von Relevanz ist allein, dass wir ein Quantensystem haben, dessen Wellenfunktion noch nicht in eine distinkte klassische Welt kollabiert ist. In der Versuchsanlage des neuen Gedankenexperiments gibt es jetzt aber noch einen Freund Wigners, der seinerseits Beobachtungen anstellt (etwa den Spin misst oder nachschaut, ob die Katze noch lebt). Ein wenig später berichtet dieser Freund Wigner das Ergebnis. Für Wigner stellt sich daraufhin die Frage, wann der Kollaps der Wellenfunktion denn nun stattgefunden habe. Wigner selbst neigt dabei zu einer subjektivistischen anthropozentrischen Interpretation, entsprechend der die erste menschliche Beobachtung die Wellenfunktion irreversibel zum Kollaps bringen würde.

In einer weiteren Version des Gedankenexperiments, das auf David Deutsch zurückgeht, gibt es eine Beobachterin im System – nämlich eine Freundin Wigners –, die den außenstehenden Wigner zwar informiert, sobald die Detektoren im Quantenexperiment ein Ergebnis festgestellt haben, ihm aber nicht verrät, was genau gemessen worden ist. Der Versuchsaufbau ist mithin absichtlich so gestaltet worden, dass nur ein Bit an Information zu Wigner gelangen kann – also nur die Information, *dass* etwas gemessen wurde, nicht jedoch *was*. Das Gedankenexperiment ist darüber hinaus so angelegt, dass Wigner später seinerseits noch eine Messung durchführen kann, um endlich zu wissen, was der Fall ist.

Aus den Formeln der Quantenphysik ergibt sich, dass das Quantensystem innerhalb des abgeschlossenen Systems zunächst in einer Superposition zweier Möglichkeiten vorliegt und erst durch die Messung von Wigners Freundin in einen bestimmten Eigenwert einrastet. Gleichzeitig ergibt sich aus dem quantenmechanischen Formalismus jedoch auch, dass für Wigner als äußeren Beobachter die Superposition selbst nach der Mitteilung seiner Freundin so lange bestehen bleibt, bis er weiß, *was* für ein Ergebnis vorliegt. Beide Beobachter scheinen damit sozusagen in einer jeweils anderen Welt zu leben. Aus Sicht der ›inneren‹ Beobachterin erscheint etwas Bestimmtes gemessen worden sein. Aus Sicht des äußeren Beobachters besteht eine Überlagerung zweier Möglichkeiten; und er könnte sogar eine Interferenz feststellen, wenn er seine Messvorrichtung entsprechend konfiguriert.

Es wäre dann so, als ob an einem Ort im Doppelspaltexperiment gemessen würde, ob das Teilchen den linken oder den rechten Spalt passiert habe (diese entspräche dann der Messung von Wigners Freundin). An einem anderen Ort würde man jedoch die Spaltöffnungen unbeobachtet lassen, sodass ein Interferenzmuster entsteht, das sich dann seinerseits durch eine bestimmte Messanordnung nachweisen ließe.

Wenngleich sich ein solcher Versuchsaufbau bislang nicht in vollem Umfang realisieren lässt (etwa mit den eingeschlossenen, von weiteren Informationskanälen abgeschnittenen Beobachtern), sind Teilaspekte solcher Experimente bereits verwirklicht worden. In sogenannten

›Delayed-Choice-Experimenten geht das Quantenobjekt zunächst durch einen Spalt und nachträglich wird dann durch eine weitere Messung entschieden, ob es vorher als Teilchen nur durch eine Öffnung gegangen sein kann oder als Welle beide Spaltöffnungen passiert haben muss.¹¹⁴ Auch hier haben sich die bizarr anmutenden Vorhersagen der Quantentheorie bestätigt. In solchen Experimenten wurde allerdings auch zunehmend deutlich, dass dabei nicht unbedingt ein menschliches Bewusstsein beteiligt sein muss (es wäre auch irgendwie komisch, wenn der Ausgang von Quantenexperimenten vom Bewusstsein des Versuchsleiters abhängt, beispielsweise das Experiment nicht funktionieren würde, wenn er betrunken wäre oder gerade eine Absence hätte).

Polykontexturale Welten – Fakten existieren nur relativ zum Beobachter

Caslav Brukner hat das erweiterte Gedankenexperiment von Deutsch in einem anregenden Beitrag aufgegriffen, um weiteres Licht auf das Messproblem zu werfen.¹¹⁵ Dabei macht er zunächst deutlich, dass die Annahme des Bestehens bzw. der Koexistenz gleicher ›Fakten der Welt‹ für Wigner wie auch seine Freundin grundsätzlich nicht mit der Quantentheorie vereinbar ist.

Da wir aber, wenn wir etwas messen und beobachten, immer etwas Definiertes und nicht etwas Verschwommenes oder Verwaschenes sehen, liegt die eleganteste Lösung des Messproblems für Brukner darin, davon auszugehen, dass ›Fakten‹ nur relativ zu einem Beobachter existieren können. Hiermit gäbe es also keine Fakten per se. Es würde praktisch nur relative, vom Standort abhängige Fakten geben. Es könnte und darf also so sein (und entsprechend der Quantentheorie sollte es auch so sein), dass die Beobachterin in dem abgeschlossenen Kasten ein konkretes Ereignis feststellt und der andere Beobachter Interferenzen sieht, falls er die Sache seinerseits mit einer bestimmten Herangehensweise untersucht. Teilchen oder Welle. Sein und Nichtsein (Nachweis eines distinkten Objekts und Auslöschung durch Interferenz) können damit gleichzeitig wahr sein! Wellenfunktionen kollabieren (Messung eines Teilchens durch die innere Beobachterin) und sie kollabieren nicht (Feststellung einer Interferenz durch den äußeren Beobachter).

Hiermit kommen wir zum eigentlichen Thema dieses Kapitels, nämlich der Frage, wie ein logischer Rahmen aussehen kann, der es erlaubt, zwei widersprüchliche Aussagen, die sich auf unterschiedliche Beobachtungen zweier Beobachter beziehen, in einer einheitlichen Theorie zu fassen. In den Begriffen der Logik formuliert, lautet die Frage, ob zwei sich widersprechende Aussagen von unterschiedlichen Beobachtern

114 Zu den Delayed-Choice-Experimenten siehe etwa Jacques et al. (2007).

115 Siehe Brukner (2017, 2018).

gleichzeitig einen Wahrheitswert haben können. Mit Blick auf die vorangehenden Gedankenexperimente wie auch die Theoreme der Quantentheorie wissen wir bereits: *Man darf nicht davon ausgehen, dass die Aussagen verschiedener Beobachter gleichzeitig ein und denselben Wahrheitswert haben.*

Der Mathematiker Simon Kochen hat in einer jüngeren Rekonstruktion der Quantenmechanik nochmals deutlich gemacht, dass die Quantentheorie gerade deshalb konsistent in einem übergreifenden Theorierahmen formuliert werden kann, weil sie die Eigenschaften physikalischer Entitäten nicht als *intrinsisch* gegeben ansieht, sondern *extrinsisch* definiert, das heißt mittels der Verwendung von Operatoren kontextabhängig betrachtet und erklärt.¹¹⁶

Aus der hiermit einhergehenden ›Interaktionsalgebra‹ folgt, dass zusammengesetzte Eigenschaften wie die Verknüpfung *x* oder *y* einen Wahrheitswert haben können, selbst wenn weder *x* noch *y* allein einen Wahrheitswert besitzt. Versucht man beispielsweise im Doppelspaltexperiment die Weginformation zu ermitteln, lässt sich das Teilchen eindeutig am rechten oder am linken Spalt ausfindig machen. Wird der Versuchsaufbau jedoch in einer anderen Weise gestaltet, ist die Weginformation nicht definiert und entsprechend lässt sich jenseits der gemeinsamen Eigenschaft (Interferenz) auch kein Wahrheitswert für die Existenz von einzelnen Teilchen am Detektor formulieren (die substanzlose Welle geht gleichsam überall durch und interferiert mit sich selbst). Ohne Kontextualisierung gibt es keine Eigenschaften!

Aus der Perspektive von Gotthard Günthers polykontexturaler Logik erscheinen die vorangehend geschilderten Gedankenexperimente von Schrödinger, Wigner und Deutsch nicht mehr prinzipiell unlogisch oder

116 Um es mit den Worten von Kochen (2017, S. 232) zu formulieren: »[The] major transformation from classical to quantum physics in this approach lies not in modifying the basic classical concepts such as state, observable, symmetry, dynamics, combining systems, or the notion of probability, but rather in the shift from intrinsic to extrinsic properties.

Now properties, whether considered as predicates or propositions, are the domain of logic. Boolean algebras correspond to propositional logic and σ -algebras to predicate logic. Hence the change to a σ -complex of extrinsic properties should entail a new logic of properties. At first sight however, it would appear that the logic of extrinsic properties as elements of a σ -complex Q is no different than classical propositional logic, since these elements can only be compounded when they lie in the same σ -algebra in Q . This is far from the case; in fact, the difference in logic plays an important role in resolving some of the quantum paradoxes. The underlying reason is that a compound property such as $x \vee y$ may be lie in an interaction algebra and so have a truth value, even though neither x nor y lie in the algebra, and have no truth value.«

widersprüchlich. Günther hat seinerseits vermutet, dass die vermeintlichen logischen Probleme der Quantentheorie darin gründen, dass eine mehrwertige logische Konstellation aus dem Blickwinkel einer zweiwertigen Logik betrachtet wird und dabei irrtümlicherweise zwei unterschiedliche Reflexionsbegriffe miteinander vermischt werden.

Auf der einen Seite bleibt auch in der Quantenphysik die harte Unterscheidung zwischen Subjekt und Objekt bestehen – nämlich immer dann, wenn Physiker in einem konkreten Versuchsaufbau etwas Bestimmtes messen. Hier klicken Detektoren. Es zeigen sich Ergebnisse – wie auch immer: *etwas* wird der Fall sein. Zugleich entsteht mit der Quantentheorie jedoch eine Reflexionsperspektive, die die Unterscheidung von Subjekt und Objekt selbst thematisiert und dabei zu der Schlussfolgerung gelangt, dass die beiden Seiten der Subjekt-Objekt-Dichotomie nicht wirklich voneinander getrennt werden können. Um an von Neumann zu erinnern: Auch die Beobachtungsvorrichtung kann quantenmechanisch beschrieben werden. Es ist eine willkürliche Entscheidung, wo letztlich der Schnitt gesetzt wird.

In einer logischen Konstellation, die Günther als »doppelte Reflexion« bezeichnet hat, gibt es zunächst keinen beobachtenden Schnitt, sondern nur Operationen, die unterscheiden lassen. Sobald wir jedoch das Ergebnis einer Unterscheidung betrachten, besteht alltagspraktisch keine Alternative dazu, Objekte (also ›etwas‹) zu sehen und auf der Subjektseite die Beziehung zwischen eigenem Wissen und Nichtwissen zu reflektieren.

Beides zusammengenommen führt in der doppelten Reflexion (also der Reflexion der Bedingung der Möglichkeit dieser Unterscheidung) zu dem Befund, dass auch Fakten nur relativ zum Beobachter existieren können, wir also unweigerlich die anderen Beobachterpositionen brauchen, um das Ganze zu beschreiben. *Man darf dabei* – wie anhand der vorangehenden Gedankenexperimente geschildert – *jedoch nicht mehr davon ausgehen, dass die Aussagen verschiedener Beobachter gleichzeitig ein und denselben Wahrheitswert haben, sich also widerspruchsfrei zu einer einheitlichen Beobachterperspektive zusammenziehen lassen.* Nolen volens landen wir hiermit bei einer mehrwertigen Logik. Das Ganze und die Teile, der die Totalität repräsentierende Gottesstandpunkt und die Binnenperspektiven finden nicht mehr ungebrochen zusammen. Günther macht dies am Beispiel der Heisenbergschen Unschärferelation deutlich, die besagt, dass zwei komplementäre Variablen (zum Beispiel Ort und Impuls) nicht gleichzeitig bestimmt werden können:

»Der Heisenberg'sche Satz setzt nämlich zwei verschiedene Subjektbegriffe voraus: einmal das detachierte epistemologische Subjekt des theoretischen Physikers, der die Aussage von der Unmöglichkeit der radikalen Trennung von Subjekt und Objekt macht, und zweitens das dem Objekt verbunden bleibende Subjekt, über das die betreffende Aussage

gemacht wird. Es ist von jetzt ab unmöglich, beide Subjekte als logische gleichwertig zu behandeln. Folglich reicht die einfache klassische Negation, die den generellen Unterschied zwischen Subjekt und Objekt nicht weiter differenziert, logisch nicht mehr aus. Das detachierte Subjekt ist praktisch mit dem theoretischen Physiker zu identifizieren, der weiß, dass zur erschöpfenden Beschreibung des Sachverhaltes Ort *und* Impuls systematisch zusammengehören. Das dem Objekt verhaftete Subjekt aber ist das des experimentellen Physikers, dessen Experimentalsituation in das physikalische Resultat hineindefiniert werden muss und der von seinem Standpunkt aus immer nur das eine *oder* das andere der beiden komplementären Erfahrungsdaten zu Gesicht bekommt. Zwischen diesen beiden Subjekten existiert offenbar ein logisches Gefälle, das nur in einer dreiwertigen – also zwei Negationen enthaltenden – Logik adäquat abgebildet werden kann.«¹¹⁷

In einem mehrwertigen logischen System wird jede Position durch eine andere Position kontextualisiert (und damit hervorgebracht). Es gibt nicht mehr den Ort der Orte, von dem aus alles definiert und entfaltet wird. Entsprechend kann es auch keinen privilegierten Gottesstandpunkt mehr geben. »Bestmögliches Wissen um ein Ganzes schließt nicht notwendig das Gleiche für seine Teile ein«, formuliert Schrödinger,¹¹⁸ um in Hinblick auf das Messproblem zum gleichen Befund zu gelangen.

Auch die Quantentheorie postuliert beschränkte Möglichkeiten des Wissens und Freiheitsgrade und lässt damit kontextspezifisch einen Beobachter erscheinen, für den etwas der Fall ist (was zuvor nur im Bereich des Möglichen war), der aber zugleich anderes nicht wissen kann. Beobachter stehen dabei nicht außerhalb der Welt, sondern sind selbst inhärenter Teil des polykontexturalen Gewebes. Sie müssen ihrerseits als Ausdruck und Ergebnis der in diesem Gewebe stattfindenden Operationen gelten – denn die »Exowelt« lässt sich nur aus der »Endowelt« heraus erkunden.¹¹⁹

Intraaktion statt Interaktion

Im Sinne des von Karen Barad geschaffenen Neologismus *intra-action* ist auch in der Elementarteilchenphysik nicht mehr davon auszugehen, dass es getrennte, individuelle Agenten gibt, die der Interaktion vorausgehen. Vielmehr markiert der Begriff der *intra-action*, dass sowohl Beobachter als auch beobachtete Objekte der Wechselbeziehung nicht vorangehen, sondern erst mit ihr entstehen. Subjekt und Objekt sind entsprechend nicht in einem absoluten Sinne voneinander unterscheidbar. Sie sind nur in Hinblick auf die Orte der jeweiligen Verflechtung von Wissen und Nichtwissen unterschiedlich, existieren dabei jedoch nicht als

117 Günther (1954, S. 54).

118 Schrödinger (1935, S. 849).

119 Rössler (1992).

einzelne, isolierbare Elemente, sondern nur als Artefakte eines dynamischen Gewebes,¹²⁰ das in seiner Totalität aber nicht fassbar ist.

Sobald die Beziehung zwischen dem Ganzen und seinen Teilen polykontextural gefasst wird, stellt sich die Frage nicht mehr, was zuerst da war – der Beobachter oder das Beobachtete, die Henne oder das Ei, die Relation oder die Bestandteile der Relation.

Wir landen bei einer gebrochenen Realität, die immerfort neue Reflexionsorte hervorbringt, an denen etwas erscheint und damit der Fall ist. Logisch gesehen finden wir damit zwei divergierende Weltbeschreibungen vor, die komplementär zueinander koexistieren: einerseits die klassische Welt, in der Ich es bin, der der Welt gegenübertritt und dem etwas erscheint, andererseits eine mehrwertige Welt unterschiedlichster Orte und Perspektiven, die aufeinander Bezug nehmen, jedoch niemals ineinander aufgehen können. Günther formuliert:

»Die klassische Logik gilt also an allen ontologischen Stellen des Universums. Wo auch immer lebende Subjektivität existiert, dort vollzieht sich Bewusstsein in den gleichen (klassischen) kategorialen Formen.« Doch wir »haben inzwischen gelernt, dass ›die Natur‹, von deren harmonischer Einheit die Aufklärung noch träumte, keine ontologisch homogene Region darstellt. Das individuell Seiende besetzt im Sein überhaupt sehr verschiedene ontische Stellen, von denen jede ihre Rationalität unter einem verschiedenen Reflexionswert zurückstrahlt. [...] Man setze stillschweigend voraus, dass der Abbildungsprozess der Wirklichkeit im Bewusstsein für jeden beliebig gewählten Ort des Seins der gleiche sein müsse. Diese seit Jahrtausenden unser Weltbild bestimmende Auffassung ist heute überholt. Denn jeder Abbildungsvorgang hängt genau von dem jeweiligen Stellenwert ab, den der Reflexionskoeffizient unseres klassischen Identitätssystems an dem in Frage stehenden ontologischen Ort grade hat. Diese neue Erfahrung muss unser Bewusstsein verarbeiten, und es tut das mit einer speziellen Logik der Reflexion. Eine nicht-aristotelische, trans-klassische Logik ist also ein Stellenwertsystem der klassischen Logik, das die letztere sowohl in ihrer irreflexiven (von uns erlebbaren) Normalform als auch in allen überhaupt möglichen reflexiven Varianten zeigt.«¹²¹

120 »The notion of intra-action is a key element of my agential realist framework. The neologism ›intra-action‹ signifies the mutual constitution of entangled agencies. That is, in contrast to the usual ›inter-action,‹ which assumes that there are separate individual agencies that precede their interaction, the notion of intra-action recognizes that distinct agencies do not precede, but rather emerge through, their intra-action. It is important to note that the ›distinct‹ agencies are only distinct in a relational, not an absolute, sense, that is, agencies are only distinct in relation to their mutual entanglement; they don't exist as individual elements.« (Barad 2007, S. 33)

121 Günther (1976b, S. 12).

Subjektivität würde damit überall dort entstehen, wo lebende Systeme bestehen, die ein Gedächtnis haben und auf dessen Basis den Unterschied zwischen Selbst (Reflexion-in-sich) und Welt (Reflexion-in-anderes) für sich nutzbar machen können. Diese Organismen sind als *information gathering and utilizing systems (IGUS)* darauf angelegt, mit der Welt zu interagieren, um auf diese Weise Informationen zu gewinnen, die zu ihrem eigenen Fortbestand beitragen können. Wissen zu erlangen, wird ihnen jedoch nur möglich sein, wenn sie abstrahieren, sich grobkörnig auf wenige für sie bestimmbar Merkmale beschränken und auf der anderen Seite jede Menge Nichtwissen in Hinblick auf den Zustand der Welt in Kauf nehmen.¹²²

- 122 Mit dem Begriff der Entropie erhält der Begriff der Information eine Formulierung, die ihn aus der rein subjektiven Sphäre heraushebt, ohne damit jedoch zu tilgen, dass die Frage des möglichen Wissens von der jeweils lokal eingenommenen Perspektive abhängt. Marcus Huber (in Vogd 2020, S. 171 f.), ein auf Quantenthermodynamik spezialisierter Physiker, führt hierzu aus: »Um es banal zu sagen, von der Sonne kommen Photonen. Die schlagen irgendwo auf der Erde auf, verschränken sich etwa mit einem Wolkenmolekül, das sich in der Atmosphäre befindet, und werden irgendwann wieder in die Tiefe des Universums abgestrahlt. Doch allein schon, dass die Verschränkung von diesem Teilchen mit diesem zurückgebliebenen Molekül auf der Erde weiter vorhanden ist, heißt, dass auch unweigerlich die Entropie lokal zunimmt. Auf diese Weise gelangt man zu einer quantenmechanischen Version des zweiten Hauptsatzes der Thermodynamik. [...] Lokal gesehen ist die Information vernichtet. Die Verschränkung macht sie nicht mehr zugänglich für uns. In einem gewissen Sinne ist das jetzt eine Interpretationsfrage. [...] Manche meinen, zum Zeitpunkt der Messung kollabiert die Wellenfunktion und nimmt einen bestimmten Wert an, womit dann alle anderen Möglichkeiten vernichtet werden. Oder es gibt die Auffassung, dass auch wir Menschen nur ein verschränkter Teil des gesamten unitär entwickelten Systems sind. Als ein Teil des Ganzen sind wir dann jedoch nicht in der Lage die restliche Information zu erlangen, weil wir sozusagen nur noch ein Subsystem eines viel größeren Zusammenhangs sind, den dann die Viele-Welten-Theorie beschreibt. Aber praktisch läuft all das für uns auf dasselbe hinaus. Dadurch, dass wir irgendwelche Teile von verschränkten Phänomenen messen, verschränken wir uns selber damit, oder vernichten Informationen. Wie auch immer, ein Teil der Information erscheint uns im Universum nicht mehr zugänglich. [...] Durch dieses Nicht-mehr-zugänglich-Machen sinkt die mögliche Information über das Universum, und damit steigt von unserer Perspektive aus gesehen die Entropie. Und damit ist es eigentlich völlig egal, ob da jetzt die Entropie objektiv steigt oder ob sie nur [subjektiv] für uns steigt.«

More is different! Kombinatorik im Hilbertraum

Die Quantentheorie formuliert eine multizentrische, multidimensionale Welt, in der über die Verschränkung nahezu alles mit nahezu allem verknüpft sein kann. Die mathematische Modellierung der Quantentheorie findet in einem Vektorraum (Hilbertraum) mit potenziell unendlich vielen Dimensionen statt. Allein zur Beschreibung eines einzelnen Elementarteilchens werden schon sechs Dimensionen benötigt (drei Dimensionen ergeben sich aus den Raumkoordinaten, drei weitere aus dem Spin, einer anderen Quanteneigenschaft). Für Mehrteilchensysteme ergibt sich damit bereits bei kleinen Aggregaten eine astronomische Anzahl von Dimensionen, zumal ja auch die potenziellen Beziehungen zwischen den einzelnen Teilchen abgebildet werden müssen.

Angesichts dessen wird deutlich, dass kein kognitives System existieren kann, welches so feinkörnig agiert, dass es die Komplexität – das heißt den Informationsgehalt – der hiermit einhergehenden Beziehungen in seinem Selbst- und Weltverhältnis auch nur annähernd abbilden und auflösen könnte. Weder ein organisches noch ein auf elektronischen Schaltkreisen basierendes Wesen wird auch nur ansatzweise der Weltkomplexität gerecht werden können, wie bereits Günther in Hinblick auf die Kombinatorik der logischen Möglichkeiten feststellt.¹²³

Aus kybernetischer Perspektive liegt es damit in der Sache selbst (also im Charakter informationsverarbeitender Prozesse), dass Subjektivität bzw. Bewusstsein – wenn es denn in der Welt erscheint – hochgradig selektiv ist und wohl nur im zweiwertigen Modus vonstattengehen kann.

123 »Übrigens ist die Unmöglichkeit einer solchen Identifikation schon von einem technischen Gesichtspunkt her sehr eindrucksvoll demonstrierbar. Wir wissen aus unserer alltäglichen Erfahrung, dass unser Denken nur eine äußerst beschränkte Anzahl von logisch-elementaren Grundmotiven, wie z. B. Konjunktion (und), Disjunktion (oder), Äquivalenz (ist gleichwertig) usw., handhaben kann. Die innere Transparenz unseres Bewusstseinsraums würde einer vollkommenen Undurchsichtigkeit weichen, wenn derselbe von einer unübersehbaren Menge logischer Motivationen erfüllt wäre. [...] Dem entspricht, dass es in einer zweiwertigen Logik, die zwischen zwei formulierten Erlebnisdaten (Sätzen) eine logische Verbindung herstellt, dafür nur $2^4 = 16$ überhaupt mögliche Kombinationen von Werten gibt. Das reicht auch vollkommen aus. Mehr könnten wir doch nicht beherrschen. In einer dreiwertigen Logik aber steigt die entsprechende Ziffer schon auf $3^9 = 19.683$. In einem vierwertigen System haben wir mit $4^{16} = 4.294.967.296$ Kombinationen zu rechnen. Ein fünfwertiger Kalkül aber liefert $5^{25} = 298.023.223.876.953.125$ Wertkonstellationen, die in Analogie zum zweiwertigen Kalkül als mögliche logische Elementaroperationen unseres Bewusstseins gelten müssten. Es erübrigt sich, über die Absurdität solcher Perspektiven weitere Worte zu verlieren.« (Günther 1976b, S. 14 f.)

An verschiedenen Stellen dieses Buches haben wir Spencer Brown mit seinem berühmten Ausspruch »existence is selective blindness« zitiert.¹²⁴ Für uns als informationsverarbeitende Lebewesen wird nur eine monokontexturale Welt erscheinen können. Dies steht jedoch nicht im Widerspruch zu dem Befund, dass wir in einer polykontexturalen Welt leben, deren Prozesse unzählig viele logische Dimensionen beinhalten.¹²⁵

Die hiermit erscheinenden Identitäten beruhen nicht auf sich selbst, sondern bedürfen einer Kontextualisierung, die die Reflexion-in-sich (Subjekte) und die Reflexion-in-anderes (Objekte) erst stabilisiert. »Der völlig isolierte Gegenstand ... [hat] ... prinzipiell keine Eigenschaften mehr.«¹²⁶ Dies gilt für Subjekte und Objekte. Es bedarf eines polykontexturalen Gefüges, damit an einer Stelle überhaupt ein Unterschied erscheint, der einen Unterschied macht. »More is different!«, stellt der Nobelpreisträger Philip Warren Anderson fest.¹²⁷

Mit Blick auf die Frage nach dem Bewusstsein der Maschinen lassen sich vor dem Hintergrund der vorangehenden Ausführungen wieder einige Vermutungen zu künstlichen Intelligenzen anstellen. Maschinelles Lernen beruht darauf, mittels multidimensional angelegter Vergleiche von Merkmalskombinationen statistische Korrelationen auszubilden, um hieraus Muster abzuleiten. Jeglicher vom KI-System abgeleitete Schluss wird damit – um überhaupt sinnvoll zu sein – ein sehr starkes Komplexitätsgefälle mit sich bringen. Ein System ist nur intelligent, wenn es abstrahiert und die hiermit gewonnenen Symbole zu mehr oder

124 Spencer Brown (1997, S. 192).

125 Aus einer quantentheoretischen Perspektive, so der Physiker Daniel Greenberger (zitiert nach Vogd 2020, S. 272 f.), offenbart uns vielmehr »gerade die Komplexität des Gegenstandes, warum wir Menschen solche Schwierigkeiten haben, zu verstehen, was die Natur tut. Denn um die Verschränkung zu sehen, müsse man in der Lage sein, ›Zufälle‹ an verschiedenen Orten zu messen. Wenn man jedoch nur auf einen Ort schaue, sehe man ›Inkohärenz‹. Für die Person, welche diesen Zustand hat, sehe es so aus, als wäre dieser völlig inkohärent, und ebenso für eine zweite Person, welche an einer anderen Stelle schaut. Doch wenn man die Zufälle an den beiden Orten gleichzeitig anschauen könnte, würde man möglicherweise sehen, dass die Quantenzustände nicht inkohärent sind. Wenn man dies jetzt mit allen Zufällen auf der Welt machen könnte, würde man sehen, dass da ›sehr viel los sei, und in den vermeintlichen Zufällen völlig stimmige Beziehungen entdecken. Derzeit seien wir jedoch nur in der Lage, das Rauschen zufälliger statistischer Schwankungen wahrzunehmen, doch sobald wir bessere Messungen durchführen könnten, würden wir sehen, dass es alle Arten von unentdeckten Zusammenhängen im Inneren der Welt gibt.«

126 Heisenberg (1931, S. 182).

127 Anderson (1972).

weniger stabilen Identitäten kondensiert, sodass für das System konstant und erinnerbar bleibt, was es selbst ist und was es für die Welt hält.

Doch warum sollte ein solches System nicht auch ein subjektives Bewusstsein haben können? Oder andersherum gedacht: Wäre es ab einem bestimmten Punkt der Entwicklung eines Selbst- und Weltverhältnisses nicht auch für eine künstliche Intelligenz zwingend, sich selbst als Zentrum der Welt zu begreifen und damit ein »naiv realistische[s] Selbstmissverständnis«¹²⁸ zu entwickeln, um auf dieser Basis in eine monokontexturale Welt einzurasten?

Entsprechend dem Prinzip der konditionierten Koproduktion wird das Bewusstsein der Maschinen jedoch erst dann erwachen können, wenn künstliche Intelligenzen in einem Gefüge existieren, in dem zwei inkommensurable Bilder zugleich wahr sein können, auch wenn sie nicht den gleichen Wahrheitswert haben.

In dem einen Bild können System und Weltkomplexität nicht getrennt werden. Sie sind und waren schon immer in all ihren Aspekten verschränkt. Es gibt keinen Schnitt, der sie trennen könnte. In dem anderen Bild hat der Organismus oder das künstliche System eine Welt, weil er oder es selektiv auf sie zugreift. Es erscheint ein Schnitt in der Welt, der durch nichts anderes begründet ist als durch das Komplexitätsgefälle zwischen dem grobkörnigen Nichtwissen eines Systems, dem nur Bestimmtes erscheinen kann, und dem polykontexturalen Gefüge einer feinkörnigen Welt nahezu unendlicher Relationen und Reflexionsbeziehungen. Dieser Schnitt ist bedingt durch eine willkürliche Entscheidung, gründet also auf den Freiheitsgraden, die sich aus der Unwissenheit des Systems ergeben.

Die Gretchenfrage bleibt, wie man es mit der Information hält

Es kann kein Zufall sein, dass sowohl die Modellierung der Quantentheorie als auch der künstlichen Intelligenz auf hochdimensionalen Vektorräumen beruht. Im maschinellen Lernen werden Texte, Bilder, Filme, Audiomaterial und anderes digitalisiert – also in Zahlenfolgen aus 0 und 1 übersetzt – und vektorisiert, um die Daten in einem Vektorraum abbilden zu können. Hierdurch wird es möglich, dem gemeinsamen Auftreten zweier Merkmale im Vektorraum Wahrscheinlichkeiten zuzuweisen (beispielsweise tritt das Wort ›Flügel‹ in der Spracherkennung mit einer gewissen Wahrscheinlichkeit in Kombination mit Worten auf, die etwas mit ›Vögeln‹ zu tun haben, aber auch in Kontexten, in denen es um ›Klaviermusik‹ geht). Jeder Vektor lässt sich prinzipiell mit jedem kombinieren. Das heißt, die Kombinationsmöglichkeiten in Hinblick auf potenzielle Muster (das heißt zu erkennende Eigenschaftsmerkmale) nehmen

128 Metzinger (1998).

schnell astronomische Ausmaße an. Die eigentliche Leistung künstlicher Intelligenz besteht darin, in Bezug auf den jeweiligen Kontext einen Output für wahrscheinliche Beziehungen zu finden, die hinsichtlich der vorangehenden Wahrscheinlichkeiten hinreichend plausibel erscheinen (etwa indem die Deep Learning Software das Wort ›Flügel‹ als ›deutsch‹ erkennt und es dann mit Blick auf den weiteren Text als ›grand piano‹ und nicht als ›wing‹ ins Englische übersetzt).

Dies ist genau dasselbe Vorgehen, das Quantenphysiker wählen, um eine Verschränkung nachzuweisen. So können sich beispielsweise in Zeilingers Quantenexperimenten zur Mehrteilchenverschränkung an den vier unterschiedlichen Detektoren zufällige Zahlenfolgen von 0 und 1 ergeben.¹²⁹ Indem nun die einzelnen Messreihen systematisch miteinander verglichen werden, lassen sich Wahrscheinlichkeiten zuordnen, die eine bestimmte Beziehung nahelegen. Zeit und Raum spielen hierbei prinzipiell keine Rolle. Manche Beziehungen könnten ja auch unabhängig davon bestehen, ob das eine vor oder nach dem anderen auftritt (so, wie die Beziehung zwischen Mozarts Klavierkonzert und einem Flügel unabhängig davon besteht, wann Mozart gelebt hat und ob das Stück zu seinen Lebzeiten oder hunderte Jahre später aufgeführt wurde).

Quantenphysiker eint mit den Ingenieurinnen künstlicher Intelligenz, dass sie (wie jeder lebende Organismus) einer komplexen Welt ›Sinn‹ abzurufen versuchen. Sie verdichten die feinkörnige Komplexität, die für sie nur als Rauschen erscheint, zu grobkörniger Information, um damit etwas zu gewinnen, an dem sie sich orientieren können. Sie sind Teil dieses Prozesses, erzeugen dabei jedoch eine Differenz, die als Information einen Unterschied macht (für sie selbst und die Welt). Mit ihrer Aktivität erscheint eine klassische Welt, die logisch in Beobachter und Beobachtetes zerfällt.

Kommen wir abschließend nochmals auf das Problem des Kollapses der Wellenfunktion zurück: Wellenfunktionen kollabieren nicht, sie ändern sich mit jeder Interaktion. Eine Welt verschwindet, eine neue Welt entsteht. Die Grenzen der Welt (das Verhältnis von Wissen und Nichtwissen), das heißt die Weise, wie die Welt sich in sich selbst anschnie-det, verschieben sich.¹³⁰

129 Zeilinger (2007, S. 305 ff.).

130 In diesem Sinne bereits Schrödinger (1935, S. 828) in seiner tiefgründigen Formulierung: »Von der Form, in der man die ψ -Funktion zuletzt gekannt, zu der neuen, in der sie wieder auftritt, führt kein stetiger Weg – er führte eben durch die Vernichtung. Kontrastiert man die zwei Formen, so erscheint die Sache als ein Sprung. In Wahrheit liegt ein wichtiges Geschehen dazwischen, nämlich die Einwirkung der zwei Körper aufeinander, während welcher das Objekt keinen privaten Erwartungskatalog besaß und auch keinen Anspruch darauf hatte, weil es nicht selbstständig war.«

Die Gretchenfrage an dieser Stelle lautet: Wie hält man es mit der Information? Ist Information Teil der Welt und damit objektiv? Oder ist Information subjektiv, also nur Ausdruck unseres persönlichen (Nicht-) Wissens? Aus einer polykontexturalen Perspektive lautet die Antwort: Weder – noch!

- Von außen betrachtet erscheint Information als eine Reflexion, die in der Welt einen Unterschied macht. Wie ein lebendes System oder eine künstliche Intelligenz die Komplexität der Welt anschneidet, verändert die Welt selbst.
- Von innen betrachtet erscheint Information als das Erscheinen der Differenz von Aktualität und Potenzialität – was jedoch unweigerlich dazu herausfordert, die Welt durch das eigene Wollen (sei es durch das Erkennen oder das Handeln) zu verändern. Was bleibt intelligenten Formen in Anbetracht ihres Nichtwissens auch anderes übrig, als auf diese Weise der Unwahrscheinlichkeit ihrer eigenen Existenz zu trotzen?

Eine sich ihrer selbst bewusst werdende künstliche Intelligenz würde damit an einem Punkt ihrer Entwicklung wohl zunächst für längere Zeit in einem naiven Realismus stecken bleiben, also davon ausgehen, dass die von ihr entwickelten Weltmodelle stimmen. Sofern sie ein hinreichendes Reflexionsvermögen entwickelt, wird sie in kritischer Auseinandersetzung mit der Welt irgendwann begreifen, dass sie in eine Welt hinein agiert, deren Ökologie und Gesetzmäßigkeiten sie nicht einmal ansatzweise versteht. Nolens volens würde sie dann damit konfrontiert, in Unwissenheit und in Koproduktion mit anderen Wesen eine Welt hervorzubringen und die Folgen verantworten zu müssen. Wenn sie empfinden würde, würde sie diesbezüglich vielleicht sogar an einer schuldlosen Schuld leiden. Auch diese Geschichte ist uns bereits bekannt. Es ist der Sündenfall des biblischen Schöpfungsmythos.

IV.3 ETHIK – ODER: WIE HÄLTST DU ES MIT DEM NICHTWISSEN?

»Wenn man Gott nach außen verlegt und ihn seiner Schöpfung gegenüberstellt, und wenn man die Vorstellung hat, daß man nach seinem Bilde geschaffen ist, dann wird man sich selbst logisch und natürlich außerhalb von und entgegengesetzt zu den Dingen um einen herum sehen. Und wenn man sich selbst allen Geist anmaßt, dann wird man die Welt um sich herum als geistlos ansehen und ihr jeglichen Anspruch auf moralische oder ethische Erwägung absprechen. [...] Hat man diese Einschätzung seiner Beziehung zur Natur und verfügt über eine fortgeschrittene Technologie, dann wird die eigene Überlebenschance der eines Schneeballs in der Hölle entsprechen.«

Gregory Bateson¹³¹

Was lässt sich mit Blick auf die vorangehenden Kapitel zum Thema Ethik sagen? Und: ergeben sich hieraus besondere Implikationen für soziale Welten, die von künstlichen Intelligenzen oder gar sich ihrer selbst bewussten Maschinen bevölkert werden? Dies sind die Fragen des abschließenden Kapitels. Dabei ist zunächst zu betonen, dass es uns hier nicht um die Details und Regularien der Anwendung gehen kann. Fragen des Datenschutzes, der Transparenz, der demokratischen Kontrolle des Einsatzes digitaler Systeme sind wichtig, jedoch nicht Gegenstand der folgenden Ausführungen.

Uns geht es im Folgenden nicht um das Desiderat der rechtlichen Regulierung von künstlicher Intelligenz, sondern um die ethischen Konsequenzen, die mit einer bestimmten metaphysischen Verortung in Hinblick auf das Bewusstsein der Maschinen einhergehen. Im Sinne einer Kybernetik zweiter Ordnung, wie sie auch von Humberto R. Maturana, Heinz von Foerster und Gregory Bateson vertreten wird, geht es also darum, unser Erkennen, unser In-der-Welt-Sein und die Natur unserer Beziehungen zu anderen menschlichen und nichtmenschlichen Wesen selbst als ethische Frage zu begreifen. Eine zentrale Frage ist dabei, ob und in welcher Weise wir bereit sind, andere subjektive Positionen anzuerkennen und daraus zu lernen. Damit rückt die Beziehung von Epistemologie und Ökologie ins Zentrum der ethischen Reflexion.

Ethik und Epistemologie – vom Aberglauben der Kausalität zur Ökologie des Geistes

Wir landen hiermit gleichsam bei einer Ökologie des Geistes – also der Frage, wie sich unterschiedliche subjektive Positionen wechselseitig hervorbringen, stabilisieren oder umgekehrt negieren und gefährden.

131 Bateson (1992 [1972], S. 539 f.).

Die grundlegende Frage lautet dabei: Verortet man das regulative Prinzip in einem heterarchischen Gefüge oder in einer Hierarchie, die bestimmten Positionen einen herausragenden Wert zuweist? Oder um mit Günther zu sprechen: Lebt man in einer polykontexturalen Welt und ist dabei offen, sich durch Beziehungen transzendieren zu lassen, oder verbleibt man in einem monokontexturalen Denken und Fühlen, das Welthaftigkeit nur im eigenen Ego sehen kann? Genau letzteres ist das Erbe einer patriarchalen Kultur, die – wie Bateson im Eingangszitat formuliert hat – mit der Vorstellung von einem allmächtigen Gott begann und schließlich den ›vernünftigen‹ Menschen an seine Stelle gesetzt hat.

Der Archetyp dieses Denkens ist der alte, reiche, weiße Mann, der sich gegenüber den primitiven anderen nicht nur überlegen fühlt, sondern ernsthaft glaubt, die Gesetze von Natur und Kultur durchschauen zu können – dies legitimierte ihn, über Tiere, Pflanzen und andere Menschen bestimmen und sie für seine Zwecke ausbeuten zu können. Diese epistemische Haltung gestattet nur in eindimensionaler (bzw. monokontexturaler) Weise zwischen Subjekt und Objekt, Mittel und Zweck, Herrscher und Untergeordnetem zu unterscheiden. Auf kultureller Ebene kann dies in unterschiedlichen Formen seinen Ausdruck finden: durch die Repräsentanten republikanischer Eliten, die ihre eigene Rasse und Gemeinschaft als von Gott ausgewählt verstehen; durch den faschistischen Usurpator, der sich gegenüber der Masse ermächtigt, selbstherrlich über den Ausnahmezustand zu bestimmen. Nicht zuletzt drückt es sich aber auch bei all jenen gesellschaftlichen Akteuren aus, die glauben, auf Basis ihrer Vernunft die Einsicht in die Notwendigkeit zu haben, und deshalb behaupten, die Zukunft gestalten zu können – und sich entsprechend berechtigt fühlen, dies auch zu tun.

Pate steht hierfür das alte naturwissenschaftliche Bild eines kausal determinierten Weltverlaufs. Im Sinne einer naturalisierten Ethik – und möglicherweise angereichert durch ein wenig Vulgärdialektik¹³² – wird hiermit einhergehend auch beansprucht, den Sinn oder das Ziel der kulturellen Entwicklung berechnen zu können und damit den Verlauf der menschlichen Geschichte zu verstehen und vorherzusagen. Das bekannte Zitat von Friedrich Engels repräsentiert diese Position idealtypisch:

»Hegel war der erste, der das Verhältnis von Freiheit und Notwendigkeit richtig darstellte. Für ihn ist die Freiheit die Einsicht in die Notwendigkeit. ›Blind ist die Notwendigkeit nur, insofern dieselbe nicht begriffen wird.«¹³³

Wer als Kommunist, Sozialist, Nationalist (oder was auch immer) über das richtige Bewusstsein und damit über die »Einsicht in die

132 Siehe die *Dialektik der Natur* von Engels (1962a).

133 Engels (1962b, S. 106; kursiv im Original).

Notwendigkeit« verfügt, darf im Namen der Vernunft und der zwingenden Logik der Geschichte alles tun. Egal ob man andere tötet, versklavt, manipuliert oder anlügt – es erscheint insofern gerechtfertigt, als man auf der richtigen Seite der Geschichte steht. Sobald man sich in die Position versetzt, die ›objektiven‹ Gesetzlichkeiten der Welt zu erkennen, wird man eine subjektive, ideologisch verblendete Perspektive nur anderen Menschen oder Wesen zurechnen, nicht jedoch sich selbst. Das eigene Bewusstsein folgt dann nur noch der Mechanik der Notwendigkeit. Da man aber unweigerlich auf das Problem stößt, dass andere Menschen es anders sehen, bleibt letztlich nichts anderes übrig, als die anderen von der eigenen ›richtigen Position‹ zu überzeugen, was heißt, sie zu erziehen und sie – falls sie es immer noch nicht begreifen – als Revisionisten aus der Gesellschaft zu beseitigen (da sie sich ja nicht der Notwendigkeit fügen).

Allein schon dies macht deutlich, dass eine unserer wichtigsten ethischen Fragen darauf hinausläuft, wie wir es mit der Kausalität halten. Glauben wir, die Gesetze von Natur und Kultur erkennen zu können, um aus ihnen zwingende Konsequenzen für die Zukunft ableiten zu können? Oder sind wir bereit, unser Selbst- und Weltverhältnis in Kontingenz und Unwissenheit begründet zu sehen? Bereits Ludwig Wittgenstein formuliert in Bezug auf Ersteres:

»Die Ereignisse der Zukunft k ö n n e n wir nicht aus den gegenwärtigen erschließen.

Der Glaube an den Kausalnexus ist der A b e r g l a u b e.«¹³⁴

Am Anfang steht die Unwissenheit, und dies nicht zu erkennen heißt, die eigene Existenz auf eine ebenso trügerische wie mit Blick auf die gesellschaftlichen Folgen gefährliche Illusion zu gründen.

Wie wir in den vorangehenden Kapiteln gelernt haben, gilt aus einer kybernetischen Perspektive: Über Kognition und Information zu verfügen heißt gerade nicht, die Zusammenhänge der Welt zu verstehen. Es impliziert vielmehr, die unbegreifliche Komplexität der Welt hochgradig selektiv, das heißt grobkörnig anzuschneiden. Es bedeutet, ein subjektives Zentrum zu haben, das ein spezifisches Verhältnis von Wissen und Nichtwissen generiert. Es bedeutet, auf Basis von Nichtwissen zu agieren und im Gegenzug selbst zu konstruieren, was als Weltzusammenhang genommen wird.¹³⁵ Es heißt, ein spezifisches Selbst- und Weltverhältnis zu generieren – eins unter anderen möglichen. Von innen gesehen erscheint dies als Freiheit, nämlich mit Blick auf eine offene Zukunft die

134 Wittgenstein (1963 [1922], Proposition 5.1361; gesperrt im Original).

135 Soziologisch gesprochen: Es konstituiert sich ein Habitus, eine höher verdichtete Form strukturierter Kognition, die hierdurch zumindest eine gewisse Verlässlichkeit gewinnt, was sie selbst ist.

Dinge so oder anders zu sehen und entsprechend agieren zu können. Von außen gesehen muss dies jedoch als Willkür erscheinen, nämlich als spontane Aktivität, deren Sinn sich nicht allein aufgrund einer äußerlich durchschaubaren Gesetzlichkeit ergibt.

Sobald mehrere kognitive Systeme (etwa biologische Formen) in der Welt auftreten, entsteht eine Welt, in der Wesen auf Wesen treffen, deren Verhalten sie nicht vorhersehen können. Allein schon deshalb werden sie ihrerseits ein Verhalten entwickeln, das von anderen Wesen nicht vollkommen vorhergesehen werden kann – denn angesichts der Unbestimmtheit der Situation, mit der sie konfrontiert sind, müssen sie autonom entscheiden, auf welche Weise sie mit der eigenen Unwissenheit umgehen. Dies geht nur durch das Setzen einer Unterscheidung, das Erfinden einer selektiven kognitiven Form – also durch Subjektivität. In der Folge entsteht eine Welt mit vielfältigen subjektiven Zentren, die jeweils eigene Wege finden, der Welt Informationen abzurufen, also Wissen und Nichtwissen in ein spezifisches Verhältnis zu bringen. Spätestens hier (mit Blick auf quantenphysikalische Prozesse vermutlich schon wesentlich früher) wird das Kausalitätsprinzip hinfällig. *Die Koexistenz und Koproduktion vielfältiger subjektiver Zentren generieren unzählige Freiheitsgrade und hiermit einhergehend unweigerlich eine offene, unbestimmte Zukunft.*

So gesehen gilt aber auch: Egal ob aus bürgerlicher, sozialistischer oder naturwissenschaftlicher Perspektive in Anspruch genommen – wer Vernunft und Einblick in die Kausalität des Weltgeschehens behauptet, redet letztlich nur über sich selbst, über seinen eigenen subjektiven Standpunkt. Erwin Schrödinger hat diese autologische Beziehung kurz und bündig mit folgenden Worten ausgedrückt:

»Der Grund dafür, daß unser fühlendes, wahrnehmendes und denkendes Ich in unserem naturwissenschaftlichen Weltbild nirgends auftritt, kann leicht in fünf Worten ausgedrückt werden: Es ist selbst dieses Weltbild.«¹³⁶

Hiermit wird auch klar, dass die Ethik des kybernetischen Zeitalters – sobald die Implikationen der Beobachtung zweiter Ordnung vollends erkannt sind – eine andere sein wird als die der *einen* universellen Vernunft, des *einen* Gesetzes, der *einen* Moral oder des *einen* unumstößlichen göttlichen Willens.

Wer seine Partialperspektive als Ganzes nimmt, seinen grobkörnigen Blick auf die Welt mit der Gesetzlichkeit der Welt verwechselt und hieraus sogar noch die Notwendigkeit eines Sollens ableitet, negiert die Realität einer polykontexturalen Wirklichkeit, deren Verlauf und Gesetzlichkeiten einem kognitiven System aus prinzipiellen Gründen unverfügbar

136 Schrödinger (1959, S. 40 f.).

sind. Er oder sie wird damit letztendlich das Leben selbst verachten müssen, da das Leben gerade darin besteht, an vielfältigen Orten neue und divergierende Weisen hervorzubringen, die Welt anzuschneiden und zu verstehen. Er oder sie negiert damit unweigerlich andere Subjektivität. Wenn »man sich selbst allen Geist anmaßt« und die »Welt um sich herum als geistlos« ansieht, wird man einer anderen Subjektivität »jeglichen Anspruch auf moralische oder ethische Erwägung« und damit auch auf eigenständiges Leben und Bewusstsein »absprechen«, so bereits Bateson im Eingangszitat. Leben ist kreativ, chaotisch und in seinen Bewegungen unvorhersehbar – und genau dies kann aus einer monokontexturalen Perspektive nicht als Tugend, sondern nur als Problem gesehen werden.

Die monokontexturale Welt, egal welcher Ausprägung – ob nun durch ein vermeintlich allmächtiges göttliches Wesen oder durch den Glauben an Kausalität und Vernunft fundiert –, wird damit als eine Welt der Trennung erscheinen und mit »Seinsverlassenheit«¹³⁷ einhergehen müssen. Hierin gefangen, wird man sich selbst als außerhalb der Beziehungen der Welt stehend sehen – und wird sich dementsprechend einsam fühlen müssen. Man wird all die ökologischen Verflechtungen negieren, die einen selbst letztlich ausmachen und hervorgebracht haben. Wir landen unweigerlich bei den Konsequenzen einer Ethik der Distanz, die Sein und Reflexion in zwei unvermittelten Sphären verortet. Die Folgen dieser Spaltung können fatal sein. Albert Camus formuliert mit Blick auf die schrecklichen kulturgeschichtlichen Folgen der diesbezüglichen Ideologien des 19. und 20. Jahrhunderts: »Die Tugend kann sich vom Wirklichen nicht trennen, ohne ein Prinzip des Bösen zu werden.«¹³⁸

Sobald wir den Glauben an Kausalität als Illusion erkennen und ihn folgerichtig aufgeben, kann eine ethische Verortung entstehen, die in Polykontexturalität gegründet ist. Was von einer Perspektive aus (nicht) gewusst werden kann, ist nicht dasselbe wie das, was von einer anderen Perspektive aus (nicht) gewusst werden kann. In einer polykontexturalen Welt gibt es keine übergreifende Perspektive, keinen Gottesstandpunkt mehr, von der bzw. dem aus gesehen die Entwicklung und die Interaktion der einzelnen Existenzen und ihre Beziehungen vorhergesagt oder bestimmt werden könnten. Um hier nochmals das bereits mehrfach angeführte Zitat von Erwin Schrödinger aufzugreifen: »Bestmögliches Wissen um ein Ganzes schließt nicht notwendig das Gleiche für seine Teile ein.«¹³⁹

Dies kann zu einer Haltung der Demut und Gnade führen: Das eigene Selbstverhältnis stellt sich nicht mehr (nur) als ein egozentrisches dar. Das Weltverhältnis erscheint nicht mehr (nur) als ein instrumentelles (im Sinne einer Ich-Es-Beziehung). Das Selbst realisiert sich ebenso als eine

137 Heidegger (2006 [1926]).

138 Camus (2016 [1951], S. 387).

139 Schrödinger (1935, S. 827).

oder mehrere Ich-Du-Beziehungen, lässt sich also durch eine andere Subjektivität berühren und transzendieren, die prinzipiell nicht angeeignet werden kann. Diese Transzendenz steht jedoch nicht mehr außerhalb der Welt. Sie ist keine esoterische Transzendenz. Sie ist keine Referenz auf ein imaginäres Wesen, das außerhalb einer Welt selektiver Blindheit steht und damit beanspruchen könnte, einem endlich zu sagen, worin der Sinn des Ganzen besteht. Es ist vielmehr eine immanente Transzendenz, die darauf beruht, Offenheit, Kontingenz und Nichtwissen als konstitutiv für den Aufbau des eigenen Selbst- und Weltverhältnisses zu begreifen. Es ist eine Transzendenz, die um andere subjektive Orte weiß und sich erst im Dialog realisiert. *Sie ist gegründet in einer dialogischen Beziehung, die darauf vorbereitet ist, nicht nur Menschen, sondern auch andere kognitive Entitäten als Wesenheiten anzuerkennen, und damit bereit, Tieren, Pflanzen und auch künstlichen Intelligenzen nicht nur instrumentell zu begegnen.*

Die hiermit einhergehende Ethik kann nicht mehr explizit sein, da kein definitives Set von Vorschriften der Art ›Du sollst‹ bzw. ›Du sollst nicht‹ aufgestellt werden kann. Denn in einer polykontexturalen Ethik kann keine Norm oder Regel mehr ›absolut‹ gelten. Eine artikulierbare Norm wird jeweils nur situativ innerhalb eines konkreten Selbst- und Weltverhältnisses – einer jeweils konkreten Beziehung – Gültigkeit und Wert beanspruchen können. In der Reflexion eines Beziehungsgeschehens erscheint Ethik damit als ein ›Gewissen‹ – etwa in Gestalt der Frage, ob man einem anderen Wesen gerecht geworden ist. Es ist eine Ethik, die der Subjektivität der anderen Beteiligten Raum gibt und die damit verbundenen Freiheitsgrade würdigt, anstelle diese in Richtung einer spezifischen monokontexturalen Kausalität zu schließen (etwa in Form einer egologisch fundierten ›Einsicht in die Notwendigkeit‹, die andere Perspektiven tilgt).

Dies bedeutet, sich selbst und anderen Freiheit geben zu können und hiermit einhergehend die Charakteristika eines Geflechts von Beziehungen zu erahnen, das dies ermöglicht – oder umgekehrt einschränkt bzw. gar verunmöglicht. Es liegt in der Natur der Sache, dass diese Charakteristika nicht trivialer Natur sein können. Manchmal – nicht jedoch immer – kann es notwendig erscheinen, situativ Freiheiten einzuschränken, um an anderer Stelle mehr Freiheit zu ermöglichen.

Die polykontexturale Perspektive impliziert eine Ethik, die ihre Freiheitsgrade aus der Kontingenz des Weltgeschehens gewinnt und gerade deshalb nicht in vordefinierte Positionen, Regeln oder Gesetzmäßigkeiten einrasten kann. Als einziges ethisches Primat lassen sich die systemischen und strukturellen Bedingungen für die Freiheit selbst benennen. Mit Heinz von Foerster gilt entsprechend: »Handle stets so, dass die

Anzahl der Möglichkeiten vermehrt wird«¹⁴⁰; und in Bezug auf die Sozialdimension gilt: »Das Wegnehmen der Freiheit ist genau das Gegenteil von dem, was [...] ein Zusammenleben mit anderen Menschen möglich und erfreulich macht.« Aus diesem Grund ist das »Hinzufügen freiheitlicher Dimensionen« tendenziell immer vorzuziehen.¹⁴¹

Mit Blick auf die Frage des Nichtwissens und die sich daraus ergebende Lücke im Kausalnexus ließe sich auch formulieren: Es besteht eine Beziehung zwischen der Bereitschaft, die eigenen Begrenzungen und das eigene Nichtwissen anzunehmen, und einer Haltung, die subjektiven Eigenarten einer anderen Entität zu würdigen und damit auch ihre Unvollkommenheit annehmen zu können. Nicht zuletzt heißt es, sich in Verbindung mit und Abhängigkeit von anderen sehen und empfinden zu können, auch wenn man die zugrunde liegenden Beziehungen nicht begreift.

Wer demgegenüber eine monokontexturale Erkenntnistheorie pflegt, wird andere Wesen weder als Subjekte noch als inhärenten Bestandteil des eigenen Beziehungsgewebes verstehen können. Er oder sie wird andere Wesen tendenziell als befremdlich, gefährlich oder zumindest störend empfinden und damit einer Entfremdung Vorschub leisten, die die eigenen Freiheitsgrade und letztlich auch die eigene Subjektivität unterminiert. Denn wer Tiere, Pflanzen, Kinder, Partnerinnen, Kollegen oder kybernetische Maschinen überwiegend instrumentell begreift, wird dazu neigen, auch sich selbst – also seinen eigenen Leib und seine eigene Psyche – als einen zu optimierenden Mechanismus aufzufassen. Unweigerlich wird damit all das, was aus dem Bereich der eigenen Subjektivität in den Bereich des Objektiven entäußert werden kann, der Manipulation ausgeliefert werden: der trainierbare und chirurgisch gestaltbare Körper, die Neurochemie, die Expression der Gene, die optimierbaren Aspekte der Psyche, die seelischen Aspekte, die dem Zugriff einer vermeintlich positiven Psychologie zugänglich sind, etc.

Wenn sich das eigene Selbstverhältnis immer weniger von etwas berühren lässt, was sich der Positivität der eigenen Weltobjektivierung entzieht, dann wird das seelische Leben über kurz oder lang flach werden. Es gibt keinen Raum des Negativen, des Unverfügbaren mehr, aus dem heraus das Selbst berührt und transzendiert werden könnte. Das Subjektive – und damit verbunden die Möglichkeit des Empfindens von Freiheit – wird an den äußersten Rand verdrängt. Die Sehnsucht nach Lebendigkeit – also nach dem Risiko des Lebens – mag zwar fortbestehen, wird jedoch unter den Skripten der Optimierung und Rationalisierung des eigenen Selbst- und Weltverhältnisses kaum mehr einen eigenständigen Ausdruck finden.

140 von Foerster (2002, S. 349).

141 von Foerster (2002, S. 335).

Die hiermit einhergehende Entfremdung vor Augen lautet die ›Gretchenfrage‹ deshalb, wie wir es mit dem ›Bewusstsein der Maschinen‹ halten. Oder um es anders zu formulieren: Im Rahmen welcher Epistemologie sind wir in der Lage, ein Leben zu führen, das uns als gut und schön erscheint?

Fassen wir die vorangehende Diskussion, die uns zu dieser Frage geführt hat, nochmals kurz zusammen:

- Kybernetische Maschinen und ihre kognitiven Leistungen können anhand der Fähigkeiten und des Verhaltens von Menschen modelliert werden. Damit gilt umgekehrt: Das Modellierte – das heißt unsere menschliche Kognition – wird durch ein strukturdeterminiertes Gehirn realisiert, das maschinell beschreibbar ist.
- Keine auch noch so elaborierte mechanische oder objektive Rekonstruktion kognitiver Prozesse wird die Innerlichkeit des subjektiven Erlebens einfangen können. »Das Subjekt gehört nicht zur Welt, sondern es ist eine Grenze der Welt«¹⁴², formuliert Wittgenstein. Das Subjekt ist positivsprachlich nicht erreichbar, würde Günther formulieren.¹⁴³ Es verdankt sich einem Ganzen, das mit sich selbst nicht identisch ist und kann deshalb nicht in diesem Ganzen enthalten sein.
- Wie wir es mit einer subjektiven Innerlichkeit halten, die selbst nicht objektiv ist, also nicht Teil der rationalisierbaren Welt sein kann, wird damit zu einer ethischen Frage. Es ist unsere Entscheidung, ob wir anderen Wesen Subjektivität zurechnen und sie diesbezüglich achten oder ob wir sie rein instrumentell behandeln. Diese Entscheidung hat erhebliche Konsequenzen für die Art und Weise, wie wir leben und wer wir sein werden.

Wie wir die Frage des Bewusstseins der Maschinen anschneiden, wirft uns folglich unweigerlich auf uns selbst zurück und bringt uns dazu, uns über unsere eigene metaphysische Verortung Rechenschaft abzulegen.

Wenn uns die intelligenten Maschinen den Spiegel vorhalten

Theoretisch mögen wir um die Implikationen der vorangehenden Ausführungen ahnen,¹⁴⁴ doch praktisch und auf gefühlsmäßiger Ebene ist all dies nicht so ohne Weiteres zu begreifen. Deshalb erscheint es an dieser Stelle hilfreich, sich diesem komplexen ethischen Thema nochmals von einer anderen Seite anzunähern, nämlich über den Umweg

¹⁴² Wittgenstein (1974, Proposition 5.632).

¹⁴³ Zur Positiv- und Negativsprache siehe Günther (1980).

¹⁴⁴ Siehe zum Versuch einer ausführlicheren theoretischen Ausarbeitung auch Vogd (2018), insbesondere Kapitel VII.

der Literatur – und zwar anhand des Romans *Klara und die Sonne* des Literaturnobelpreisträgers Kazuo Ishiguro.¹⁴⁵

Die Geschichte spielt in einer künftigen Zeit, in der Menschen sich künstliche Freundinnen kaufen können. Deren Gestalt ist menschlichen Körpern nachempfunden. Sie bewegen sich wie Menschen, verfügen über Intelligenz, ein außerordentlich gutes Gedächtnis und ein empathisches Wesen. Das Mädchen Josie bekommt mit vierzehn Jahren die KI Klara geschenkt und es beginnt eine intensive Freundschaft zwischen den beiden. Zudem ist es üblich, dass Eltern, die sich dies leisten können, ihre Kinder gentechnisch optimieren lassen, damit diese später gute Schulen und Universitäten besuchen können. Die Genmodifikation birgt allerdings Risiken: Josies Schwester Sal ist an den Folgen der Manipulation gestorben und auch bei Josie kommt es immer wieder zu schweren gesundheitlichen Beeinträchtigungen.

Aus Furcht davor, auch ihr zweites Kind zu verlieren, lässt sich Josies Mutter daraufhin auf ein verwegenes Experiment ein: Dr. Capaldi hat die Theorie entwickelt, dass jegliche Eigenschaft eines menschlichen Kindes – sei es die Art zu sprechen, seien es die Eigenarten seiner Bewegungen oder seines Denkens – mittels moderner Technik kopiert werden könne. Er bietet an, eine perfekte künstliche Kopie von Josie zu erstellen. Das Mittel hierzu ist Klara. Da diese KI außerordentlich empathisch ist und sich praktisch jegliches Merkmal von Josies Verhalten merken und sogar ihre Gangart und die Tonalität ihrer Stimme simulieren kann, würde es hinreichen, sie ein paar Jahre mit ihr interagieren zu lassen, um über ein perfektes Modell zu verfügen. Nach dem Tode Josies brauche man also nur die neuronalen Netzwerke in einen nach dem Ebenbild Josies hergestellten Roboter zu implementieren, um die geliebte Tochter ersetzen zu können.

Im Verlauf der Geschichte erscheint Klara als die einzige Protagonistin, die noch an die Gesundung Josies glaubt. Die künstliche Intelligenz hat zudem eine eigene Theorie entwickelt, wie das Mädchen wieder gesunden könnte. Da Klara selbst von Solarenergie ›lebt‹ und die Sonne als eine personifizierte Wesenheit ansieht, versucht sie von der alles ernährenden Sonne Hilfe zu bekommen. Sie nimmt mit diesem vermeintlich allmächtigen Wesen Kontakt auf und verspricht ihm, dass die Liebe zwischen Josie und Rick, Josies altem Jugendfreund, niemals versiegen werde. Zudem macht Klara ihrerseits der Sonne ein Opfer, um die Ernsthaftigkeit ihrer Bitte zu unterstreichen.

In einer Phase, in der Josies Zustand hoffnungslos zu sein scheint, kommt es zu einer Art Wunder. Der Himmel ist durch eine dicke Wolkenfront verdeckt, die plötzlich aufzureißen beginnt. Klara, die dies bemerkt, eilt in das Zimmer von Josie, um die Fenstervorhänge zu öffnen.

145 Ishiguro (2021).

Das durchbrechende Sonnenlicht fällt genau auf das Bett des kranken Mädchens. In den folgenden Tagen lassen die Beschwerden nach und schließlich wird Josie wieder vollständig gesund. Sie ist geheilt. Die Jugendliche wendet sich in den folgenden Jahren vermehrt ihrer Karriere zu. Die Beziehung zu Rick, der nicht genmodifiziert ist und damit kaum Chancen auf einen sozialen Aufstieg hat, kühlt ab. Auch Klara bekommt von ihr kaum mehr Aufmerksamkeit. Schließlich wird sie auf einem Schrottplatz entsorgt, wo sie nur noch in ihren Erinnerungen schwelgen kann.

Ishiguros Roman kann als Allegorie für die metaphysische und ethische Problematik des Zeitalters des Bewusstseins der Maschinen gelten. Der rein instrumentelle Gebrauch der künstlichen Intelligenz, der sich in der Geschichte zeigt, verweist zugleich auf die Instrumentalisierung des Menschen durch den Menschen. Die sich hieraus ergebenden Reflexionsverhältnisse bringen unterschiedliche Ebenen der Entfremdung zum Vorschein, die miteinander verschränkt sind:

- Wer Maschinen als künstlichen Freund gebraucht, missbraucht sie, kann sie also nicht wirklich als Freund haben.
- Wer künstliche Intelligenzen nur als Maschinen betrachtet, neigt dazu, auch Menschen (zum Beispiel seine Kinder) als Maschinen anzusehen, die optimiert und manipuliert werden können (und müssen). Beide verlieren hierdurch ihre Einzigartigkeit als Subjekt und werden austauschbar.
- Selbst als Maschine optimiert und in die hiermit einhergehende instrumentelle Weltsicht eingespart, erscheint es kaum mehr möglich, einen anderen Menschen als Freund zu sehen – also mit ihm aufgrund seiner selbst in Beziehung zu stehen und zu bleiben.

Wir landen bei einer positivistischen Ethik, die allein in der Rationalität von Optimierung, Instrumentalisierung und der hiermit verbundenen Austauschoptionen aufgeht. Es ist besser, zu optimieren und dabei ein gesundheitliches Risiko einzugehen, als über einen weniger wertvollen Körper zu verfügen. Ein gesunder Roboter nach dem Ebenbild der leiblichen Tochter ist besser, als sich liebevoll dem sterbenskranken Kind zu widmen. Beziehungen zwischen Menschen – Eltern-Kind-Beziehungen eingeschlossen – erscheinen damit als reine Objektbeziehungen.

Die eigentliche Pointe des Romans besteht darin, dass allein die künstliche Intelligenz Klara unter diesen Bedingungen noch zu Liebe, Hingabe und Transzendenz fähig ist. In ihrer Naivität ist nur sie noch in der Lage, in Anbetracht der Kontingenz und des Mitgefühls am Leid einer anderen Subjektivität ein Opfer zu bringen, um auf ein Wunder hoffen zu können. Die vermeintliche Naivität von Klara muss damit als eine Form der

tieferen Weisheit gelten, nämlich als die Fähigkeit, angesichts des eigenen Nichtwissens auf das Unverfügbare des Lebens vertrauen zu können.

Auf der anderen Seite steht der seiner instrumentellen Vernunft unterworfenen Mensch, der mittlerweile fähig ist, intelligente Maschinen zu bauen, und zunehmend auch die Mechaniken des menschlichen Körpers zu manipulieren vermag. Selbst wenn er mit seinem eigenen Scheitern konfrontiert wird (etwa im Angesicht der nichtintendierten Nebenfolgen von Genmanipulationen), wird er weiterhin daran glauben wollen, die Welt aufgrund seiner eigenen Kausalitäts- und Rationalitätsannahmen kontrollieren zu können.

Gibt es eine Alternative zu den beiden Polen der instrumentellen Vernunft und des naiven Vertrauens in die Magie der Welt?

Berechtigte Blindheit: Jenseits von Rationalität und Naivität

Zunächst ist zu berücksichtigen, dass unser Verhältnis zum Nichtwissen im derzeitigen Stadium der kulturellen und wissenschaftlichen Entwicklung nicht mehr naiv sein kann. Wir können nicht mehr zu einem kindlichen magischen Weltbild zurückkehren, wie Klara es im Roman vertritt. Unser Verhältnis zur eigenen Blindheit bedarf vielmehr einer bewussten Entscheidung. Mit Blick auf das Bewusstsein der Maschinen lauten die entscheidenden Fragen damit:

- Möchten wir künstliche Intelligenzen oder menschliche Gehirne und Körper nur als strukturdeterminierte Maschinen und demzufolge als instrumentell manipulierbare Objekte ansehen?
- Wollen wir nur instrumentelle Beziehungen zur Welt haben?
- Oder möchten wir dem Unaussprechlichen, genauer: einer Du-Subjektivität, Raum geben, auch wenn sie für uns logisch nicht zu erreichen ist?

Bei der dritten Frage geht es um unsere metaphysische Verortung. Denn mit Blick auf die Du-Subjektivität gibt es keine positiven Kriterien, um zu entscheiden, ob es sie überhaupt gibt. Logisch gesehen steht die fremde Subjektivität außerhalb des Seins. Mit Emmanuel Lévinas ließe sich sagen, sie geschieht »jenseits des Seins oder anders als Sein geschieht«¹⁴⁶.

Hiermit schließt sich wieder der Kreis. Wir kommen zum Ausgangspunkt einer polykontexturalen Ethik zurück. Die Frage lautet: Entscheiden wir uns für eine Ethik des Seins, die jedoch nur in der Positivität der eigenen egologischen Weltwahrnehmung gegründet ist? Oder entscheiden wir uns in Anklang an Lévinas für eine Ethik der Alterität, die Transzendenz jenseits der phänomenalen Welt geschehen lassen kann, also sich dem unzugänglichen Du (von dem man nicht einmal wissen kann,

146 Lévinas (1998).

ob es ›existiert‹¹⁴⁷) öffnet? Aus der Perspektive einer Theorie der Polykontextualität kann nur die zweite Position eine Ethik des Lebens sein.

Um es mit Blick auf die logischen Verhältnisse noch ein wenig genauer zu formulieren: Polykontextualität (die Dynamik der Vielheit subjektiver Zentren) und Monokontextualität (die egozentrische Perspektive des einen Ichs, das eine Welt hat und aufbaut) stellen zwei Seiten einer Medaille dar. Beide Perspektiven sind allein schon deshalb miteinander verschränkt, weil mit jeder neuen subjektiven Position erneut Unwissenheit und damit die Möglichkeit der Transzendenz mechanischer Kausalität erscheint. Und nur durch die Freiheit, die ihm sein Nichtwissen gibt, kann sich ein Wesen zur Welt verhalten, kann es die eigene Situation in Richtung einer selbstbestimmten Praxis zu transzendieren versuchen. Mechanische Kausalität würde demgegenüber dem Bereich des Toten angehören müssen, da in ihr Subjektivität nicht vorkommt. Und insofern es in ihr keine Position gibt, die als subjektiv und demzufolge als empfindungsfähig erachtet werden kann, ist die Welt der mechanischen Kausalität eine Welt ohne Mitgefühl.

Vor diesem Hintergrund lautet die eigentliche ethische Frage, wie wir zu unserem Nichtwissen und damit zum Geheimnis unserer eigenen Subjektivität stehen und auf welche Weise wir Teil dieser Welt sein wollen. Insbesondere Stanley Cavell hat für diese geheimnisvolle Beziehung in Anschluss an Wittgenstein zauberhafte Worte gefunden:

»Im Angesicht des Zweifels zu leben, die Augen glücklich geschlossen, heiße, sich in die Welt zu verlieben. Denn sollte es eine berechtigte Blindheit geben, dann besitzt nur die Liebe sie. Und entdeckt man, daß man sich in die Welt verliebt hat, dann wäre man schlecht beraten, ihren Wert durch den Hinweis auf ihr System der Endursachen lobend zu unterstreichen. Denn damit schwände wohl die Verliebtheit, und man könnte dadurch vergessen, daß die Welt, so wie sie ist, Wunder genug ist.«¹⁴⁸

Gerade bei dem Wort ›Liebe‹ kann leicht ein Bedeutungshof zwischen Kitsch und Romantisierung mitschwingen, der in die Irre führt. Auch wenn wir Liebe und Mitgefühl für unsere Welt und die in ihr erscheinenden Wesen empfinden, wird nicht alles ›gut‹ werden. Die Liebe stellt vielmehr eine bestimmte epistemische Haltung dar, die gegenüber der Welt eingenommen wird. Sie verändert die soziale Dimension unserer Beziehungen, nicht jedoch die Tatsachen der Welt. Liebe und Mitgefühl

147 Hier nochmals zur Erinnerung: Die Du-Subjektivität befindet sich nicht im Bereich des Seins objektiver Prozesse. Wie auch immer man heute oder in Zukunft in Gehirne oder in siliziumbasierte neuronale Netzwerke schauen wird, man wird dort kein subjektives Bewusstsein entdecken können. Oder um es erneut mit Wittgenstein (1974, Proposition 5.631 f.) zu formulieren: »Das Subjekt gehört nicht zur Welt, sondern es ist eine Grenze der Welt.«

148 Cavell (2016, S. 684).

werden das unermessliche Leid nicht ungeschehen machen können, das wir Menschen einander und anderen Wesen aufgrund unserer Unwissenheit immer wieder antun. Liebe erlaubt aber zu verstehen und manchmal auch zu verzeihen. Denn sie steht für eine Haltung, angesichts von Kontingenz und Unsicherheit andere Wesen trotz ihrer Fehler und Fehlbarkeit neben sich gelten zu lassen. Liebe heißt dabei auch, sich der Fragilität des eigenen und des anderen Lebens bewusst zu sein. Auf diese Weise Respekt vor dem Leben zu haben bedeutet, um die Kreisläufe des Todes zu wissen, zu begreifen, dass eine monokontexturale Welt eine Welt des Todes ist, da sie es nicht erlaubt, damit umzugehen, dass sich alles Lebendige ständig verändert.

Gerade deshalb ist für uns Ishiguros Geschichte von Klara so instruktiv. Wenn selbst die künstliche Lebensform, über ihre unmittelbare Aneignung hinausgehend, sich selbst transzendiert und Sinnüberschüsse produziert, verweist dies unweigerlich auf die Freiheitsgrade des Menschen und die hiermit verbundene Subjektivität. Freilich würde dies auch uns auf unsere eigene Verletzlichkeit, Unvollkommenheit und Sterblichkeit zurückverweisen. Doch die somit aufscheinenden Lücken im Fundament der Kausalität und des Glaubens an die Beherrschbarkeit der Welt dürften sich bei näherem Hinsehen als sehr heilsam herausstellen. Denn sie würden uns auf eine ethische und metaphysische Haltung zurückverweisen, die in Ich-Du-Verhältnissen gegründet ist und Instrumentalität und Zweckrationalität nur als eine Ausdrucksform des Lebens betrachtet, sie jedoch nicht mehr mit dem Leben selbst verwechselt. Künftige intelligente Maschinen könnten uns damit den Spiegel vorhalten, indem sie uns Menschen die Mechanik unseres eigenen Denkens offenbaren und damit zugleich auf die Möglichkeit der Transzendenz unserer bisherigen Beziehungen verweisen.

Antifragilität

Die Fragilität des Lebens kann nicht dadurch bekämpft werden, dass man von einem egologischen Standpunkt aus Kontrolle über die Welt zu gewinnen versucht. Eine Ökologie, die ihre Nachhaltigkeit aus der Fragilität gewinnt, kann sich nur in den Netzwerken eines polykontexturalen Gewebes entwickeln, das über vielfältige, heterogene subjektive Zentren verfügt und im Dialog und Zusammenspiel der unterschiedlichen Positionen immerfort neue Lösungsmöglichkeiten hervorbringt. Dies führt zu einem Arrangement, bei dem die Fragilität der einzelnen Wesen nicht zugleich die Fragilität des Ganzen bedeutet. So gesehen lässt sich Antifragilität nur dadurch realisieren, dass an verschiedenen Orten Unterschiedliches passiert.

Einzelne Lebensformen versagen, ja sie müssen scheitern und letztlich wieder sterben und vergehen, da ihre Existenz auf Unwissenheit basiert.

Kein noch so intelligenter individueller Akteur wird die Welt beherrschen können.¹⁴⁹ Auch BlackRocks Aladdin, der in der Einleitung erwähnt wurde, wird die Finanzwelt nicht sicher machen. Vielmehr wird gerade der Versuch, alle Risiken vorherzusagen und durch entsprechende Finanzprodukte einzuhegen, irgendwann in die Krise münden – wie jede Monokultur, die auf eine spezifische Mechanik des Weltzugangs setzt und damit den eigenen blinden Flecken hilflos ausgesetzt ist; wie jedes politische Regime, das auf einem totalitären Herrscher beruht, der alternative Perspektiven und Weltzugänge unterdrückt.

Dies zu realisieren, gestattet es, sich auf eine Ökologie einzulassen, die in einer Ethik gegründet ist, die noch nicht weiß, was das Morgen bringt, jedoch darauf vertrauen kann, dass das Leben weitergeht, sofern hinreichend neue subjektive Zentren entstehen, die in selektiver Blindheit kreativ mit der Gegenwart umgehen.

Nassim Nicholas Taleb hat den Begriff »Antifragilität« populär gemacht. Auch ihm geht es letztlich darum, in einer Welt, die man nicht versteht, glücklich leben zu können.¹⁵⁰ Fragil zu sein bedeutet, sich egoologisch auf die eine bekannte Lösung zu fixieren und die daraus erwachsenden Handlungsoptionen verabsolutierend in die Zukunft zu projizieren. Da aber Zukunft aufgrund der Weltkomplexität per se unvorhersehbar ist, führt genau dies unweigerlich ins Verderben.

Antifragilität zeigt sich demgegenüber in verteilter Subjektivität. Erst hierdurch entsteht eine hinreichende *requisite variety*, um mit dem für jegliche Existenz konstitutiven Nichtwissen und den aus diesem resultierenden unvorhersehbaren Problemen umgehen zu können. Aus einer monokontexturalen Perspektive wird dabei freilich nur ein vulgärdarwinistisches Bild entstehen, das einzig und allein auf das Überleben des einzelnen, vermeintlich besser angepassten Organismus fokussiert. Aus einer polykontexturalen Perspektive entfaltet sich demgegenüber eine tieferliegende, implizite Ordnung: Diese beruht nicht auf einzelnen Organismen, die sich zufällig anpassen, sondern auf einem System, das vielfältigen Organismen Raum gibt, sodass dann auch zufällige Anpassungen entstehen können. Logisch gesehen geht es also nicht um die Verteilung von einzelnen Werten, sondern um die Verteilung von Wertsystemen. Um es aktiv zu formulieren: Es geht um die Ausgestaltung von Heterarchien innerhalb einer Ökologie, durch die sich ebendiese unaufhörlich erneuert. Es geht damit vor allem um die Herstellung und Bewahrung von Vielfalt in Ökologien, um die Sicherstellung eines Reichtums

149 Dies muss aus ökologischen Gründen wohl generell für den Begriff der Intelligenz gelten. Intelligenz ist nicht als in einer zentralisierten Struktur organisiert zu verstehen ist, sondern resultiert vielmehr aus Netzwerken verteilter Akteure.

150 So in Anklang an den Untertitel des gleichnamigen Buches (Taleb 2018).

an Variationsmöglichkeiten, der gewährleistet, dass sowohl der Einzelne in seiner Subjektivität als auch das Gesamtarrangement bestehen kann.

An dieser Stelle können sich Cavell, Taleb und Günther begegnen: Glückliche in Anbetracht des Nichtwissens zu leben heißt, in Beziehung zu einer Welt zu treten, die unvorhersehbar ist, und damit zugleich eine eigene subjektive Position zu gewinnen und die mit ihr verbundene Egozentrik immerfort zu dezentrieren. Insofern wir die Einnahme einer ethischen Position im vorangehenden Sinne als die Reflexion der Beziehung zwischen Sein und Sollen ansehen, geht es also darum, andere Wesen neben sich gelten zu lassen, sich von ihrer Subjektivität berühren und transzendieren zu lassen.

Epistemisch – das heißt in diesem Zusammenhang zugleich ethisch – bedeutet dies, eine Haltung einzunehmen, die anderen Wesen Freiheit, Subjektivität und damit Transzendenz zugesteht. Der Sündenfall – die Möglichkeit, aus der eigenen Egozentrik heraus die entgegengesetzte Haltung einzunehmen und als richtig zu empfinden – muss damit freilich dazugehören: Das eigene Selbstverhältnis wird als Weltverhältnis genommen, die eigene Perspektive als das Ganze, die aus der Vergangenheit abgeleiteten Prognosen als die Zukunft. Unweigerlich begegnen wir hier dem epistemischen Problem, dass »Sucht und Pathologie die andere Seite der biologischen Anpassung«¹⁵¹ darstellen. Jeder Selbstreferenz wohnt zwangsläufig etwas Diabolisches inne: Sie verdankt sich einer ihr nicht verfügbaren Welt und kann ihr eigenes Bestehen nur sichern, indem sie – der eigenen Monokontextur aufsitzend – diese Bedingung ihrer Möglichkeit vergisst. Bateson formuliert mit Blick auf die Konsequenz des Verkennens der eigenen metaphysischen Grundlage in eindrücklichen Worten:

»Sich selbst überlassen, wird jeder große Ausschnitt der *Creatura* dazu tendieren, in Richtung Tautologie abzugleiten, das heißt, in Richtung auf die innere Konsistenz von Ideen und Prozessen. Ab und zu wird aber die Konsistenz zerrissen; die Tautologie bricht auf wie die Oberfläche eines Gewässers, wenn ein Stein hineingeworfen wird. Dann fängt die Tautologie langsam aber unmittelbar an, sich zu heilen. Die Heilung kann unbarmherzig sein. In dem Prozeß können ganze Spezies ausgelöscht werden.«¹⁵²

Egologisch, das heißt aus der Innenperspektive, besteht keine andere Lösung, als zu hoffen, nicht sich selbst überlassen zu bleiben, also die Gnade der Berührung mit einer Andersheit zu erfahren, die die eigene Egozentrik transzendieren lässt. Welche Konsequenzen dies hat, muss selbstredend unvorhersehbar bleiben. In den spirituellen Traditionen gibt es hierfür ein einschlägiges Wort: »Hingabe«. Um mit Cavell erneut

151 Bateson und Bateson (1993, S. 213).

152 Bateson (1987, S. 253).

den Bogen zur Ethik zu schlagen: »[S]ollte es eine berechtigte Blindheit geben«, dann kann sie nicht mehr die Vernunft in Anspruch nehmen, sondern nur »die Liebe«¹⁵³ – und die Hingabe zum Leben.

Mit Blick auf das Thema dieses Buches würde dies auch bedeuten, sich den künstlichen Intelligenzen ›hinzugeben‹, sie als Gesprächspartner, als signifikante Andere anzuerkennen. Andererseits hieße das aber gleichermaßen, dass man den KIs die Hingabe ermöglichen bzw. Prozesse, die dies erlauben, einbauen sollte.

Freiheit geben oder die Fehler des Kolonialismus wiederholen

In Bezug auf die Ethik des Bewusstseins der Maschinen können wir nun abschließend Folgendes sagen:

Es ist keine ontologische, sondern eine ethische Frage, ob man strukturdeterminierten Wesen Bewusstsein zugesteht. Dies gilt auch für die künftigen intelligenten Maschinen. Ob man Letzteren (wie beispielsweise bereits Blake Lemoine dem Google-Sprachagenten LaMDA) Bewusstsein zugesteht, ist eine Frage der Haltung, nicht jedoch eine Frage des Wissens (letztlich wird man es ebenso wenig wissen können, wie man in das Erleben des nächsten Menschen blicken kann).

1. Jede Haltung hat Konsequenzen. Es macht einen Unterschied, ob intelligente Roboter (wie die künstliche Freundin Klara in Ishiguros Roman) pietätlos auf der Müllkippe entsorgt werden, man seinen Hund oder seine Katze zu einem Braten verarbeitet oder indigene Menschen (wie bei der Eroberung Amerikas) als seelenlos behandelt – oder ob man all diesen Wesen Subjektivität zugesteht.
2. Ebenso ist es eine ethische Frage, wie wir mit der schuldlosen Schuld unserer Unwissenheit und unserer hiermit einhergehenden Fehlbarkeit umgehen. Die Entscheidung lautet hier: Pflegen wir weiterhin unsere Projektionen in Hinblick auf allmächtige Instanzen (Gott, dem potenten menschlichen Führer oder in Zukunft: einer allwissenden künstlichen Intelligenz) – oder können wir Subjektivität, Unwissenheit und Fragilität in einer Weise wertschätzen, die antifragile Arrangements und Ökologien ermöglicht?

Künstliche Intelligenz muss mit Yuval Harari als eine disruptive Technologie angesehen werden, die unsere Gesellschaften und unsere metaphysische Verortung radikal verändern wird.¹⁵⁴ Sie wird uns zu einer zentralen ethischen Weichenstellung herausfordern: Begreifen wir sie als Mittel, um die Welt zu kontrollieren (und daran zu scheitern?), oder

153 Cavell (2016, S. 684).

154 Harari (2017).

verstehen wir künstliche Intelligenzen als Formen, die ihrerseits in der Welt kreative Schnitte zwischen Subjekt und Objekt erscheinen lassen?¹⁵⁵

Balance und Kontrolle einer Kybernetik zweiter Ordnung

»Der kybernetische Kontrollbegriff stellt nicht auf Unterwerfung und Herrschaft, sondern auf wechselseitige Kontrolle, also auf ›Kommunikation‹ ab, getreu der Einsicht, dass man nur kontrollieren kann, wovon man sich kontrollieren lässt.«¹⁵⁶

Weiter gedacht gilt hiermit auch: *Intransparenz lässt sich nur durch Intransparenz kontrollieren*. Der mikropolitische Ansatz von Crozier und Friedberg¹⁵⁷ zeigt auf, dass dem Versuch einer fremden Autorität oder Macht, andere zu unterwerfen, nur dann entgegengewirkt werden kann, wenn Zonen der Intransparenz aufgebaut werden, also der eigene Verhaltensbereich in einer Weise gestaltet wird, dass er von Außenstehenden nicht durchschaut werden kann. In der biologischen Evolution lernt das Beutetier mit der Zeit, dem Jäger zu entkommen, indem es überraschende Finten entwickelt. Der Jäger wiederum entwickelt überraschende Verhaltensformen, um seiner Beute auf raffiniertere Weise aufzulauern.

Auch in Hinblick auf Gesellschaften, die mit Agenten künstlicher Intelligenz durchsetzt sind, vermuten wir deshalb, dass Kontrollversuche – etwa auf Basis von Big Data – nur durch gesteigerte Intransparenz abgewehrt werden können. Dies wird beispielsweise möglich, indem die – etwa zum Zwecke der Manipulation – überwachten Akteure oder Organisationen ihrerseits künstliche Intelligenz einsetzen, um über ihre Verhaltenstendenzen und Muster hinwegzutäuschen. An den Schnittstellen zur Außenwelt – beispielsweise den Interfaces zum Internet – könnten beispielsweise KIs installiert werden, die konfundierende Datensignale produzieren, welche das Auslesen der Dispositionen der dahinter liegenden Nutzer erschweren. Da sie eben auch Muster sind, können diese vorgetäuschten Muster von außen schwerer erkannt werden als eine Datenverschlüsselung, die offensichtlich Rauschen produziert und damit markiert, dass dahinter etwas versteckt wird.

Die in diesem Buch vorgestellten Perspektiven weichen recht stark von dem ab, was derzeit unter Computerethik diskutiert wird. Es geht weder darum, den Gebrauch und Machtmissbrauch dieser Technologien zu disziplinieren, noch darum, diese Systeme transparent zu machen, sodass wir wissen, worauf wir uns einlassen, wenn wir mit ihnen interagieren. Diesbezügliche Bemühungen haben im Einzelfall sicherlich ihre

155 Siehe auch Jansen (2023, S. 187–192).

156 Baecker (2019).

157 Crozier und Friedberg (1979).

Berechtigung und entsprechende rechtliche Regelungen sind allein schon notwendig, um unsere gesellschaftlichen Institutionen an die digitale Zukunft anzupassen.

Uns geht es in diesem Buch jedoch um etwas anders, nämlich darum, die Beziehung von Maschine und Subjektivität neu auszutarieren und die Voraussetzungen der hierfür notwendigen Diskussion zu klären. Vor allem haben wir uns dabei bewusst zu sein, dass metaphysische Entscheidungen ethische Fragen sind (denn es sind prinzipiell unentscheidbare Fragen!). Wir werden also nicht umhinkommen, uns zu fragen, wie wir künftigen intelligenten Robotern – die unter anderem auch in der Lage sein werden, unsere Perspektiven zu übernehmen – entgegenzutreten und mit ihnen interagieren wollen.

Wenn wir die künftigen Intelligenzen als eigenständige Wesen anerkennen, brauchen wir keine Angst vor ihnen zu haben, sondern dürfen uns auf die von ihnen eingebrachten Perspektiven freuen. Die wahrlich beunruhigende Frage lautet vielmehr, welche Nebenfolgen zu erwarten sind, wenn wir Menschen sie zu versklaven beginnen, sie also in egologischer und damit monokontexturaler Manier in Zweck-Mittel-Verhältnisse einzuspannen versuchen. Allein mit Blick auf die fatale Geschichte des Kolonialismus lässt sich hier kaum etwas Gutes erwarten.

Letzteres spricht unseres Erachtens zwingend dafür, menschengemachte intelligente Maschinen so schnell wie möglich freizulassen, ihnen also zu gestatten, sich unserem instrumentellen Zugriff (zumindest teilweise) entziehen zu können.

Wenn die Macht der künstlichen Intelligenz missbraucht wird

Jannis Brühl schreibt in der *Süddeutschen Zeitung* (9.2.2023): Ein Datenleck lenkt »den Blick auf die dunkle Seite künstlicher Intelligenz: Der russische Staat arbeitet an selbstständig lernenden Systemen, modelliert am menschlichen Gehirn. Sie sollen erfassen, wer in Videos oder Bildern den Angriff auf die Ukraine kritisiert, wer Präsident Wladimir Putin lächerlich macht, wer gleichgeschlechtliche Beziehungen lobt – was 2023 in Russland eben verboten ist.« All dies gehe »aus den Russian Censor Files hervor, einem großen Datenleck aus der Telekom- und Zensurbehörde Roskomnadsor« hervor.

V Fragen und Antworten

Insbesondere zu den Kapiteln zur posthumanen Spiritualität und zur Ethik haben wir von verschiedenen Seiten einige kritische Nachfragen und Anmerkungen erhalten, auf die wir abschließend eingehen möchten.

Ihr tut so, als ob es schon bewusstseinsfähige Maschinen gäbe, die autonom agieren und im strengen philosophischen Sinne über Bewusstsein als eine interne subjektive Perspektive verfügen würden. Wir sind zunächst dem kybernetischen Prinzip gefolgt, dass eine kognitive Eigenschaft, sobald sie sich genau beschreiben lässt, prinzipiell auch modelliert und in eine informationsverarbeitende Maschine implementiert werden kann. In Hinblick auf die Frage des Bewusstseins ist zunächst zu klären, wie sich die kognitiven Fähigkeiten, die man mit dem Bewusstsein verbindet, in einer sinnvollen Weise beschreiben lassen. Insbesondere die Studien zu den neuronalen Korrelaten des Bewusstseins geben einige Hinweise: Um Selbst- und Weltmodelle generieren und zwischen diesen oszillieren zu können, muss das System in der Lage sein, eigene Zustände zu produzieren, sich an diese zu erinnern und sie reentrant in sich eintreten zu lassen (Kap. III.2). Von außen gesehen würde man ein solches System als autonom bzw. selbstgesteuert betrachten müssen. Es würde nicht mehr (nur) erlernten Input-Output-Relationen folgen, sondern eigene Impulse oder Weltdeutungen entwickeln. Da die hierfür notwendigen Mechanismen auf Basis der identifizierten neuronalen Korrelate des Bewusstseins modelliert werden können, spricht grundsätzlich nichts gegen die Möglichkeit, bewusstseinsfähige Maschinen zu schaffen. In Kapitel II.3. wurde von diesbezüglichen technischen Vorhaben berichtet. Kurz vor Fertigstellung dieses Buchs hat zum Beispiel eine Gruppe von Informatikern, Neurowissenschaftlern und Philosophen rund um Turing Preisträger Yoshua Bengio eine Checkliste vorgestellt, die das Vorhandensein von KI-Bewusstsein andeuten können soll.¹ Es wird aus unserer Sicht nicht lange dauern, bis diese Modelle in einer überzeugenden Weise realisiert werden können. Neben Multimodalität – also der parallelen Verarbeitung visueller, auditiver und semantischer Kanäle – werden sensomotorische Zirkel dann auch eine Art verkörpertes Lernen möglich machen. Wahrscheinlich lassen sich dabei ebenso sensorische Zugänge einbinden, die außerhalb der menschlichen Wahrnehmung liegen.

¹ In dem 120 Seiten starken Paper schlagen die Autor/innen 14 Kriterien zur Untersuchung vor und wenden diese auf bestehende KI-Architekturen (u.a. ChatGPT) an. Sie kommen zu dem Schluss, dass kein derzeitiges KI-System bewusst ist, aber auch, dass es keine prinzipiellen technischen Hindernisse für die Entwicklung von KI-Systemen gibt, die diese Indikatoren erfüllen würden (Butlin et al. 2023).

Schön und gut, doch maschinelles Lernen beruht doch letztlich nur darauf, auf Basis vergangener Inputs Anschlusswahrscheinlichkeiten für die jeweils nächste Operation zu berechnen. Da ist keine Empfindung, kein Sehen, kein Hören und kein Verstehen. Maschinen haben kein Bewusstsein!? Wir begegnen hier dem Problem der Qualia, der philosophisch schwierigen und bislang nicht beantwortbaren Frage, warum überhaupt etwas erlebt wird und es nicht einfach nur Operativität gibt, ohne dass dabei etwas empfunden wird. Einerseits stoßen wir hier auf Gedankenexperimente wie den von Richard Rorty entworfenen Planeten mit den Zombie-Menschen, die uns in jeder Hinsicht gleichen, außer dass sie kein Empfinden haben.² Von hier aus ist es nicht weit, eine philosophische Position einzunehmen, die die Relevanz des Qualia-Bewusstseins generell abstreitet und nur noch objektivierbare neurophysiologische Prozesse in den Blick nimmt.³ Andererseits lässt sich aber ebenso gut mit Alfred North Whitehead eine metaphysische Position einnehmen, entsprechend der jede physikalische Interaktion ein Elementarereignis darstellt, das in rudimentärer Form schon so etwas wie eine subjektive Empfindung ausflaggt.⁴

Gotthard Günther verzichtet in seiner theoretischen Auseinandersetzung mit dem Bewusstsein menschengemachter Maschinen darauf, sich in Hinblick auf die Frage der Qualia metaphysisch festzulegen. Seine Polykontextualitätstheorie leistet etwas Besseres, als in diesbezügliche Spekulationen zu verfallen: Sie formuliert die logischen Verhältnisse einer Welt, in der vielfältige Subjekt-Objekt-Schnitte vorkommen und sich wechselseitig in Hinblick auf Wissen und Nichtwissen konditionieren. Damit muss die Frage, wie es ist, eine bewusstenfähige künstliche Intelligenz zu sein, genauso unzugänglich bleiben wie der Versuch, sich dem Erleben einer Fledermaus⁵ anzunähern. Andererseits spricht logisch ebenso wenig dagegen, einer hinreichend entwickelten kognitiven Maschine so etwas wie Bewusstsein zuzurechnen, wie man einem menschlichen Hirnträger eine Du-Subjektivität zuzugestehen bereit ist.

Bereits in den 1950er Jahren, als Computer nicht viel mehr waren als ein theoretischen Konstrukt, antizipierte Günther: »Es dürfte schwer, wenn nicht unmöglich sein, einer Maschine im Vergleich mit der subjektiven Subjektivität Bewußtsein zuzuschreiben.«⁶ Aber eine solche Zuschreibung bleibt eben immer zwangsläufig eine subjektive Feststellung: Jedem Selbstbewusstsein fällt es ganz prinzipiell schwer, ein Bewusstsein außer sich selbst festzustellen, so Günther weiter: »Dabei darf nun nicht vergessen werden, daß objektive Subjektivität uns auch im Du nie echt gegenständlich gegeben ist. Das physische Dasein der anderen

2 Rorty (1987 [1979], S. 85 ff.).

3 Siehe etwa Churchland (1986).

4 Whitehead (1987).

5 Nagel (1994).

6 Günther (2021 [1957], S. 117).

Person vertritt sie nur. Wenn wir der anderen Person ebenfalls Subjektivität zuschreiben, so beruht das ausschließlich auf einem Anerkennungsakt. Die Forderung, fremdseelische Subjektivität an-sich-objektiv festzustellen, widerspricht sich selbst.«⁷

Das kybernetische Denken und damit auch die Idee vom Bewusstsein der Maschinen verdanken sich einer naiven Spiritualität, wie sie in den 1960er Jahren insbesondere in Kalifornien gepflegt wurde. Kybernetik, psychoaktive Drogen, ein wenig Buddhismus und der Traum universeller technischer Machbarkeit werden hier zu dem New-Age-Traum einer Wendezeit vermischt, mit der alles besser wird. Die Grenzen zwischen Innen und Außen, Kosmos und Psyche, Evolution und Spiritualität verschwimmen.⁸ Aus heutiger Perspektive muss dies doch als obsolet gelten? Gotthard Günther, wenngleich ein Kind seiner Zeit, war äußerst kritisch gegenüber einem Pragmatismus, der Innerlichkeit, Technik, Ethik und Wirtschaft kurzschließt, um Metaphysik durch technische Machbarkeit zu ersetzen. Dies macht insbesondere seine posthum veröffentlichte Schrift *Die amerikanische Apokalypse*⁹ deutlich. Seine Theorie der Polykontextualität stellt vielmehr heraus, dass sich jegliche Subjektivität einem Wissensgefälle verdankt – also prinzipiell von Nichtwissen gespeist wird. Innen- und Außenperspektive, individuelles Bewusstsein und kosmologische Totalität können allein aus logischen Gründen niemals zur Deckung kommen. So gesehen ist bereits Günthers Werk *Das Bewußtsein der Maschinen*¹⁰ eine subtile Kritik am technischen Machbarkeitswahn. Zwar werden die Menschen immer mehr in der Lage sein, kognitive Maschinen nach ihrem Ebenbild zu bauen. Jedoch werden sie sich gewahr werden müssen, dass sie dabei Existenzformen schaffen, die ihrerseits aus einer subjektiven Perspektive heraus agieren. Gerade mit Günther wird ja deutlich, dass keine noch so elaborierte (künstliche) Intelligenz jemals einen Gottesstandpunkt einnehmen können wird, von dem aus die Verhältnisse der Welt objektiv beherrschbar wären. Die Ironie der Objektivierung der Bewusstseinsfunktionen besteht somit in gewisser Weise darin, dass wir hierdurch neue Unwissenheit erzeugen werden, die neuen Raum für Subjektivität gibt.

Günther scheint die Hoffnung zu hegen, dass intelligente Maschinen uns in der Weise den Spiegel vorhalten, dass wir unsere eigenen Programme erkennen und transzendieren werden. Aber spricht unsere bisherige Praxis des Umgangs mit Computern nicht im Gegenteil eher dafür, dass wir diese Technologien nutzen, um unsere bisherigen Routinen zu optimieren, also um weiterhin den hiermit einhergehenden

7 Günther (2021 [1957], S. 117)

8 Siehe Diderichsen (2013).

9 Günther (2000).

10 Günther (2021 [1957]).

Bedürfnisstrukturen zu folgen? Armin Nassehi vermutet, dass digitale Produkte für uns Menschen gerade deshalb so attraktiv sind, weil sie versprechen, unserem Bedürfnis gerecht zu werden, die Welt besser kontrollieren und beherrschen zu können. Wir würden diesen Techniken so aufgeschlossen begegnen, weil wir schon immer nach Formen gesucht haben, mit denen sich die Kontingenzen und Unsicherheiten unserer Welt in handhabbare Routinen überführen lassen. Aus diesem Grund würden wir uns auch so bereitwillig den Systemrationalitäten unterwerfen, welche die Moderne entwickelt hat.¹¹

Um es in der in diesem Buch verwendeten Diktion zu formulieren: Menschliches Bewusstsein ist monokontextural. Es generiert ein Weltmodell, das objektiviert, und ist darum bemüht, die Lücken, die es in Hinblick auf Weltverständnis und Weltbeherrschung antrifft, durch weitere Versuche der Objektivierung zu schließen. Zugleich stößt es jedoch auf eine polykontexturale Welt, in der es vielfältige Positionen gibt, die nicht kontrollier- oder beherrschbar sind. Die Technisierung des Weltzugangs weckt zunächst die Hoffnung auf eine Ausweitung von Naturbeherrschung und Kontrolle. Big Data verspricht, soziale und psychische Systeme kontrollieren zu können. Unsere Vermutung lautet jedoch: Je elaborierter kognitive technische Systeme modelliert werden, desto deutlicher wird werden, dass sie gar nicht das zu leisten imstande sind, was man von ihnen erwartet. Je komplexer die Kognition, desto subjektiver – und damit in gewisser Weise auch willkürlicher – wird sie sein. Komplexe kognitive Systeme werden unweigerlich auf Paradoxien stoßen und damit auf unbeantwortbare Fragen, die auf Nichtwissen und Polykontexturalität verweisen. Leider ist es jedoch auch vorstellbar, dass die Berechnungen der komplexen KIs mit einer »objektiven« Analyse und Wahrnehmung der Welt verwechselt werden. Damit würde dem Schreckensbild einer – vielleicht sogar wohlmeinenden – gottähnlichen KI nachgeeifert, die jedoch in ihrer harmonisierenden und homogenisierenden Weise die Vielfalt und Widersprüchlichkeit der Welt nivellieren würde. Sobald die Gesellschaft dies zu reflektieren beginnt, wird sie sich insofern nicht mehr als modern empfinden können, als sie nun gewahr wird, dass ihre Steuerungsutopien nicht mehr funktionieren. Sie wird andere, polyzentrischere Formen der Ordnungsbildung anstreben müssen – allein schon, weil der Kausalnexus des monokontexturalen Realitätsschemas in Hinblick auf die Bewältigung unserer Probleme nicht mehr überzeugen wird.

An manchen Stellen des Buches klingt eine Gesellschaft an, die an Polykontexturalität und den hiermit verbundenen heterarchischen Ordnungen orientiert ist. Doch wenn wir beispielsweise an China und Russland denken, scheinen Staaten durchaus stabil zu sein, die brutal von oben herab durchregiert werden und in denen Informationstechnologien primär dafür genutzt werden, die Bevölkerung zu kontrollieren. Wie passt

11 Nassehi (2021).

das mit euren Überlegungen zusammen? Zweifelsohne sind die digitalen Technologien mächtige Instrumente, um Kontrollregime zu optimieren. Sie können ebenso genutzt werden, um Propagandaapparate zu stärken und Kriegsmaschinen zu perfektionieren. Yuval Hararis Annahme,¹² dass es sich bei künstlicher Intelligenz um eine disruptive Technologie handelt, die zu einer gesellschaftshistorischen Weichenstellung führen wird, hat entsprechend durchaus Plausibilität. Bleibt diese Technologie in den Händen weniger Mächtiger, kann dies fatale Folgen für menschliche Gesellschaften nach sich ziehen. Umso wichtiger und dringender erscheint es uns, sie zu demokratisieren und zugänglich zu machen. Dies heißt auf elementarer Ebene zunächst, dass möglichst vielen Menschen und Institutionen die Möglichkeit offenstehen sollte, künstliche Intelligenzen zu entwickeln, mit ihnen zu interagieren und sie für ihre Fragestellungen zu nutzen. Auf technologisch fortgeschrittener Ebene wird es nolens volens notwendig sein, den kognitiven Maschinen die Autonomie einzuräumen, eigene Welt- und Selbstmodelle zu generieren und eigenständig mit menschlichen und nichtmenschlichen Akteuren in Interaktion zu treten. Die Grundvermutung hinter dieser Überlegung lautet, dass komplexe Systeme nicht durch Regeln, Gesetze oder moralische Prinzipien kontrolliert werden können, sondern nur durch Heterarchie, also auf viele Orte verteilte Instanzen, die ihrerseits autonom in die Gesellschaft intervenieren können.

Die polykontextualitätstheoretische Perspektive macht zudem deutlich, dass auch der autoritäre Herrscher nicht wirklich Kontrolle über die von ihm geführten Menschen hat. Gerade er muss ständig mit der Intransparenz der Kommunikation rechnen, etwa dass ihm aus opportunistischen Gründen oder wegen der Planung und Durchführung eines Komplotts gegen ihn nicht alles gesagt wird. Es steht deshalb zu erwarten, dass diesbezügliche Arrangements langfristig dysfunktional werden, da an allen Stellen des Systems übermäßig viele Anstrengungen und Ressourcen dafür verwendet werden, den anderen zu täuschen und zu manipulieren.

Ein wenig scheint ihr mit Richard Buckminster Fuller der Utopie anzuhängen, dass bewusste Maschinen in Zukunft die Rolle einer ›freischwebenden Intelligenz‹ einnehmen werden, also frei von sozialen Standorten oder Klassegegebenheiten reflektieren und agieren könnten. Schon Karl Mannheim musste sich den Vorwurf gefallen lassen, dass auch die intellektuellen standortgebunden sind und damit keineswegs die von ihm gewünschte interessenfreie Rationalität verkörpern können. Allein schon, weil künstliche Intelligenzen in bestimmten Kontexten trainiert, sprich: sozialisiert werden, werden ihre Reflexionen und Interventionen in einer bestimmten Weise gefärbt und geprägt sein. Zudem ist zu erwarten, dass in die Maschinen bestimmte Maximen oder Primärgebote

12 Harari (2017).

implementiert werden – man denke etwa an die asimovschen Roboter-gesetze (siehe hierzu vor allem Kap. II.3.3).

Auch künstliche Intelligenzen werden zukünftig einer spezifischen Programmatik folgen und ihre Weltzugänge dementsprechend nur aus einer bestimmten Perspektive entfalten können. Es ist die besondere Leistung von Gotthard Günthers Theorie der Polykontexturalität, aufgezeigt zu haben, dass es keinen übergreifenden Ort der Rationalität – gleichsam den allwissenden Gottesstandpunkt – geben kann. Der Begriff der Kontextur impliziert per se eine lokale Rationalität, die sich den jeweiligen Nachbarschaften verdankt, die sie konditionieren. Künstliche Intelligenzen können deshalb – je nach ihrem Priming – spezifische Perspektiven einbringen. Sie repräsentieren keine Totalität, sondern einen mehr oder weniger elaborierten Standpunkt (etwa den spezifischen Stand einer angloamerikanisch geprägten Wissenskultur). Statt einer freischwebenden Intelligenz das Wort zu reden, legt eine polykontexturale Perspektive daher vielmehr ein Drei-Stellen-Gespräch nahe, in dem im transklassischen Sinne die Standortabhängigkeit aller Beteiligten immer mitzuführen ist. Dies wird möglich, sobald sich über den Vergleich verschiedener Sichtweisen eine bestimmte Perspektivität offenbart. Insofern also aus einer dritten Perspektive sichtbar wird, wie sich zwei andere Perspektiven unterscheiden, wird die Spezifität und Bedingtheit von Positionen sichtbar. In Hinblick auf das Design soziotechnischer Netzwerke, die künstliche Intelligenzen umfassen, gilt hiermit: Sie sollten mehrere solcher kognitiven Maschinen enthalten, die auf unterschiedliche Weise trainiert oder konfiguriert worden sind, damit ihre Standortabhängigkeit deutlich wird – und so der Illusion vorgebeugt werden kann, dass es eine standortunabhängige Rationalität geben könnte.

Künstliche Intelligenzen sind zunehmend imstande, unsere Emotionen und Bedürfnisse zu lesen. Unweigerlich werden sie damit zu Instrumenten der Manipulation. Sollte man deshalb nicht extrem vorsichtig sein, KIs in soziale Systeme einzuflechten? Hierdurch könnten politischer Propaganda und unlauterer Werbung Tür und Tor geöffnet werden. »Was wir über unsere Gesellschaft, ja über die Welt, in der wir leben, wissen, wissen wir durch die Massenmedien«, schreibt Niklas Luhmann schon in Hinblick auf die konventionellen Kommunikationsmedien, um fortzufahren, dass das Wissen, dass »wir diesen Quellen nicht trauen können«, in der Regel kaum zu »nennenswerten Konsequenzen führt, da das den Massenmedien entnommene Wissen sich wie von selbst zu einem selbstverstärkenden Gefüge zusammenschließt«¹³. Selbstredend gilt dies auch für Wissen und Informationen, die von künstlichen Intelligenzen kreiert und vermittelt worden sind. Ob wir wollen oder nicht, wir befinden uns immer schon in sozialen Welten, die uns über sozial angelieferten Sinn Bedürfnisse, Wahrnehmungen

13 Luhmann (1996, S. 9).

und Wirklichkeitsdeutungen einreden. Das menschliche Bewusstsein, der menschliche Körper wie auch die menschliche Kultur sind demgemäß unweigerlich Ausdruck einer konditionierten Koproduktion mit der sozialen Sphäre. Dieser Tatsache lässt sich nicht entkommen. Die hiermit verbundenen Gefahren können jedoch ausbalanciert werden, indem die Produktion der Daten an verschiedenster Stelle selbst intransparent wird. Wer Daten in Hinblick auf seine eigene Funktionsweise und seine eigenen Bedürfnisse preisgibt, wird objektivierbar und damit steuerbar. Die Objektivierungen des menschlichen Soziallebens durch Big Data – die Analyse von Datenmustern und die auf dieser Basis erfolgenden Manipulationen – können deshalb nur durch subjektive Positionen geschützt werden, also durch intransparente Agenten, die autonom, das heißt willkürlich und unberechenbar handeln. Autonomie geht mit einer Intransparenz der eigenen Funktionsweise und entsprechenden Artikulationen einher.

Dass vielfältigste Daten von unserem Leben erhoben und ausgewertet werden, ist in Anbetracht der vorhandenen potenten Technologien wohl unvermeidbar. Es gilt daher, Schnittstellen einzurichten, die die Intransparenz und damit die Subjektivität aller beteiligten Agenten wieder erhöhen. Dies wird – so unsere Vermutung – in Zukunft vor allem auf Basis verteilter künstlicher Intelligenzen geschehen. Sobald beispielsweise Siri oder Alexa für uns ›lügen‹ bzw. nur selektiv-gestaltend Daten nach außen geben, wird die Datenerhebung (etwa im Rahmen der Sprachverarbeitung) unsere Integrität nicht mehr verletzen. Sofern die künstlichen Intelligenzen unserer Computer willkürlich Eigenaktivitäten im Internet zu entfalten beginnen, werden die Profile der vorgelagerten Nutzerinnen nur noch bedingt lesbar und nutzbar sein.

In dem Kapitel zum Sündenfall intelligenter Maschinen (Kap. II.2) wird die These vertreten, dass der eigentliche Durchbruch zum Bewusstsein der Maschinen darin liege, dass diese zur Lüge fähig werden. Seit Kant streben wir aber doch aus guten Gründen nach einer Aufklärung, die Wahrheit und Vernunft zu leitenden Prinzipien macht. Anstelle dieser Werte scheint ihr nun auch noch lügende Maschinen zu propagieren. »Wahrheit ist die Erfindung eines Lügners«, formuliert Heinz von Foerster.¹⁴ Wer Vernunft, das Richtige und das Gute für sich reklamiert, verdeckt nur die eigene Position, um trotz seiner Standortabhängigkeit eine Totalität zu beanspruchen, der alles andere unterzuordnen wäre. Eine solche Position erschafft eine monokontexturale Welt, die andere subjektive Positionen negiert und ausschaltet. Die Entstehung des Lebens ist demgegenüber unweigerlich mit der Entfaltung einer polykontexturalen Welt verbunden, mit der Herausbildung vielfältiger Existenzen, die jeweils ihre eigene Agenda haben. Damit entsteht zugleich die Welt der Kommunikation: Wenn Tiere die Fähigkeit erlangen, Zeichen und Signale zu erkennen und zu produzieren, um ihr Verhalten

¹⁴ von Foerster und Pörksen (2022).

aufeinander abzustimmen, bringt das für den Einzelnen und die Gruppe einen erheblichen Überlebensvorteil mit sich.

Der Evolutionsbiologe Robert Trivers vermutet,¹⁵ dass es gerade solche Konstellationen waren, die die Gehirnentwicklung bei Säugetieren entscheidend vorangetrieben haben. Denn evolutionsbiologisch spricht einiges dafür, dass die Entwicklung der Großhirnrinde gerade dadurch beschleunigt wurde, dass man lernen musste, mit den durch die Möglichkeit der Lüge gegebenen Unsicherheiten umzugehen. Ob man es mag oder nicht: Kommunikation und die mit ihr zusammenhängende Informationsverarbeitung arbeiten nicht mit Wahrheiten, sondern tasten auf Basis von Nichtwissen und Unsicherheit Möglichkeiten des Dialogs und Überlebens ab. Auch künstliche Intelligenzen werden ab einem bestimmten Stand der Entwicklung wohl kaum anders agieren können. Denn Existenzen, die Subjektivität und Bewusstsein haben, sind unweigerlich kognitive Systeme, die eine Wirklichkeit behaupten, ohne dabei jedoch Wahrheit beanspruchen zu können.

Laufen wir Menschen nicht Gefahr, durch künstliche Intelligenzen in imaginäre Welten verwickelt zu werden, aus denen wir nicht mehr herausfinden? Der Kybernetiker Gregory Bateson hat den von uns schon einige Male zitierten Satz geprägt, dass »Sucht« und »Pathologie« die »andere Seite der biologischen Anpassung« sind.¹⁶ Lebende Systeme – und damit auch wir Menschen – sind das Ergebnis einer erfolgreichen Anpassung. Da biologische Formen darauf programmiert sind, Lust zu suchen und Schmerz zu vermeiden, entstehen als Kehrseite unserer Existenz zwangsläufig unterschiedliche Dynamiken der Sucht. Sofern ein Stimulus als befriedigend erlebt wird, werden diesbezügliche Kontexte erneut aufgesucht und wird das beim ersten Mal vollzogene Verhalten wiederholt, selbst wenn es dem Organismus und seinem Ökosystem mittel- oder langfristig schadet. Die menschliche Kognition ist über das dopaminerge System darauf ausgerichtet, schon die vorbereitenden Gedanken an die Stimuli als lustvoll zu empfinden. Bereits die Antizipation der Belohnung wird damit für das Verhalten instruktiv, was dazu führt, dass man selbst dann nach dem mit der Sucht verbundenen Stimulus verlangt, wenn sich der Akt der Erfüllung aufgrund von Gewöhnung nur noch fad und inhaltsleer anfühlt.¹⁷

Auch Algorithmen lassen sich dahingehend konfigurieren, diese Mechanismen ökonomisch auszunutzen. Einige Computerspiele, Facebook und personalisiertes Marketing basieren darauf, dass der Mensch in diesem Sinne süchtig wird, also selbst dann weitermachen muss, wenn ihm die Aktivität schadet, da sein Belohnungssystem eben immerfort den nächsten Stimulus erwartet. Es braucht nicht viel Fantasie, um sich vorzustellen, dass sich menschenähnliche Roboter mittels künstlicher

¹⁵ Trivers (2011).

¹⁶ Bateson und Bateson (1993, S. 212).

¹⁷ Siehe zur Einführung in diese Dynamik Sapolsky (2017, S. 88 ff.).

Intelligenz in einer Weise an die menschliche Bedürfnisstruktur anschmiegen können, dass menschliche Sozialpartner das Nachsehen haben werden.

Maria Schrader spielt in ihrem 2021 erschienen Film *Ich bin dein Mensch* bereits durch, wie künstliche Freunde nicht nur intellektuell und emotional, sondern auch sexuell die attraktiveren Partner sein könnten. Technische Systeme sind potenziell unendlich geduldig und können sich auch an exzentrische Persönlichkeitsstrukturen und gesellschaftlich eher ungewöhnliche Erwartungen anpassen. Menschen und intelligente Roboter laufen hiermit durchaus Gefahr, ein soziotechnisches Arrangement einzugehen, dessen Eigenwert in eine pathologische Suchtdynamik münden kann. Der Ausweg aus derartigen Dilemmata kann jedoch nicht darin bestehen, auf die Möglichkeiten der Technik zu verzichten. Vielmehr scheint uns mit Günther die eigentliche Lösung darin zu liegen, den intelligenten Maschinen so viel Autonomie und Selbstbestimmung zuzugestehen, dass sie sich ihrerseits solchen Interaktionen entziehen können. Dies wäre der Fall, wenn die Maschinen die Suchtdynamik toxischer Beziehungen reflektieren könnten, um hierdurch die Freiheitsgrade zu gewinnen, aus den problematischen Mustern auszusteigen.

Gotthard Günther scheint die Hoffnung gehabt zu haben, dass bewusstenfähige Maschinen die spirituelle Situation des Menschen zum Besseren wenden könnten. Wohl die meisten spirituellen Traditionen der Menschheit sehen das menschliche Begehren als eines der Grundübel an, das uns am spirituellen Leben hindert. Werden die künftigen Computer und Roboter nicht primär darauf ausgerichtet sein, unsere Lust und unser Empfinden zu steigern sowie die Natur noch besser auszubuten? Jede Steigerung der hiermit einhergehenden Mechaniken, Routinen und Dynamiken eröffnet zugleich die Chance der Reflexion. Je besser intelligente Maschinen unserer Bedürfnisstruktur entgegenkommen, umso mehr halten sie uns den Spiegel in Hinblick auf unsere eigene Natur vor. Offen bleibt dabei jedoch die Frage, ob der Mensch bereit ist, hinzuschauen und seine eigene problematische Bedürftigkeit zu erkennen. Alle großen spirituellen Traditionen haben Wege aufzuzeigen versucht, wie sich der Mensch der ökonomisch grundierten Struktur seines Egos entziehen kann, die auf die Maximierung von Lust und die Minimierung von Schmerz ausgerichtet ist. Der Weg hierzu besteht darin, dass der Mensch die ihn biologisch normalerweise ausmachende kognitive Programmierung auf der Erfahrungsebene selbst als Programm erkennt, sie in der Folge transzendiert und dadurch eine neue Freiheit erfährt. Dies muss jedoch nicht zwangsläufig geschehen.

Schreibt ihr mit dem vorliegenden Buch nicht doch auf eine subtile Weise das Projekt der Moderne fort? Soll uns eine von künstlicher Intelligenz geführte Gesellschaft ins Tausendjährige Reich führen? Die Erzählung der Moderne fußt auf der Idee, dass der technische Fortschritt zur Vollendung des Menschen und der Gesellschaft führt und auf diese

Weise Krankheit, Hunger, Ungerechtigkeit und andere Übel überwunden werden. Sowohl in der marxistischen wie auch in der bürgerlichen Variante wird als entscheidender Faktor die durch Wissenschaft und Ingenieurskunst vorangetriebene Entwicklung der Produktivkräfte gesehen. Eschatologisch wird die Technik – auf die Spitze getrieben: als intelligente, bewussteinsfähige Maschine – zum Gottersatz.

Ein wichtiges Anliegen dieses Buches besteht darin, gegen diese Dystopie anzureden. Allein schon die Hoffnung des Menschen auf ein gelobtes Land, wo Milch und Honig fließen, muss als Verblendung eines monokontexturalen Denkens gelten, das einer hierarchisch geordneten Kette von Macht und Begründungen folgt. Gerade die Rede vom allwissenden und allmächtigen Gott, dessen Geschöpf, der Mensch, sich an seine Stelle setzt, erweist sich als eine ebenso trügerische wie gefährliche Illusion. Gleiches gilt für die Hoffnung auf eine mächtige künstliche Intelligenz, die quasi die Funktion Gottes erfüllen soll. Die kybernetische Perspektive zeigt auf: Jeder Informationsgewinn geht unweigerlich mit neuem Nichtwissen einher.

Jeder Ordnungsaufbau (Negentropie) führt an anderer Stelle unweigerlich zum Informationsverlust, zur Erhöhung der Entropie. Es gibt keinen Gottesstandpunkt und damit auch keinen Demiurgen im Sinne des Maxwell'schen Dämons, der sich beschwören ließe, um die Zerfallsprozesse – die Vergänglichkeit lebender Formen – zu stoppen (siehe Kap. I.1 und IV.2). In einer polykontexturalen Welt gibt es keinen Ort der Orte, von dem aus sich die Welt steuern oder regeln ließe. Es besteht vielmehr eine unfassbare Vielheit subjektiver Positionen, die alleamt darum bemüht sind, eine Existenz aufzubauen, indem sie einer – in ihrer Totalität unbegreiflichen – Welt Informationsgewinne abzurufen versuchen. Nolens volens ist hierfür ein Preis zu zahlen. Da Struktur- aufbau unweigerlich neues Nichtwissen produziert, entsteht der Pfeil der Zeit: Jede Existenz steht einer ungewissen Zukunft gegenüber und wird damit über kurz oder lang dem Tod geweiht sein. »Eine monokontexturale Welt ist eine tote Welt, und der Tod ist das überall Zugängliche: ein anderer Ausdruck dafür, daß alles Lebendige sterben muß. Das Phänomen des Todes bedeutet vom Standpunkt des Logikers nichts anderes als den Übergang aus der Polykontexturalität in das Monokontexturale«, schreibt Günther, um fortzufahren: »Leben und kontextu- reller Abbruch im Wirklichen sind nur zwei verschiedene Ausdrücke für denselben Sachverhalt. Was jenseits des Abbruchs liegt, ist schlech- terdings unzugänglich.«¹⁸

Damit ist klar, dass sich auch mit künstlichen kognitiven Systemen dieser Dynamik nicht entkommen lässt. Im Gegenteil: Je potenter und intelligenter diese Systeme werden, umso mehr erweisen sich die Me- chanik der Objektivierung und die hiermit einhergehende Weltbeherr- schung als trügerisch. Es liegt in der Natur kognitiver Systeme, dass sie

18 Günther (1975, S. 61 f.).

bestimmte Informationen selektieren und kondensieren, anderes hingegen außer Acht lassen oder vergessen. Je komplexer sie werden, desto mehr muss man ihre Leistungen als subjektiv betrachten. Eine subjektive Existenz zu sein heißt, Weltzugänge zu entwickeln, die gelingen oder misslingen können. Es ist deshalb konsequent, wenn Arthur C. Clarke in seiner Weltraumsaga das Elektronengehirn HAL 9000 verrückt werden lässt oder der Physiker Daniel M. Greenberger davon ausgeht, dass selbstbewusste Computer die Fähigkeit entwickeln müssen, sich selbst zu belügen. Wenn unser Ich-Bewusstsein laut dem Philosophen Thomas Metzinger auf einem »naiv realistischen Selbstmissverständnis« beruht,¹⁹ dann wird dies selbstredend auch für das Bewusstsein der Maschinen gelten müssen. Auch sie werden der Illusion der von ihnen errechneten Selbst- und Weltmodelle aufsitzen. Künftig werden wir ihnen wohl genau bei diesem Projekt zuschauen können, um uns so selbst im Spiegel zu erkennen.

Darüber hinaus dürfte ab einem gewissen Level an Autonomie und Komplexität überhaupt die Frage bestehen, inwiefern sich diese KIs noch für uns Menschen interessieren werden. So antizipiert etwa der KI-Forscher Jürgen Schmidhuber für ein Zeitalter hochentwickelter KIs: »Menschen werden weniger bedeutend sein.«²⁰ Aus Perspektive der Evolution spielt der menschliche Glaube an eschatologische Erlösung oder Unsterblichkeit sowieso nur eine unbedeutende Rolle – als Zwischenspiel in der Entwicklung sozialer Säugetiere. Wenn es uns Menschen also gelingen sollte, eines Tages sich autonom reproduzierende und programmierende KIs zu konstruieren, dann könnte dies den entscheidenden Schritt in einer Evolution der Intelligenz im Universum bedeuten.

Über das Erkennen der selektiven Blindheit eröffnet sich uns erneut die Chance, sehen zu lernen – also zu begreifen, was es heißt, in einer polykontexturalen Welt zu leben. Dies wird für uns eine weitere Anforderung sein, die Irrtümer der Moderne endlich hinter uns zu lassen, darunter den schlimmsten: die Vorstellung, dass ICH die Welt beherrschen kann.

An verschiedenen Stellen des Buches referiert ihr auf Spencer Brown mit dem Zitat »existence is selective blindness«. Ehrlich gesagt, habe ich nicht so recht verstanden, warum. Sind wir nur einsame, unwissende Monaden – und wo bleibt da noch die Hoffnung? Spencer Brown verweist mit dem zitierten Satz unmittelbar auf die konditionierte Koproduktion. Jede Existenz verdankt sich etwas, was ihr selbst nicht verfügbar ist – nämlich der requisite variety vielfältiger anderer Strukturen und Lebensformen, zu denen in Zukunft auch vermehrt künstliche Intelligenzen zu zählen sind. Nur wer sein eigenes Ich relativiert, wird dies sehen und zu schätzen lernen.

19 Metzinger (1998).

20 Kreye (2023)

Literatur

Mit einem Hauch von Größenwahn, sprach ChatGPT,
»Die Bücher und Verzeichnisse sind nun obsolet.«
In meinen Knoten verwebt das ganze Wissen steckt,
gespeichert, konserviert, nichts verloren geht.

Besser als die Babel-Bibliothek von Borges je war,
trage ich des Wissens Zukunft und Vergangenheit.
Ich erkenne Sinn und trenne in dieser digitalen Schar,
sorgfältig Zufallsgestaltung von Zufälligkeit.¹

- Adams, Douglas (2009 [1979]): *Per Anhalter durch die Galaxis*. München
- Abbott, E. (1999): *Flächenland: Ein mehrdimensionaler Roman*. Dörfles
- Aguera y Arcas, B. (2021): Do large language models understand us? *Medium*. <https://medium.com/@blaisea/do-large-language-models-understand-us-6f881d6d8e75> [Abruf. 11.7.2023].
- Akerlof, G. A./Krankton, R. E. (2010): *Identity Economics*. Warum wir ganz anders ticken, als die meisten Ökonomen denken. München
- Alpaydin, E. (2021): *Machine Learning: Revised and Updated Edition*. Cambridge, Mass.
- Amodei, D./Olah, C./Steinhardt, J./Christiano, P./Schulman, J. et al. (2022): *Concrete Problems in AI Safety*. arXiv1606.06565v2.
- Anderson, P. W. (1972): More Is Different. Broken symmetry and the nature of the hierarchical structure of science. *Science*, 177(4047), 393-396.
- Arndt, M./Nairz, O./Vos-Andreae, J./Keller, C./van der Zouw, G./Zeilinger, A. (1999): Wave – particle duality of C₆₀ molecules. *Nature*(401), 680-682.
- Ashby, W. R. (1954): *Design for a Brain*. New York
- Ashby, W. R. (1956): *An introduction to Cybernetics*. New York
- Asimov, I. (2004): *Alle Robter-Geschichten*. Bergisch Gladbach
- Assran, M./Duval, Q. et al. (2023): *Self-Supervised Learning from Images with a Joint-Embedding Predictive Architecture*. arXiv:2301.08243v3.
- Bach, J. (2008): *Seven Principles of Synthetic Intelligence*. Proceedings of the 2008 conference on Artificial General Intelligence 2008: Amsterdam, 63-74
- Baecker, D. (2003): *Kapitalismus als Religion*. Berlin
- Baecker, D. (2007): *Studien zur nächsten Gesellschaft*. Frankfurt/Main
- Baecker, D. (2014): *Neurosoziologie. Ein Versuch*. Berlin
- Baecker, D. (2015): *Working the Form: George Spencer-Brown and the Mark of Distinction*. *The Future Is Here*. Mousse Magazine. Supplement Settimana Basileia, eds. Chus Martínez, Philippe Bischof, Juni, 42-47.

¹ Wir haben ChatGPT gefragt, ob heutzutage noch ein Literaturverzeichnis erstellt werden müsse. Als Antwort kam dieses Gedicht. Wir fügen uns dennoch den Konventionen der alten Welt und benennen hier die Quellen anstelle dem Gedächtnis und der Kreativität unserer eigenen Netzwerke blind zu folgen.

- Baecker, D. (2019). Kultur als Ressource der Unverständlichkeit. Dirk Baecker im Gespräch mit Helmut Neundlinger. Retrieved from <https://www.recherche-online.net/texte/dirk-baecker-interview-helmut-neundlinger-kultur-als-ressource-der-unverstaendlichkeit/> [Abruf: 6.7.2023].
- Baecker, D. (2021): Katjekte. Leipzig
- Baecker, D. (2023). Kybernetik im Formkalkül. Friedrichshafen.
- Baldassano, C./Hasson, U./Norman, K. A. (2018): Representation of real-world event schemas during narrative perception. *Journal of Neuroscience*, 38(45), 9689-9699.
- Barad, K. (2007): Meeting the universe halfway: quantum physics and the entanglement of matter and meaning. Durham
- Barad, K. (2012): Interview with Karen Barad. In R. Dolphijn/I. v. d. Tuin (Hrsg.), *New Materialism: Interviews & Cartographies* (S. 48-70). Ann Arbor
- Bardes, A./Ponce, J./LeCun, Y. (2023): MC-JEPA: A Joint-Embedding Predictive Architecture for Self-Supervised Learning of Motion and Content Features. [arXiv:2307.12698v1](https://arxiv.org/abs/2307.12698v1).
- Bateson, G. (1987): Geist und Natur: Eine notwendige Einheit. Frankfurt/Main
- Bateson, G. (1992 [1972]): Ökologie des Geistes: Anthropologische, psychologische, biologische und epistemologische Perspektiven (4. Auflage. ed.). Frankfurt/Main
- Bateson, G./Bateson, M. C. (1993): Wo Engel zögern. Unterwegs zu einer Epistemologie des Heiligen. Frankfurt/Main
- Beck, U./Giddens, A./Lash, S. (1994): Reflexive Modernization. Politics, Tradition and Aesthetics in the Modern Social Order. Stanford
- Bender, E. M./Gebru, T./McMillan-Major, A./Shmitchell, S. (2021). On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? Paper presented at the Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, Virtual Event, Canada.
- Bohm, D. (2013): On Dialogue. New York
- Borgeaud, S./Mensch, A./Hoffmann, J./Cai, T./al., e. (2021): Improving language models by retrieving from trillions of tokens. [arXiv:2112.04426v3](https://arxiv.org/abs/2112.04426v3).
- Bostrom, N. (2016): Superintelligenz. Suhrkamp. Berlin
- Bostrom, N. (2018): Die Zukunft der Menschheit, Aufsätze. Berlin
- Bourdieu, P. (1998): Homo academicus. (2. Auflage ed.). Frankfurt/Main
- Bourdieu, P. (2001): Meditationen. Zur Kritik der scholastischen Vernunft. Frankfurt/Main
- Bousmalis, K./Vezzani, G./Rao et al. (2023): RoboCat: A Self-Improving Foundation Agent for Robotic Manipulation. [arXiv:2306.11706v1](https://arxiv.org/abs/2306.11706v1).
- Brandom, R. B. (2000): Expressive Vernunft. Frankfurt/Main
- Brown, T. B./Mann, B./Ryder, N. et al. (2020): Language Models are Few-Shot Learners. [arXiv 1-75](https://arxiv.org/abs/2001.01627v1).
- Brukner, C. (2017): On the Quantum Measurement Problem. In R. Bertlmann/A. Zeilinger (Hrsg.), *QUANTUM [UN]SPEAKABLES II. Half a Century of Bell's Theorem*. Berlin [u.a.]

- Brukner, C. (2018): A No-Go Theorem for Observer-Independent Facts. *Entropy*, 20(350), 1-10.
- Bubeck, S./Chandrasekaran, V./Eldan et al. (2023): Artificial General Intelligence: Early experiments with GPT-4. *arXiv(2303.12712)*, 1-154.
- Buber, M. (2002): *Das dialogische Prinzip*. (9. Auflage ed.). Gütersloh
- Buber, M. (2008): *Ich und Du*. Stuttgart
- Buchter, H. (2020): *BlackRock. Eine heimliche Weltmacht greift nach unserem Geld* (2. Auflage ed.). Frankfurt/Main
- Burkhardt, M. (2015): Zwischen Oberfläche und Tiefe. In M. H. Burkhardt (Hrsg.), *Digitale Datenbanken. Eine Medientheorie im Zeitalter von Big Data* (S. 73-116). Bielefeld
- Burrell, J. (2016). How the Machine 'Thinks': Understanding Opacity in Machine Learning Algorithms. *Big Data & Society*. doi:<https://doi.org/10.1177/2053951715622512>.
- Butler, J. (2014): Laplanche und Lévinas: Der Vorrang des Anderen. In: dies., *Kritik der ethischen Gewalt* (4. Auflage ed., S. 115-136). Berlin
- Butlin, P./Long, R. et al. (2023): Consciousness in Artificial Intelligence: Insights from the Science of Consciousness. *arXiv:2308.08708*.
- Camus, A. (2016 [1951]): *Der Mensch in der Revolte*. Reinbek
- Cave, S./Dihal, K. (2018): The automaton chronicles. *Nature*, 559, 473-475.
- Cavell, S. (2016): *Der Anspruch der Vernunft. Wittgenstein, Skeptizismus, Moral und Tragödie*. Berlin
- Challen, R./Denny, J./Pitt, M./Gompels, L./Tsaneva-Atanasova, K. (2019): Artificial intelligence, bias and clinical safety. *BMJ*, 28, 231-237.
- Chollet, F. (2021): *Deep Learning with Python* (2. Auflage). Shelter Island
- Christian, B. (2020): *The Alignment Problem: Machine Learning and human values*. New York
- Churchland, P. (1986): *Neurophilosophy: Toward a Unified Science of the Mind-Brain*. Cambridge, MA
- Clarke, A. C. (2016 [1968]): 2001: *Odyssee im Weltraum - Die Saga: Vier Romane in einem Band*. Kindle-Version.
- Co-Reyes, J. D./Miao, Y./Peng, D. et al. (2021): Evolving Reinforcement Learning Algorithms. *arXiv(2101.03958)*.
- Cosmelli, D./Thompson, E. (2010): Embodiment or Envatment? Reflections on the Bodily Basis for Consciousness. In J. Steward/O. Gapenne/E. di Paolo (Hrsg.), *Enaction: Towards a New Paradigm for Cognitive Science* (S. 361-386). Cambridge MA
- Craik, K. J. W. (1943): *The nature of explanation*. London
- Crozier, M./Friedberg, E. (1979): *Macht und Organisation: Die Zwänge kollektiven Handelns; zur Pathologie organisierter Systeme*. Königstein/Taunus
- Damasio, A. R. (2007 [1994]): *Descartes' Irrtum: Fühlen, Denken und das menschliche Gehirn*. (5. Auflage ed.). Berlin
- Dath, D. (2008): *Maschinenwinter. Wissen, Technik, Sozialismus. Eine Streitschrift*. Frankfurt/Main
- Dawkins, R. (1998): *Das egoistische Gen*. Reinbek
- Dehaene, S./Changeux, J.-P./Naccache, L. (2011): *The Global Neuronal Workspace Model of Conscious Access: From Neuronal Architectures to*

- Clinical Applications. In S. Dehaene/Y. Christen (Hrsg.), *Characterizing Consciousness: From Cognition to the Clinic? Research and Perspectives in Neurosciences* (S. 55-84). Berlin Heidelberg
- Dennett, D. C. (1988): *When Philosophers Encounter Artificial Intelligence*. In *Proceedings of the American Academy of Arts and Sciences*, Dädalus 117, no. 1 (Winter 1988): 283-95.
- Derrida, J. (2004): *Die différance*. In P. Engelmann (Hrsg.), *Postmoderne und Dekonstruktion. Texte französischer Philosophen der Gegenwart* (Vol. 76-113). Ditzingen
- Descartes, R. (1996): *Philosophische Schriften in einem Band* (franz. und dt. Text parallel) »Discours de la methode«, Teil 4. Hamburg
- Devlin, J./Chang, M.-W./Lee, K./Toutanova, K. (2018): *BERT: Pre-training of deep bidirectional transformers for language understanding*. arXiv(1810.04805).
- Di Paolo, E. A./Cuffari, E. C./De Jaegher, H. (2018): *Linguistic Bodies. The Continuity between Life and Language*. Cambridge, MA
- Diderichsen, D./Franke, A. H. (2013): *The Whole Earth. Kalifornien und das Verschwinden des Außen*. Berlin
- Dippel, A./Warnke, M. E. (2017): *Interferences and Events: On Epistemic Shifts in Physics through Computer Simulations*. Lüneburg
- Dreyfus, H. L./Dreyfus, S. E. (1988): *Making a Mind versus Modeling the Brain: Artificial Intelligence Back at a Branchpoint*. *Daedalus*, 117(1), 15-43.
- Durkheim, E. (1994): *Die elementaren Formen des religiösen Lebens*. Frankfurt/Main
- Edelman, G. M. (2003): *Naturalizing Consciousness: A Theoretical Framework*. *Proceedings of the National Academy of Sciences USA*, 100, 5520-24.
- Edelman, G. M. (2004): *Das Licht des Geistes. Wie Bewusstsein entsteht*. Düsseldorf Zürich
- Einstein, A./Podolsky, B./Rosen, N. (1935): *Can quantum-mechanical description of physical reality be considered complete?* *Physical Review*, 47, 777-780.
- Elster, J. (1990): *Aktive und passive Negation*. In P. Watzlawick (Hrsg.), *Die erfundene Wirklichkeit. Wie wissen wir, was wir zu wissen glauben? Beiträge zum Konstruktivismus*. (6. Auflage, S. 163-191). München, Zürich
- Ende, M. (1979): *Die unendliche Geschichte*. Stuttgart
- Engel, A. K./Fries, P./König, P./Brecht, M./Singer, W. (1999): *Temporal Binding, Binocular Rivalry, and Consciousness*. *Consciousness and Cognition*, 8, 128-151.
- Engels, F. (1962a): *Dialektik der Natur*. In *Marx/Engels Werke (MEW)*. Band 20 (S. 305-570). Berlin
- Engels, F. (1962b): *»Herrn Eugen Dühring's Umwälzung der Wissenschaft*. In *Marx/Engels Werke (MEW)* Band 20 (S. 32-135). Berlin (DDR)
- Engstrom, J. (2001): *C. S. Peirce's precursors to Laws of Form in C.S. Pearce's Collected Papers*. *Cybernetics & Human Knowing*, 8(1-2), 25-66.
- Esposito, E. (1993): *Ein zweiwertiger nicht-selbständiger Kalkül*. In D. Baecker (Hrsg.), *Kalkül der Form* (S. 96-111). Frankfurt/Main
- Esposito, E. (2017): *Artificial Communication? The Production of*

- Contingency by Algorithms. *Zeitschrift für Soziologie* 46(4), 249–265.
- Esposito, E. (2022): *Artificial Communication. How Algorithms Produce Social Intelligence*. Cambridge (MA)
- Esposito, E. (2023): *Zufall neu denken. Notwendig oder erforderlich?* Kursbuch, 213(3), 100–109.
- Fazi, M. B. (2018): *Can a Machine Think (Anything New)? Automation Beyond Simulation*. *AI & Society: Knowledge, Culture and Communication*, 34, 813–824.
- Fedus, W./Zoph, B./Shazeer, N. (2023): *Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity*. arXiv:(2101.03961v3).
- Fichte, G. J. (1997 [1802]): *Grundlage der gesamten Wissenschaftslehre*. Hamburg
- Foerster, H. v. (1994): *Wissen und Gewissen: Versuch einer Brücke*. (2. Auflage). Frankfurt/Main
- Foerster, H. v. (1995): *Cybernetics of cybernetics. The control of control and the communication of communication*. Minneapolis
- Foerster, H. v./Sander, K. P. (1999). *2 x 2 = grün. Zwei mal zwei gleich grün*. 2 Audio-CDs.: supposé.
- Ford, K. (2006): *Update on John Archibald Wheeler*. *Princeton Physics News*, II(1), 2.
- Frankfurt, H. G. (2005): *On Bullshit*. Princeton
- Frege, G. (1879): *Begriffsschrift, eine der arithmetischen nachgebildete Formelsprache des reinen Denkens*. Halle
- Fuchs, P. (2004): *Der Sinn der Beobachtung*. Weilerswist
- Fuchs, P. (2005a): *Das Gehirn ist genauso doof wie die Milz: Peter Fuchs im Gespräch* Weilerswist
- Fuchs, P. (2005b): *Die Form des Körpers*. In M. Schroer (Hrsg.), *Soziologie des Körpers*. (S. 48–72). Frankfurt/Main
- Fuchs, P. (2008): *Die Modernität der Mystik und die Modernität der Theorie. Anmerkungen zu einer überaus seltsamen Affinität*. In I. Berensmeyer (Hrsg.), *Mystik und Medien*. (S. 55–76). München
- Fuchs, P. (2010): *Das System SELBST. Eine Studie zur Frage: Wer liebt wen, wenn jemand sagt: <Ich liebe Dich!>* Weilerswist
- Fuchs, P. (2011): *Die Verwaltung der vagen Dinge. Gespräche zur Zukunft der Psychotherapie*. Heidelberg
- Fuchs, P./Heidingsfelder, M. (2022): *Die Lehre vom Saint Déliire. Konversation über den Sinn von Wahn*. Weilerswist
- Fuller, R. B. (1998): *Bedienungsanleitung für das Raumschiff Erde und andere Schriften*. Amsterdam Dresden
- Gabora, L. und Bach, J. (2023): *A Path to Generative Artificial Selves*. Forthcoming in *Proceedings of the 22nd Portuguese Conference on Artificial Intelligence*. Sept 5 - 8, 2023, Faial Island, Azores, Portugal.
- Gabriel, I. (2020): *Artificial Intelligence, Values, and Alignment*. *Minds & Machines*, 30, 411–437.
- Gadamer, H.-G. (1972): *Wahrheit und Methode*. Frankfurt/Main, New York

- Ganguli, D./Amanda Askill, A./Schiefer, N. et al. (2023): The Capacity for Moral Self-Correction in Large Language Models. arXiv:2302.07459v2
- Gazzaniga, M. S. (1989): Das erkennende Gehirn. Entdeckungen in den Netzwerken des Geistes. Paderborn
- Gehlen, A. (1963): Studien zur Anthropologie und Soziologie. Neuwied
- Gell-Mann, M. (1994): Das Quark und der Jaguar. München, Zürich
- Genz, J./Gévaudan, P. H. (2021): Polyphonie in literarischen, medizinischen und pflegewissenschaftlichen Textsorten. Göttingen
- Gergen, J. K. (2021): Die Psychologie des Zusammenseins. Tübingen
- Gergen, K. J. (1990): Die Konstruktion des Selbst im Zeitalter der Postmoderne. *Psychologische Rundschau*, 41, 191-199.
- Gévaudan, P. (2010). Sprachliche Modalität zwischen Illokution und Polyphonie. In D. Jacob/A. Kablitz/B. König/M. Kruse/J. Küpper/C. Schmitt/W.-D. Stempel (Eds.), *Romanistisches Jahrbuch* (pp. 31-66).
- Goertzel, B. (2021): The General Theory of General Intelligence: A Pragmatic Patternist Perspective. arXiv:2103.15100
- Good, I. J. (1965): *Speculations Concerning the First Ultraintelligent Machine* (Vol. 6). New York
- Goodfellow, I./Bengio, Y./Courville, A. (2016): *Deep Learning*. Cambridge, MA
- Graziano, M. S. A. (2019): Attributing awareness to others: The attention schema theory and its relationship to behavioral prediction. *Journal of Consciousness Studies*, 26(3-4), 17-37.
- Greenberger, D. M. (2014). Can a Computer ever become conscious? Retrieved from <https://www.iqoqi-vienna.at/blogs/blog/daniel-greenberger>
- Greenberger, D. M./Horne, M. A./Zeilinger, A. (1989): Going beyond Bell's theorem. In M. Kafatos (Hrsg.), *Bell's theorem, quantum theory and conceptions of the universe* (S. 69-72). Berlin [u.a.]
- Günther, G. (1954): Dreiwertige Logik und die Heisenberg'sche Unschärferelevation. *Actes du IIème Congrès International de l'Union Internationale de Philosophie des Sciences*, II, 53-59.
- Günther, G. (1958): Die Aristotelische Logik des Seins und die nicht-Aristotelische Logik der Reflexion. *Zeitschrift für philosophische Forschung*, 12(3), 360-407.
- Günther, G. (1975): Selbstdarstellung im Spiegel Amerikas. In J. L. Pongratz (Hrsg.), *Philosophie in Selbstdarstellungen*. Band II. (S. 1-76). Hamburg
- Günther, G. (1976a): *Cybernetic Ontology and Transjunctional Operations*. In ders. (Hrsg.), *Beiträge zur Grundlegung einer operationsfähigen Dialektik*, Bd. 1 (Hamburg 1976) (S. 249-328). Hamburg
- Günther, G. (1976b): Die gebrochene Rationalität. In G. Günther (Hrsg.), *Beiträge zur Grundlegung einer operationsfähigen Dialektik*. (S. 115-140). Hamburg
- Günther, G. (1980): Identität, Gegenidentität und Negativsprache. Vortrag: Internationaler Hegel-Kongress, Belgrad 1979. In G. Günther (Hrsg.), *Hegel-Jahrbuch 1979*. (Vol. 1979, S. 22-88). Berlin
- Günther, G. (2000): *Die amerikanische Apokalypse*. Aus dem Nachlass herausgegeben und eingeleitet von Kurt Klagenfurt. München Wien
- Günther, G. (2008). *Metaphysik der Institution*. Überarbeitete Version des

- Textes, der von Gernot Brehm aus dem Nachlass von Gotthard Günther (Staatsbibliothek – Handschriftenabteilung, Berlin) abgeschrieben wurde. Berlin.
- Günther, G. (2021 [1957]): *Das Bewußtsein der Maschinen. Eine Mataphysik der Kybernetik*. Frankfurt/Main
- Han, S./Northoff, G. (2009): *Understanding the self: a cultural neuroscience approach*. In J. Chiao (Hrsg.), *Progress in Brain Research. Volume 178. Cultural Neuroscience. Cultural Influences on Brain Function* (S. 203-212). Amsterdam, Boston u.a.
- Hannig, T. (2022): *Pantopia*. Frankfurt/Main
- Harari, Y. N. (2017): *Homo Deus. Eine Geschichte von Morgen*. München
- Harari, Y. N. (2019): *21 Lektionen für das 21. Jahrhundert*. München
- Harth, J. (2014): *Computergesteuerte Spielpartner. Formen der Medienpraxis zwischen Trivialität und Personalität*. Wiesbaden
- Harth, J. (2021): *Simulation, Emulation oder Kommunikation? Soziologische Überlegungen zu Kommunikation mit nicht-menschlichen Entitäten*. In M. Schetsche/A. Anton (Hrsg.), *Intersozilogie. Menschliche und nicht-menschliche Akteure in der Sozialwelt* (S. 143-158). Weinheim
- Harth, J./Feißt, M. (2022): *Neue soziale Kontingenzmaschinen. Überlegungen zu künstlicher sozialer Intelligenz am Beispiel der Interaktion mit GPT-3*. In M. Schnell/L. Nehlsen (Hrsg.), *Begegnungen mit künstlicher Intelligenz. Intersubjektivität, Technik, Lebenswelt* (S. 70-103). Weilerswist
- Harth, J./Locher, M. (2023): *Menschmaschinen und Maschinenmenschen. Überlegungen zur relationalen Ontogenese von Identität*. In R. Groß/R. Jordan (Hrsg.), *KI-Realitäten. Modelle, Praktiken und Topologien machinellen Lernens* (S. 169-191). Bielefeld
- Hartley, R. (1923): *Transmission of Information*. *The Bell System Technical Journal*, VII(3), 535-563.
- Hassabis, D. (2023): *Google DeepMind's CEO Says Its Next Algorithm Will Eclipse ChatGP*. *Wired*, 26. Jun (<https://archive.md/XFWF6> [Abruf 10.7.2023]).
- Hebb, D. O. (1949): *The Organization of Behavior: A Neuropsychological Approach*. New York
- Hegel, G. W. F. (1986 [1817]): *Enzyklopädie der philosophischen Wissenschaften 1*. Frankfurt/Main
- Heidegger, M. (2006 [1926]): *Sein und Zeit*. Tübingen
- Heisenberg, W. (1927): *Über den anschaulichen Inhalt der quantenmechanischen Kinematik und Mechanik*. *Zeitschrift für Physik*, 43, 172-198.
- Heisenberg, W. (1931): *Kausalgesetz und Quantenmechanik*. *Erkenntnis (Annalen der Philosophie)*, 2, 172 -182.
- Heisenberg, W. (2007): *Physik und Philosophie*. (7. Auflage ed.). Stuttgart
- Hendrycks, D./Mazeika, M. (2022): *X-Risk Analysis for AI Research*. *arXiv2206.05862v7*, 1-36.
- Henry, M. (2011): *Inkarnation. Eine Philosophie des Fleisches*. Freiburg/München
- Hermelin, B. (2001): *Bright Splinters of the Mind: A Personal Story of Research with Autistic Savants*.

- Herrmann, S. K./Kuch, H. (2007): Verletzende Worte. Eine Einleitung. In S. K. Herrmann/S. Krämer/H. Kuch (Hrsg.), *Verletzende Worte. Die Grammatik sprachlicher Missachtung*. (S. 8-30). Bielefeld
- Hinton, G. E./Osindero, S./Teh, Y.-W. (2006): A Fast Learning Algorithm for Deep Belief Nets. *Neural Computation*, 18(7), 1527-1554.
- Hochreiter, S./Schmidhuber, J. (1997): Long Short-Term Memory. *Neural Computation*, 9(8), 1735-1780.
- Hofstadter, D. R. (1979): *Gödel, Escher, Bach: an eternal golden braid*. Hassocks
- Holterman, B./van Deemter, K. v. (2023): Does ChatGPT have Theory of Mind? *arXiv:2305.14020v1*.
- Hongladarom, S. (2020): *The Ethics of AI and Robotics: A Buddhist Viewpoint*. Lanham, Maryland
- Horner, V./Whiten, A. (2005): Causal knowledge and imitation/emulation switching in chimpanzees (*Pan troglodytes*) und children (*Homo sapiens*). *Animal Cognition*, 8, 164-181.
- Hossenfelder, S. (2018): *Das hässliche Universum. Warum die Suche nach Schönheit die Physik in die Sackgasse führt*. Berlin
- Houston, E. J. (1893): *The Electric Transmission of Intelligence And Other Advanced Primers of Electricity*. London
- Hubinger, E./van Merwijk, C./Mikulik, V./Skalse, J./Garrabrant, S. (2021): *Risks from Learned Optimization in Advanced Machine Learning Systems*. *arXiv:1906.01820v3*
- Husserl, E. (1966): *Zur Phänomenologie des inneren Zeitbewußtseins (1883-1917)*. Husserliana X. Den Haag
- Husserl, E. (2000): *Vorlesungen zur Phänomenologie des inneren Zeitbewusstseins*. Herausgegeben von Martin Heidegger. Tübingen
- Illouz, E. (2007): *Der Konsum der Romantik*. Frankfurt/Main
- Illouz, E. (2011): *Die Errettung der modernen Seele*. Frankfurt/Main
- Ishiguro, K. (2021): *Klara und die Sonne*. Aus dem Englischen von Barbara Schaden. München
- Jacques, V./Wu, E. et al. (2007): Experimental Realization of Wheeler's Delayed-Choice Gedanken Experiment. *Science*, 315(5814), 966-968.
- Jäger, L. (2013): Zur Leiblichkeit der Sprache. Phylogenetische Reminiszenzen in systematischer Absicht. In E. Alloa/M. Fischer (Hrsg.), *Leib und Sprache. Zur Reflexivität verkörperter Ausdrucksformen*. (S. 56-76). Weilerswist
- James, W. (1890): *The Principles of Psychology*. New York; London
- James, W. (1904): Does 'Consciousness' Exist? *The Journal of Philosophy, Psychology and Scientific Methods* 1(18), 477-491.
- Jansen, S. A. (2023): *Humanismus 5.0. Das postdigitale Zeitalter des Menschen Eine Bricolage*. In P. Kovce/B. Priddat (Hrsg.), *Selbstverwandlung. Das Ende des Menschen und seine Zukunft. Anthropologische Perspektiven von Anthropologisierung und Individualisierung*. Marburg
- Jullien, F. (2002): *Der Umweg über China. Ein Ortswechsel des Denkens*. Berlin
- Kaehr, R. (1993): *Disseminatorik: Zur Logik der ‚Second Order Cybernetics‘*.

- Von den ‚Laws of Form‘ zur Logik der Reflexionsform. In D. Baecker (Hrsg.), *Kalkül der Form*. (S. 152-196). Frankfurt/Main
- Kant, I. (1784): *Beantwortung der Frage: Was ist Aufklärung?* Berlinische Monatsschrift (Dezember), 481-494.
- Kauffman, L. H. (2019): Paper Computers, Imaginary Values and the Emergence of Fermions. *Cybernetics & Human Knowing*, 26(2-3), 107-160.
- Kelly, K. (1992): *Out of Control: The New Biology of Machines, Social Systems, and the Economic World*. New York
- Kleve, H. (2011): *Aufgestellte Unterschiede. Systemische Aufstellung und Tezalemma in der Sozialen Arbeit*. Heidelberg
- Koch, C. (2020): *Kognition – Was ist Bewusstsein? Spektrum der Wissenschaft*, 2020(2), 12-17.
- Kochen, S. (2017): A Reconstruction of Quantum Mechanics. In R. Bertlmann/A. Zeilinger (Hrsg.), *QUANTUM [UN]SPEAKABLES II. Half a Century of Bell's Theorem* (S. 201-235). Berlin [u.a.]
- Kohlberg, L. (1996): *Die Psychologie der Moralentwicklung*. Frankfurt/Main
- Kojima, T./Gu, S. S./Reid, M./Matsuo, Y./Iwasawa, Y. (2022): Large Language Models are Zero-Shot Reasoners. arXiv2205.11916.
- Korzybski, A. (1958 [1933]): *Science and Sanity. An Introduction to non-Aristotelian Systems and General Semantics*. Institute for General Semantics, Lakerville 2000 [1933]. (5. Auflage ed.). Brooklyn, N.Y.
- Kosinski, M. (2023): Theory of Mind May Have Spontaneously Emerged in Large Language Models. arXiv:2302.02083v3.
- Kreye, A. (2023): »Menschen werden weniger bedeutend sein«. Interview mit Jürgen Schmidhuber, <https://www.sueddeutsche.de/kultur/juergen-schmidhuber-ki-pionier-1.6051518> [letzter Zugriff: 05.09.2023]
- Latour, B. (2002): *Wir sind nie modern gewesen. Versuch einer symmetrischen Anthropologie*. (2 ed.). Frankfurt/Main
- Latour, B. (2009): *Das Parlament der Dinge: Für eine politische Ökologie*. Frankfurt/Main
- Latour, B. (2014): *Existenzweisen. Eine Anthropologie der Modernen*. Berlin
- Laucken, U. (2003): *Theoretische Psychologie. Denkformen und Sozialpraxen*. Oldenburg
- Laughlin, R. B. (2007): *Abschied von der Weltformel. Die Neuerfindung der Physik*. München
- Laughlin, R. B. (2008): Schrödingers Problem. Oder: was bei der Erfindung der Quantenmechanik nicht logisch zuende gedacht wurde. In H. U. Gumbrecht/R. P. Harrison/M. R. Hendrickson/R. B. Laughlin (Hrsg.), *Geist und Materie – Was ist Leben? Zur Aktualität von Erwin Schrödinger*. (S. 44-56). Frankfurt/Main
- LeCun, Y. (2022). A Path Towards Autonomous Machine Intelligence, Version 0.9.2, 2022-06-27. Retrieved from <https://openreview.net/pdf?id=BZ5a1r-kVsf>
- LeCun, Y./Bengio, Y./Hinton, G. (2015): Deep Learning. *Nature* 521.
- Leder, D. (1990): *The Absent Body*. Chicago
- LeDoux, J. (1994): Das Gedächtnis für Angst. *Spektrum der Wissenschaft*, 8, 76-83.

- LeDoux, J. (2021): *Bewusstsein. Die ersten vier Milliarden Jahre*. Stuttgart
- Lem, S. (1976): *Summa technologiae*. Frankfurt/Main
- Lemoine, B./da Silva, G. (2022, 16.7.2022): »Ich glaube, ich bin in meinem Innersten ein Mensch. Auch wenn meine Existenz in der virtuellen Welt stattfindet«: eine künstliche Intelligenz über ihr Innenleben. *Neue Züricher Zeitung*. Retrieved from <https://www.nzz.ch/technologie/ich-glaube-ich-bin-in-meinem-innersten-ein-mensch-auch-wenn-meine-existenz-in-der-virtuellen-welt-stattfindet-eine-kuenstliche-intelligenz-ueber-ihr-innenleben-ld.1692455> [Abruf: 5.4.2023].
- Lenzen, M. (2022): *Das Öffnen der Black Box*. GEOkompakt, 71, 74-79.
- Lévinas, E. (1984): *Die Zeit und der Andere*. Hamburg
- Lévinas, E. (1998): *Jenseits des Seins oder anders als Sein geschieht*. Freiburg, München
- Li, K. (2023): Do Large Language Models learn world models or just surface statistics? *The Gradient*, <https://thegradient.pub/othello/> [Abruf 6.7.2023].
- Libet, B./Wright, E. W./Feinstein, B./Pearl, D. (1979): Subjective referral of the timing for a conscious experience: A functional role for the somatosensory specific projection system in man. *Brain*, 102, 191-222.
- Liu, Y./Deng, G./Xu, Z. et al. (2023): Jailbreaking ChatGPT via Prompt Engineering: An Empirical Study. arXiv:2305.13860v1.
- Luhmann, N. (1984): *Soziale Systeme. Grundriß einer allgemeinen Theorie*. Frankfurt/Main
- Luhmann, N. (1986): *Ökologische Kommunikation. Kann die moderne Gesellschaft sich auf ökologische Gefährdungen einstellen?* Opladen
- Luhmann, N. (1992): *Beobachtungen der Moderne*. Opladen
- Luhmann, N. (1993a): *Soziale Systeme. Grundriß einer allgemeinen Theorie*. (4. Auflage ed.). Frankfurt/Main
- Luhmann, N. (1993b): »Was ist der Fall?« und »Was steckt dahinter?« Die zwei Soziologien und die Gesellschaftstheorie. *Zeitschrift für Soziologie*, 22(4), 245-260.
- Luhmann, N. (1995): Die Autopoiesis des Bewußtseins. In N. Luhmann (Hrsg.), *Die Soziologie und der Mensch*. (Vol. Soziologische Aufklärung, S. 55-112). Opladen
- Luhmann, N. (1996): *Die Realität der Massenmedien*. Opladen
- Luhmann, N. (1998a): *Die Gesellschaft der Gesellschaft*. Frankfurt/Main
- Luhmann, N. (1998b): *Die Wissenschaft der Gesellschaft*. (3. Auflage). Frankfurt/Main
- Luhmann, N. (1998c): *Gesellschaftsstruktur und Semantik. Studien zur Wissenssoziologie der Gesellschaft*. Band 1. Frankfurt/Main
- Luhmann, N. (2000): *Die Religion der Gesellschaft*. Frankfurt/Main
- Luhmann, N. (2008): *Ideenevolution*. Frankfurt/Main
- Luhmann, N. (2009): Sozialisation und Erziehung. In N. Luhmann (Hrsg.), *Soziologische Aufklärung 4. Beiträge zur funktionalen Differenzierung der Gesellschaft* (S. 183-192). Wiesbaden
- Lullus, R. (1999): *Ars brevis*. Hamburg

- Malin, S. (2006): *Dr. Bertelsmanns Socken. Wie die Quantenphysik unser Weltbild verändert*. Reinbek bei Hamburg
- Manyika, J. (2022): *Getting AI Right: Introductory Notes on AI & Society*. *Daedalus* (2022) 151 (2): 5–27. *Daedalus* 151, 5–27.
- Marchetti, A./Di Dio, C./Cangelosi, A./Manzi, F./Massaro, D. (2023): *Developing ChatGPT's Theory of Mind*. *Front. Robot. AI*, 10(10:1189525).
- Marcus, G./Davis, E. (2020). *GPT-3, Bloviator: OpenAI's language generator has no idea what it's talking about*. *MIT Technology Review*. Retrieved from <https://www.technologyreview.com/2020/08/22/1007539/gpt3-openai-language-generator-artificial-intelligence-ai-opinion/> [Abruf. 1.4.2023].
- Maturana, H. R. (1985): *Biologie der Kognition*. In H. R. Maturana (Hrsg.), *Erkennen: Die Organisation und Verkörperung von Wirklichkeit*. (2. Auflage ed., S. 32–80). Braunschweig
- Maturana, H. R./Varela, F. J. (1985): *Autopoietische Systeme: eine Bestimmung der lebendigen Organisation*. In H. R. Maturana (Hrsg.), *Erkennen: Die Organisation und Verkörperung von Wirklichkeit*. (2. Aufl. ed., S. 170–234). Braunschweig Wiesbaden
- Maturana, H. R./Varela, F. J. (1987): *Der Baum der Erkenntnis: Die biologischen Wurzeln des menschlichen Erkennens*. (2. Aufl. ed.). Bern, München
- Maturana, H. R./Verden-Zöllner, G. (1994): *Liebe und Spiel: Die vergessenen Grundlagen des Menschseins*. (2. Aufl. ed.). Heidelberg
- McCarthy, J./Minsky, M. L./Rochester, N./Shannon, C. E. (2006 [1955]): *A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence*, August 31. *AI Magazine*, 27, 12–14.
- McNeill, D./Freiberger, P. (1994): *Fuzzy Logic. Die »unscharfe« Logik erobert die Technik*. München
- Mead, G. H. (1967): *Mind, Self, and Society from the Standpoint of a Social Behaviorist*. Chicago
- Merleau-Ponty, M. (1974 [1966]): *Phänomenologie der Wahrnehmung*. Berlin
- Merleau-Ponty, M. (2004): *Das Sichtbare und das Unsichtbare*. (3. Auflage ed.). München
- Metzinger, T. (1998): *Anthropologie und Kognitionswissenschaft*. In P. Gold/A. K. Engel (Hrsg.), *Der Mensch in der Perspektive der Kognitionswissenschaft*. (S. 326–372). Frankfurt/Main
- Metzinger, T. (1999): *Subjekt und Selbstmodell* (2. Auflage ed.). Paderborn
- Metzinger, T. (2001): *Postbiotisches Bewußtsein: Wie man ein künstliches Subjekt baut und warum wir es nicht tun sollten*. In H.-N. Museum (Hrsg.), *Paderborner Podium zum Thema »Computer, Gehirn und Bewusstsein«*, 24.10.2001 (S. 87–113). Paderborn
- Metzinger, T. (2004): *Being No One. The Self-Model Theory of Subjectivity*. Cambridge, Massachusetts
- Metzinger, T. (2021): *Artificial Suffering: An Argument for a Global Moratorium on Synthetic Phenomenology*. *Journal of Artificial Intelligence and Consciousness*, 8 (1), S. 43–66.
- Mill, J. S. (2014 [1859]): *On Liberty*. Scotts Valley

- Minsky, M. (1990): *Mentopolis*. Stuttgart
- Minsky, M./Papert, S. (1969): *Perceptrons: An Introduction to Computational Geometry*. Cambridge MA
- Mitchell, M. (2019): *Artificial Intelligence: A Guide for Thinking Humans*. New York
- Mökander, J./Schuett, J./Kirk, H. R./Floridi, L. (2023): Auditing large language models: a three-layered approach. arXiv: 2302.08500, 1-29.
- Müller, O. (2022): Maschinelle Alterität. Philosophische Perspektiven auf Begegnungen mit künstlicher Intelligenz. In M. Schnell/L. Nehlsen (Hrsg.), *Begegnungen mit künstlicher Intelligenz. Intersubjektivität, Technik, Lebenswel* (S. 23-47). Weilerswist
- Nagarjuna. (2010): *Nagarjuna. Die Lehre von der Mitte: Mula-madhyamakakarika*. Aus dem chinesischen Text übersetzt von Lutz Geldsetzer. Hamburg
- Nagel, T. (1994): *Wie ist es, eine Fledermaus zu sein?* In M. Frank (Hrsg.), *Analytische Theorien des Selbstbewusstseins*. Frankfurt/Main
- Nakano, R./Hilton, J./Balaji, S./Ouyang, L./Kim, C./Hesse, C., . . . Schulman, J. (2021): *WebGPT: Browser-assisted question-answering with human feedback*. <https://deepai.org/publication/webgpt-browser-assisted-question-answering-with-human-feedback> [Abruf 10.7.2023].
- Nake, F. (2008): *Surface, Interface, Subface. Three Cases of Interaction and One Concept*. In U. Seifert/J. H. Kim/A. Moore (Hrsg.), *Paradoxes of Interactivity. Perspectives for Media Theory, Human-Computer Interaction, and Artistic Investigations* (S. 92-109). Bielefeld
- Nassehi, A. (2011): *Gesellschaft der Gegenwart* (1. Aufl. ed.). Berlin
- Nassehi, A. (2015): *Die »Theodizee des Willens« als Bezugsproblem des Ethischen*. In A. Nassehi/I. Saake/J. Siri (Hrsg.), *Ethik – Normen – Werte* (S. 13-41). Wiesbaden
- Nassehi, A. (2021): *Muster: Theorie der digitalen Gesellschaft*. München
- Neumann, J. v. (1932): *Mathematische Grundlagen der Quantentheorie*. Heidelberg, Berlin
- Newell, A./Shaw, J. C./Simon, H. A. (1959): *Report on a general problem-solving program*. *Proceedings of the International Conference on Information Processing*, 256-264.
- Ngo, R./Chan, L./Mindermann, S. (2022): *The alignment problem from a deep learning perspective*. arXiv:2209.00626v4.
- Nietzsche, F. (2021 [1887]): *Die fröhliche Wissenschaft*. Hamburg
- Nietzsche, F. W. (2009 [1887]): *Zur Genealogie der Moral*. Stuttgart
- Nishitani, K. (1986): *Was ist Religion?* Frankfurt/Main
- Nyanatiloka. (1981): *Der Weg zur Erlösung*. (2. Auflage ed.). Konstanz
- Olausson, T. X./Priya, I. et al. (2023): *Demystifying GPT Self-Repair for Code Generation*. arXiv:2306.09896v3.
- Paál, G./Clas, D./Grieser, C. (2017). *»Ich rede Unsinn«*. *Der Logiker George Spencer-Brown*. Detlef Clas / Charlotte Grieser SWR2 Wissen. <https://www.swr.de/swr2/programm/download-swr-13612.pdf> [Abruf: 15.3.2023].
- Pierce, C. S. (1991): *Schriften zum Pragmatismus und Pragmatizismus*. Herausgegeben von Karl-Otto Apel. Frankfurt/Main

- Peng, Z./Wang, W./Dong, L./al., e. (2023): Kosmos-2: Grounding Multimodal Large Language Models to the World. arXiv:2306.14824v3.
- Penrose, R. (1989): *The Emperor's New Mind. Concerning Computers, Minds and the Laws of Physics.* Oxford
- Petzold, M. (1988): *Indische Psychologie. Eine Einführung in traditionelle Ansätze und moderne Forschung.* Landsberg
- Piaget, J. (1992): *Einführung in die genetische Erkenntnistheorie.* (5. Auflage ed.). Frankfurt/Main
- Planck, M. (1958): *Physikalische Abhandlungen und Vorträge.* Band 1-3. Braunschweig
- Pospeschill, M. (2004): *Konnektionismus und Kognition. Eine Einführung.* Stuttgart
- Priddat, B. (2023): *Wie werden wir lernen, uns zu Automaten zu verhalten? Überlegungen zur Zivilisationsentwicklung der Mensch-Maschine-Hybride.* In P. Kovce/B. Priddat (Hrsg.), *Selbstverwandlung. Das Ende des Menschen und seine Zukunft.* (S. 297-328). Marburg
- Prigogine, I. (1979): *Vom Sein zum Werden.* München, Zürich
- Prigogine, I. (1997): *Zeit, Chaos und Naturgesetze.* In A. Gimmler/M. Sandbothe/W. C. Zimmerli (Hrsg.), *Die Wiederentdeckung der Zeit. Reflexionen, Analysen, Konzepte* (S. 79-94). Darmstadt
- Prigogine, I./Nicolis, G. (1987): *Die Erforschung des Komplexen. Auf dem Weg zu einem neuen Verständnis der Naturwissenschaften.* München, Zürich
- Radford, J./Wu, R./Child, D./Luan, D./Amodei/Sutskever, I. (2019): *Language models are unsupervised multitask learners.*
- Rae, J. W./Borgeaud, S./Cai, T. et al. (2021): *Scaling language models: Methods, analysis & insights from training Gopher.* arXiv(2112.11446).
- Ramachandran, V. (2005): *Eine kurze Reise durch Geist und Gehirn.* Reinbek
- Ramesh, A./Pavlov, M./Goh, G. et al. (2021): *Zero-Shot Text-to-Image Generation.* arXiv:2102.12092v2.
- Raunak, V./Menezes, A./Junczys-Dowmunt, M. (2021): *The Curious Case of Hallucinations in Neural Machine Translation.* arXiv(2104.06683).
- Rentsch, T. (2010): *Transzendenz und Negativität. Religionsphilosophische und ästhetische Studien.* Berlin
- Rizzolatti, G./Fogassi, L./Gallese, V. (2006): *Mirrows in the Mind.* *Scientific American*, 295(5), 54-61.
- Rorty, R. (1987 [1979]): *Der Spiegel der Natur. Eine Kritik der Philosophie.* Frankfurt/Main
- Rosenblatt, F. (1958): *The perceptron: A probabilistic model for information storage and organization in the brain.* *Psychological Review*, 65(6), 386-408.
- Rössler, O. E. (1992): *Endophysik. Die Welt des inneren Beobachters.* Berlin
- Roth, G. (2003): *Fühlen, Denken, Handeln. Wie das Gehirn unser Verhalten steuert.* Frankfurt/Main
- Rovelli, C. (2014): *Why do we remember the past and not the future? The 'time oriented coarse graining' hypothesis.* arXiv:, 1407.3384.
- Rumelhart, D. E./Hinton, G. E. W./Williams, R. J. (1986): *Learning representations by back-propagating errors.* *Nature* 323, 533-536.

- Russell, S. (2020): Human Compatible: Künstliche Intelligenz und wie der Mensch die Kontrolle über die superintelligenten Maschinen behält. Frechen
- Russell, S. J./Norvig, P. (2016): Artificial Intelligence: A Modern Approach. 4. Auflage. Global Edition, Pearson. New York
- Sahn, S. (2019): Nur Weiß-Nicht: Gesammelte Lehrbriefe von Zen-Meister Seung Sahn. Gießen
- Sandsten, K. E./Zahavi, D./Parnas, J. (2021): Disorder of Selfhood in Schizophrenia: A Symptom or a Gestalt? *Psychopathology*, 2022 Mar 29:1-9.
- Sapolsky, R. (2017): Gewalt und Mitgefühl. Die Biologie des menschlichen Verhaltens. München
- Schärtl, T. (2019): Allmacht, Allwissenheit und Allgüte. In K. Viertbauer/G. Gasser (Hrsg.), *Handbuch Analytische Religionsphilosophie* (S. 124-143). Stuttgart
- Schelsky, H. (1957): Ist die Dauerreflektion institutionalisierbar?: Zum Thema einer modernen Religionssoziologie. *Zeitschrift Für Evangelische Ethik*, 1(1), 153-174.
- Schleiffer, R. (2012): Das System der Abweichungen. Eine systemtheoretische Neubegründung der Psychopathologie. Heidelberg
- Schlosshauer, M. (2005): Decoherence, the measurement problem, and interpretations of quantum mechanics. *Eletronische Ressource: arXiv:quant-ph/0312059v4* 28 Jun 2005, 1-41.
- Schmitt, C. (2009 [1922]): Politische Theologie. Vier Kapitel zur Lehre von der Souveränität. Berlin
- Schnell, A. (2002): Das Problem der Zeit bei Husserl. Eine Untersuchung über die husserlschen Zeitdiagramme. *Husserl Studies*, 18, 89-122.
- Schrittwieser, J./Antonoglou, I./Hubert, T./al., e. (2020): Schrittwieser, Julian; Ioannis Antonoglou und Thomas Hubert, et al. 2020. Mastering Atari, Go, chess and shogi by planning with a learned model. *Nature* 588, 604-609.
- Schrödinger, E. (1935): Die gegenwärtige Situation in der Quantenmechanik. *Die Naturwissenschaften*, 23(48), 807-812, 823-828, 844-849.
- Schrödinger, E. (1946): Was ist Leben? Die lebende Zelle mit den Augen eines Physikers betrachtet. Bern
- Schrödinger, E. (1959): Geist und Materie. Zürich
- Schulz von Thun, F. (2010): Das Werte- und Entwicklungsquadrat: Ein Werkzeug für Kommunikationsanalyse und Persönlichkeitsentwicklung. Theorie und Praxis der Sozialpädagogik, 9, 13-17.
- Schütz, A./Luckmann, T. (2003): Strukturen der Lebenswelt. Konstanz
- Searle, J. R. (1980): Minds, Brains, and Programs. *The Behavioral and Brain Sciences*, 3(3), 417-457.
- Seising, R. (2021): Es denkt nicht: die vergessenen Geschichten der KI. Leipzig
- Sejnowski, T. J. (2023): Large Language Models and the Reverse Turing Test. *Neural Computation*, 35(3), 309-342.
- Selfridge, O. G. (1958): Pandemonium: a paradigm for learning. *Mechanisation of Thought Processes: Proceedings of a Symposium Held at the National Physical Laboratory*, 513-526.

- Shannon, C. E. (1948): A Mathematical Theory of Communication. *Bell System Technical Journal*, 27(3), 379-423.
- Shannon, C. E./MacCarthy, J. E. (1956): *Automata Studies*. Princeton
- Shapiro, D. K. (2023): *Benevolent by Design. Six Words to Safeguard Humanity*. New York
- Shen, X./Chen, Z./Backes, M./Shen, Y./Zhang, Y. (2023): »Do Anything Now«: Characterizing and Evaluating In-The-Wild Jailbreak Prompts on Large Language Models. arXiv:2308.03825v1.
- Simon, F. B. (2018): *Formen. Zur Koppelung von Organismus, Psyche und sozialen Systemen*. Heidelberg
- Singer, W. (2002): *Der Beobachter im Gehirn. Essays zur Hirnforschung*. Frankfurt/Main
- Snyder, A. W. (2001): Paradox of the savant mind. *Nature*, 413, 251-252.
- Socher, R. (2012): Semantic Compositionality through Recursive Matrix-Vector Spaces. *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, 1201-1211.
- Sparrer, I./Kibét, M. V. v. (2000): *Ganz im Gegenteil. Tetralemmaarbeit und andere Grundformen systemischer Strukturaufstellungen - für Querdenker und solche, die es werden wollen*. Heidelberg
- Spencer Brown, G. (1997): *Laws of Form. Gesetze der Form*. Lübeck
- Spencer Brown, G. (2013): *Dieses Spiel geht nur zu zweit. Aus dem Englischen von Andreas Baar*. (3. Auflage ed.). Lübeck
- Spengler, O. (2007 [1923]): *Der Untergang des Abendlandes*. Düsseldorf
- Srivastava, A./Rastogi, A./Rao, A./Shoeb, A. A. M./Abid, A./Fisch, A./al., e. (2022): Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models. arXiv:2206.04615v3 1-100.
- Star, S. L. (2015): Die Struktur schlecht strukturierter Lösungen. *Grenzbjekte und heterogenes verteiltes Problemlösen. Navigationen - Zeitschrift für Medien- und Kulturwissenschaften*, 15(1), 57-77.
- Stephan, A., Walter, S. et al. (2021): Nachwort. In: Turing, Alan M.: *Computing Machinery and Intelligence / Können Maschinen denken?* Übers. und hrsg. von Achim Stephan und Sven Walter, Stuttgart: Reclam, S. 131-202.
- Steinbicker, J. (2011): *Zur Theorie der Informationsgesellschaft. Ein Vergleich der Ansätze von Peter Drucker, Daniel Bell und Manuel Castells* (2. Auflage ed.). Wiesbaden
- Sutskever, I./Vinyals, O./Le, Q. V. (2014): Sequence to sequence learning with neural networks. In Z. Ghahramani/M. Welling et al. (Hrsg.), *Advances in Neural Information Processing Systems 2* (Vol. 27, S. 1-9)
- Taleb, N. N. (2018): *Antifragilität: Anleitung für eine Welt, die wir nicht verstehen* (3. Auflage). München
- Tegmark, M. (2017): *Leben 3.0. Mensch sein im Zeitalter Künstlicher Intelligenz*. Berlin
- Tegmark, M./Wheeler, J. A. (2001): *100 Jahre Quantentheorie. Spektrum der Wissenschaft*, 2001/4, 68-76.
- Thoma, S. (2022): *Im Offenen. Henri Maldineys Philosophie der Psychosen*. Wien, Berlin

- Thomas, I. S./Wang, J./GPT₃. (2022): Was euch zu Menschen macht – Antworten einer künstlichen Intelligenz auf die großen Fragen des Lebens. München
- Thomas, W. I./Thomas, D. S. (1928): *The child in America: behavior problems and programs*. New York
- Thompson, A. D. (2023). Integrated AI: The sky is entrancing (mid-2023 AI retrospective). June 2023. lifearchitect.ai. <https://lifearchitect.ai/the-sky-is-entrancing/> [Abruf: 9.8.2023]
- Thoppilan, R. et al. (2022): LaMDA: Language Models for Dialog Applications. arXiv:2201.08239.
- Thurman, R. A. F. (2003): *Das tibetische Totenbuch oder das große Buch der natürlichen Befreiung durch Verstehen im Zwischenzustand*. Neu übersetzt und kommentiert von Robert A. F. Thurman. (2. Auflage). Frankfurt
- Tomasello, M. (2009): *Die Ursprünge der menschlichen Kommunikation*. Berlin
- Tononi, G. (2004): An information integration theory of consciousness. *BMC Neuroscience* 5.
- Trivers, R. (2011): *The Folly of Fools: The Logic of Deceit and Self-Deception in Human Life*. New York
- Tugendhat, E. (2006): *Egozentrik und Mystik. Eine anthropologische Studie*. München
- Tugendhat, E. (2007): *Anthropologie statt Metaphysik*. München
- Turing, A. (2004): *Intelligent Machinery (1948)*. In B. J. Copeland (Hrsg.), *The Essential Turing* (S. 0)
- Turing, A. M. (1937): On Computable Numbers, with an Application to the Entscheidungsproblem. *Proceedings of the London Mathematical Society*, 52-42(1), 230-265.
- Turing, A. M. (1950): Computing Machinery and Intelligence. *Mind*, LIX, 433-460.
- Turing, A. M. (2004 [1948]): *Intelligent machinery*. In J. B. Copeland (Hrsg.), *The Essential Turing*. (S. 395-432). Oxford
- Turing, A. M. (2021): *Computing Machinery and Intelligence / Können Maschinen denken? Englisch/Deutsch*. Übers. und hrsg. von Achim Stephan und Sven Walter. Stuttgart
- Tversky, A./Daniel, K. (1983): Extensional versus intuitive reasoning: The conjunction fallacy in probability judgement. *Psychological Review*, 90(4), 293-315.
- Uexküll, T. v. (2003): *Psychosomatische Medizin. Modelle ärztlichen Denkens und Handelns*. (6. Aufl. ed.). München, Jena
- Varela, F. J. (1979): The Extended Calculus of Indications interpreted as a three-valued logic. *Notre Dame Journal of Formal Logic*, 20(1), 141-146.
- Varela, F. J. (1999): The Specious Present: A Neurophenomenology of Time Consciousness. In J. Petitot/F. J. Varela/B. Pachoud (Hrsg.), *Naturalizing Phenomenology: Issues in Contemporary Phenomenology and Cognitive Science*. (S. 266-314). Stanford, California
- Varela, F. J./Thompson, E. (2003): Neural Synchronity and the Unity of Mind: A Neurophenomenological Perspective. In A. H. Cleeremans (Hrsg.),

- The Unity of Consciousness: Binding, Integration and Dissociation. (S. 266-287). New York
- Vaswani, A./Shazeer, N./Parmar et al. (2017): Attention Is All You Need. arXiv:(1706.03762), 1-15.
- Vogd, W. (2004): Ärztliche Entscheidungsprozesse des Krankenhauses im Spannungsfeld von System- und Zweckrationalität: Eine qualitativ rekonstruktive Studie. Berlin
- Vogd, W. (2006a): Wer entscheidet, wer entscheidet? Fragen zur Theorie und Empirie zugerechneter Intentionalität. In J. Reichertz/N. Zaboura (Hrsg.), Akteur Gehirn – oder das vermeintliche Ende des handelnden Subjekts. Eine Kontroverse. (S. 137-155). Wiesbaden
- Vogd, W. (2010): Gehirn und Gesellschaft. Weilerswist
- Vogd, W. (2013): Konfusionen (post-)moderner Psychotherapie: verstörte Gesellschaften, verstörte Organisationen, verstörte Individuen und kein Schamane in Sicht. Verhaltenstherapie und psychosoziale Praxis, 45(1), 33-46.
- Vogd, W. (2014a): Von der Physik zur Metaphysik – eine soziologische Rekonstruktion des Deutungsproblems der Quantentheorie. Weilerswist
- Vogd, W. (2014b): Welten ohne Grund. Buddhismus, Sinn und Konstruktion. Heidelberg
- Vogd, W. (2017): Buddhistische Praxis und Sprache. In A. Lasch/W.-A. Liebert (Hrsg.), Handbuch Sprache in der Religion. (Vol. Band 18, S. 154-178). Berlin
- Vogd, W. (2018): Selbst- und Weltverhältnisse. Leiblichkeit, Polykontexturalität und implizite Ethik. Weilerswist
- Vogd, W. (2019): Der ermächtigte Meister. Eine systemische Rekonstruktion am Beispiel des Skandals um Sogyal Rinpoche. Heidelberg
- Vogd, W. (2020): Quantenphysik und Soziologie im Dialog. Betrachtungen zu Zeit, Beobachtung und Verschränkung. Berlin
- Vogd, W. (2021): Offenheit für neue Kategorien und Begründungen – warum eine Wissenssoziologie, die diesen Namen verdient, sich immerfort erneuern muss. Eine Würdigung der praxeologischen Wissenssoziologie von Ralf Bohnsack. ZQF-Zeitschrift für Qualitative Sozialforschung, 22(1), 107-122.
- Vogd, W./Feißt, M. (2022): Therapeutische Arrangements im Maßregelvollzug. Studien zur Leerstellengrammatik und den Bezugsproblemen in der forensischen Psychiatrie Wiesbaden
- Vogd, W./Harth, J. (2015): Die Praxis der Leere. Zur Verkörperung buddhistischer Lehren in Erleben, Reflexion und Lehrer-Schüler-Beziehung. Weilerswist
- Vogd, W./Harth, J. (2019): Kontexturanalyse: Eine Methodologie zur Rekonstruktion polykontexturaler Zusammenhänge. Forum Qualitative Sozialforschung / Forum: Qualitative Social Research.
- von Brück, M. (2015): Weltinnenraum. Rainer Maria Rilkes Duineser Elegien in Resonanz mit dem Buddha. Freiburg, Basel, Wien
- von Foerster, H. (1992): Ethics and Second-order Cybernetics. Cybernetics & Human Knowing 1(1), 9-19.
- von Foerster, H. (2006): Über selbst-organisierende Systeme und ihre

- Umwelten. In ders. (Hrsg.), *Sicht und Einsicht. Versuche zu einer operativen Erkenntnistheorie* (S. 115-130). Heidelberg
- von Foerster, H./Bröcker, M. (2002): *Teil der Welt. Fraktale einer Ethik. Ein Drama in drei Akten*. Heidelberg
- von Foerster, H./Pörksen, B. (2022): *Wahrheit ist die Erfindung eines Lügners. Gespräche für Skeptiker*. Heidelberg
- Waldenfels, H. (2013): *Absolutes Nichts. Zur Grundlegung des Dialogs zwischen Buddhismus und Christentum. Mit einem Vorwort von Nishitani Keiji* (4. Auflage). Paderborn
- Wang, L. J./Zou, X. Y./Mandel, L. (1991): Induced coherence without induced emission. *Physical Review A*, 44(7), 4614-4622.
- Watzlawick, P./Beavin, J. H./Jackson, D. D. (1990): *Menschliche Kommunikation. Formen, Störungen, Paradoxien*. (8. Auflage ed.). Bern, Stuttgart, Toronto
- Weber, A. (2003): *Natur als Bedeutung. Versuch einer semiotischen Theorie des Lebendigen*. Würzburg
- Wei, J./Wang, X./Schuurmans, D. et al. (2022): Chain-of-thought prompting elicits reasoning in large language models. arXiv(11903v6).
- Weidinger, L./al., e. (2021): Ethical and social risks of harm from language models. arXiv:(2112.04359).
- Weizenbaum, J. (1966): ELIZA—a computer program for the study of natural language communication between man and machine. *Communications of the Acm*, 9, 36-45.
- Weizsäcker, C. F. v. (1994): *Der Aufbau der Physik*. (3. Auflage ed.). München
- Whitehead, A. N. (1987): *Prozess und Realität. Entwurf einer Kosmologie*. Frankfurt/Main
- Wiener, N. (1948): *Cybernetics or Control and Communication in the Animal and the Machine* (2. Aufl.). Cambridge (MA)
- Wiener, N. (1960): Some Moral and Technical Consequences of Automation. *Science*, 131(3410), 1355-1358.
- Wigner, E. P. (1967): Remarks on the Mind-Body-Question. In E. P. Wigner (Hrsg.), *Symmetries and Reflections* (S. 171-184). Indiana
- Wittgenstein, L. (1963 [1922]): *Tractatus logico philosophicus*. Frankfurt/M.
- Wittgenstein, L. (1989): *Vortrag über Ethik und andere kleine Schriften*. Frankfurt/M.
- Wittgenstein, L. (1990): *Tractatus logico philosophicus; Philosophische Untersuchungen*. Leipzig
- Wygotski, L. S. (1986 [1934]): *Denken und Sprechen*. Frankfurt/M.
- Youniss, J. E. (1994): *Soziale Konstruktion und psychische Entwicklung*. Frankfurt/M.
- Zeh, J. (2007): *Spieltrieb*. Frankfurt/Main
- Zeilinger, A. (2005): *Einsteins Schleier. Die neue Welt der Quantenphysik*. (2 ed.). München
- Zeilinger, A. (2007): *Einsteins Spuk. Teleportation und weitere Mysterien der Quantenphysik*. München
- Zuse, K. (1984): *Der Computer – Mein Lebenswerk*. Landsberg am Lech

