



DiDaT Grobplanung zum Vulnerabilitätsraum (VR) 06

## **Vertrauenswürdigkeit und Zuverlässigkeit digitaler Daten und Informationen**

Sebastian Hallensleben (VDE), Roland. W. Scholz (Donau Uni Krems), Andreas Kaminski (HLRS Universität Stuttgart), Julio Lambing (VEZL)

Inputs durch Dirk Helbing (ETH Zürich), Karl-Heinz Simon (CESR Universität Kassel), Malte Reissig (IASS), Dirk Marx (btu Cottbus), Sabine Thürmel (TU München)

### **1 Gegenstand, Ziele und Leitfrage**

Fälschungen von Texten und Fotos sind nicht neu – sei es in der Werbung, für politische Manipulationen, bei Betrügereien oder für andere Zwecke. Wir haben uns darauf eingestellt, Texten und Fotos mit gesunder Skepsis zu begegnen. Das gilt insbesondere im digitalen Raum.

Für Videos konnte man dagegen bisher annehmen, dass sie tatsächlich ein reales Geschehen zeigen. „Fälschungen“ waren dort nur in engem Rahmen möglich, beispielsweise durch geschicktes Schneiden, eine falsche Zuordnung oder den Einsatz eines professionellen Filmstudios. Videos galten bisher als das weitgehend unbestechliche digitale Äquivalent des Augenscheins. Seit 2018 sind jedoch mit künstlicher Intelligenz ausgestattete Werkzeuge (v.a. Deep Fake<sup>1</sup>) verfügbar, mit denen praktisch jedermann beliebige Video- und Audioaufnahmen fälschen kann. Mit entsprechender Rechenleistung sind diese Fälschungen sogar in Echtzeit möglich, d.h. ein angeblicher Live-Fernsehaufttritt einer prominenten Persönlichkeit kann während des Programms gesteuert werden. Die öffentliche Aufmerksamkeit ruht derzeit vor allem auf im WWW veröffentlichten Videos, in denen demonstriert wird, wie

man US-Politiker Worte in den Mund legen kann. Entsprechend haben erste Staaten und Unternehmen Maßnahmen gegen die Verbreitung solcher irreführender Videos ergriffen.<sup>2</sup> Der politische Fokus überdeckt jedoch, dass die überragende Mehrzahl solcher Videos derzeit eine neue Form nicht einvernehmlicher Pornografie zeigen, bei denen computer-generierte Gesichter von Prominenten auf die Köpfe der Sexdarsteller geschnitten werden, mit Abermillionen von Zuschauern.<sup>3</sup> Die allermeisten der von solchen „morph porn“ betroffenen Personen sind öffentlich agierende Frauen. Zunehmend sind aber auch vollkommen unprominente Frauen betroffen. Die Fälschungen bergen also nicht nur das Potential für privaten Rufmord, Cyberstalking oder auch Identitätsdiebstahl eingesetzt werden zu können, mit verheerenden Folgen für die Betroffenen.<sup>4</sup> Sie demonstrieren auch ihre Zugänglichkeit und Eignung für den Einsatz in Alltagsauseinandersetzungen, etwa bei der Fälschung von rechtlichen Beweisen.

Gefälschte Videos sind also nur ein augenfälliges Symptom für die generellen neuen Möglichkeiten zur Fälschung von einflussreichen Informationen. Politisch einflussreicher waren bisher sogenannte „Shallowfakes“, also mit einfachen

Mittel gefälschte Videos.<sup>5</sup> Auch überzeugende „Fotos“ nichtexistenter Menschen sowie glaubwürdige Texte lassen sich mittlerweile mit minimalem Aufwand in großer Menge und mit zahlreichen Stellschrauben generieren<sup>6</sup>. Im Einzelhandel sind Fake Reviews ein weitverbreitetes und ernsthaftes Phänomen der Konsumententäuschung.<sup>7</sup>

Parallel zu dieser technologischen Entwicklung sinkt der Einfluss der traditionellen Massenmedien und ihrer Filter- und Verifizierungsfunktion für Informationen. Sie können nur schwer mithalten mit der Dynamik der Online-Medien, in denen Inhalte, egal, ob echt oder gefälscht, können sich rasend schnell verbreiten, teilweise gezielt vorangetrieben durch kommerzielle Dienstleister<sup>8</sup>. Die Werbewirtschaft und manche politischen Akteure haben sich bereits auf diese neuen Verbreitungsmöglichkeiten eingerichtet. Über gezielte Einflussnahmen beispielsweise der russischen *Internet Research Agency (IRA)*<sup>9</sup> sowie Wahlmanipulation durch *Cambridge Analytica*<sup>10</sup> ist ausführlich berichtet worden. In Gabun und Malaysia spielten Vorwürfe zum Einsatz von Deep Fakes eine relevante Rolle in schweren politischen Krisen.<sup>11</sup>

Eine Flut falscher Informationen hat das Potenzial, Fakten in der Wahrnehmung zu verdrängen. Dies geschieht nicht nur durch eine bewusste Entscheidung von Rezipienten, dieser oder jener Information eher zu vertrauen, sondern auch durch eine unbewusste Überlagerung bereits abgespeicherten Wissens.<sup>12</sup>

Der hier verwendete Begriff „Vertrauen“ zeigt den Kern der umrissenen Problematik an: Vertrauen in und Vertrauenswürdigkeit von Informationen sind seit jeher für Menschen eine notwendige und zugleich heikle Angelegenheit. Die Entwicklung der Informationstechnik im Zeitalter der neuen sozialen Medien und der künstlichen Intelligenz hat die generelle Zeugenschafts- und Vertrauensproblematik verschärft. (Siehe: *Begriffliche und epistemologische Einordnung von Vertrauen*). Es stellt sich die Frage: Wie können die Zuverlässigkeit digitaler Informationen sowie IT-gestützte Vertrauensinfrastrukturen in naher Zukunft in Deutschland so gestaltet werden, dass ein fakten- und wertebasierter öffentlicher, wissenschaftlicher und politischer Diskurs möglich bleibt, um eine Disruption der Grundlagen von Demokratie und Rechtsstaat zu verhindern?

### **Leitfrage und inhaltliche Abgrenzung von VR6**

Wie können die Zuverlässigkeit digitaler Informationen sowie IT-gestützte Vertrauensinfrastrukturen in naher Zukunft in Deutschland so gestaltet werden, dass ein fakten- und wertebasierter öffentlicher, wissenschaftlicher und politischer Diskurs möglich bleibt, um eine Disruption der Grundlagen von Demokratie und Rechtsstaat zu verhindern? Wie sieht eine Kombination aus sozialen und technischen Ansätzen aus, die eine Verifizierung von Fakten unterstützt? Wie kann auch künftig mündige politische Meinungsbildung ablaufen? Welche Anreizsysteme können vorgeschlagen werden, mit denen Wahrheitsfindung und -verbreitung präferiert werden? Wie lassen sich kluge Netze des Vertrauens knüpfen? Wie kann dies auf angemessene Weise technisch unterstützt werden?

## **Begriffliche und epistemologische Einordnung von Vertrauen**

Die Philosophin Annette Baier hat als Arbeitsdefinition vorgeschlagen, Vertrauen als akzeptierte Verletzbarkeit zu begreifen<sup>(1)</sup>. Folgt man diesem theoretisch einflussreichem Vorschlag, muss eine Analyse des hier untersuchten Vulnerabilitätsraums berücksichtigen, dass Menschen einerseits auf elementare Weise, die Möglichkeit verletzt zu werden, nicht vermeiden können, dass sie andererseits jedoch lernen können und müssen, auf angemessene Weise Vertrauen zu gewähren oder zu entziehen.

### **Die allgemeine Problematik von Vertrauen und Zeugenschaft**

Auf den ersten Blick mag die Annahme naheliegen, dass in modernen Gesellschaften die Bedeutung von Vertrauen abnimmt: Es könnte so scheinen, als ob Verwissenschaftlichung und Technisierung Vertrauen ablösen. Man muss und braucht anderen nicht zu vertrauen, weil man durch Technik und Wissenschaft Transparenz- und Kontrollmöglichkeiten gewonnen hat, die an die Stelle von Vertrauen treten. Doch bereits Gründerfiguren der Soziologie wie Max Weber und Georg Simmel waren anderer Meinung: Technik und Wissenschaft selbst vergrößern die *Abhängigkeit* des Einzelnen von den Leistungen *anderer*, was den Vertrauensbedarf erhöht.

Diese Abhängigkeit von anderen wird besonders deutlich und relevant im Bereich von **Information und Wissen**. Das meiste, von dem wir annehmen, dass wir es wissen, haben wir *durch Andere* erfahren: dass es Bakterien gibt und einige davon unserer Gesundheit abträglich sein können; dass die Erde keine Scheibe ist und um die Sonne zirkuliert; wie weit Stuttgart und Frankfurt entfernt sind und wann der nächste Zug abfahren soll; welche Partei bei der letzten Wahl gewonnen hat, wer sie gewählt hat und wer jetzt die Regierung stellt. Selbst wann und wo wir geboren wurden wissen wir nicht qua eigener Erkenntnis.

Die Vertrauenswürdigkeit von Informationen ist eng verwoben mit dem erkenntnisbezogenen **(epistemischen) Abhängigkeit von Anderen**. Der klassische Wissensbegriff der Philosophie geht davon aus, dass Wissen eine (a) Überzeugung ist, die (b) wahr, zusätzlich aber (c) gerechtfertigt sein muss – man muss gute Gründe für die Überzeugung haben, damit man einen Wissensanspruch erheben darf, eine Überzeugung darf nicht bloß zufällig wahr sein. Dass Andere eine mögliche Quelle und Begründung eigener Wissensansprüche sind, ruft folgende Frage auf: Wie begründen und rechtfertigen wir es (c), dass wir etwas zu wissen glauben, weil Andere es uns gesagt haben – oder anders ausgedrückt: weil Andere es bezeugt haben.

Es ist nun eine offene Frage, wie wir (c) auffassen können. Eine mögliche Begründung, die aber vieles offen lässt, ist die folgende Antwort: Wir können einem Sprecher glauben, wenn wir ihm *vertrauen*. Das damit aufgeworfene Problem können wir als Zeugenschafts- und Vertrauensproblem bezeichnen. Das allgemeine Problem lässt sich durch die Frage ausdrücken: Wann sind wir darin gerechtfertigt, den Anderen zu glauben? Wann sind wir darin gerechtfertigt, den Anderen zu vertrauen, wenn sie etwas behaupten?

### **Problemverschärfung durch die Digitalisierung der Gesellschaft**

Die ohnehin bestehende, alltäglich aufscheinende Problematik von Zeugenschaft und Vertrauen verschärft sich durch die Digitalisierung der Gesellschaft, insbesondere durch die Entwicklung und alltägliche Verbreitung des Internets. Wir erhalten erstens viel mehr Informationsangebote aus zweitens mehr möglichen Quellen, die wir nicht persönlich kennen, zu drittens Themen, bei denen wir häufig weniger Möglichkeiten haben, sie durch eigene Erfahrung kritisch zu prüfen. Zudem scheint die Digitalisierung der Kommunikation neue Möglichkeiten ihrer Manipulation zu eröffnen. Ein Beispiel hierfür ist die auf Daten und Algorithmen basierte Selektion von Nachrichten, ein anderes Beispiel die erwähnten Deep Fake Videos. Sie verschärfen die allgemeine Frage, wann wir darin gerechtfertigt sind, Anderen zu glauben. Denn wir sind in mehr Themen, die für uns relevant sind, abhängig von dem, was (relativ anonyme) Andere uns mitteilen. Und die technischen Möglichkeiten der Manipulation scheinen gewachsen zu sein.

(1) Annette Baier (1986): Trust and Antitrust; in: Ethics; Vol. 96, No. 2 (Jan., 1986); S. 231-260; S. 235

Anders gefragt: Wie lassen sich kluge Netze des Vertrauens knüpfen? Wie kann dies auf angemessene Weise technisch unterstützt werden?

**Die Fragestellung des Vulnerabilitätsraums wird also in mehrfacher Hinsicht thematisch eingrenzt:**

Selbstverständlich besteht die Möglichkeit, dass durch die Digitalisierung veranlasste Schwächung der Vertrauenswürdigkeit von Informationen auch ökonomische Effekte als Unseens hat, die relevante Auswirkungen auf einzelne Geschäftsbranchen oder auf Konsumgewohnheiten haben. (Siehe die Frage der Vertrauenswürdigkeit von Produktreviews.) Zur pragmatischen Beschränkung des Untersuchungsfelds sollte jedoch die Gefährdung von Demokratie und Rechtsstaat im Vordergrund stehen.

Auch eine geografische Eingrenzung ist aus pragmatischen Gründen geboten. Die bisher aufgezeigte Problemlage als auch die im Abschnitt 2 dargelegten Vulnerabilitäten sowie die darauf antwortenden, möglichen Lösungsansätze bestehen zwar global bzw. international, zumindest insoweit als Onlineinformationen und Onlinediskurs zugänglich sind und bereits eine signifikante Rolle in Politik und Gesellschaft spielen (dies schließt lediglich einige wenige geografisch abgelegene Räume aus).

Angeichts der sich in der Stakeholder-Analyse (Abschnitt 3) abzeichnenden Akteurskonstellation, der Sprachabhängigkeit von Onlineinformationen sowie der Jurisdiktionsbezogenheit von Lösungsansätzen sollte zunächst Deutschland, fallweise auch die DACH-Region betrachtet werden. Im Verlauf des Pro-

jekts kann und soll jedoch herausgearbeitet werden, inwieweit sich Lösungen auch als Impulse bzw. Piloten für Europa (im Sinne der Europäischen Union) eignen.

Drittens soll hier nicht jede Form von Fälschung bei der Übermittlung von Informationen betrachtet werden. Kriminelle Angriffe auf die Datenübertragung bei Geldüberweisungen, der Übermittlung von Steuerbefehlen in der öffentlichen Stromversorgung oder der Koordination von öffentlichen und privaten Transport können sicherlich eine gewichtige Gefahr für eine Gesellschaft darstellen. Sie unterscheiden sich jedoch insofern von der hier geschilderten Thematik, als dass sie eine illegitime Manipulation der Übermittlungskanäle beinhalten und bestehende Informationen in der Übertragung durch falsche ersetzen. Solche Angriffe auf Datenübermittlungswege sind jedoch etwas anderes als Angriffe durch gefälschte Informationen, die erst **durch ihre Wahrnehmung und ihre diskursive Verhandlung in sozialen Räumen ihre schädliche Wirkung erhalten.**

Andererseits liegt der Fokus von VR6 auf der Information, unabhängig davon, wo sie im digitalen Raum vorliegt. In Abgrenzung zu VR05 betrachtet dieser Vulnerabilitätsraum also nicht (oder nur am Rande) die Funktion und Struktur spezifischer Bereiche des digitalen öffentlichen Raums (z.B. Social Media, klassische Nachrichtenmedien, Plattformen für nutzergenerierte Inhalte oder Fachpublikationen), sondern abstrahiert zur Information an sich sowie deren ursprünglicher Quelle.

## 2 Welche nicht intendierten, unbeabsichtigten Nebenfolgen sind von Interesse und warum?

### Gefährdung des demokratischen Gemeinwesens als Konsequenz

Wenn selbst Videos überzeugend gefälscht werden können und damit die letzte Bastion der Tatsachenprüfung durch Augenschein fällt, dann sind sämtliche Onlineinhalte fragwürdig. **Wahrheit und Lüge werden ununterscheidbar.** Nicht nur gefälschte Informationen können für echt gehalten werden, sondern auch echte Informationen können in Zweifel gezogen werden. Der Politiker, der bei der Annahme von Bestechungsgeldern gefilmt wurde, kann solche Videobelege künftig problemlos als Fälschung abtun. Aviv Ovdya hat für diese Entwicklung den Begriff „Infokalypse“ geprägt,<sup>13</sup> der inzwi-

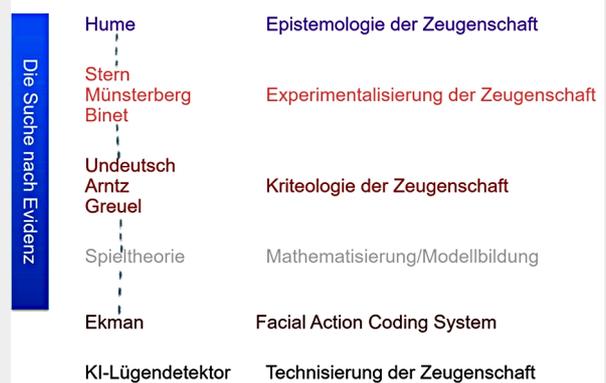
schen auch in Form einzelner Zwischenrufe an die breitere Öffentlichkeit gelangt ist.<sup>14</sup>

Traditionelle Massenmedien können zwar weiterhin versuchen, als Filter und Prüfer zu agieren, sind aber der emotionalen Wirkung und der rasend schnellen Verbreitung einer gefälschten „Nachricht“ gegenüber machtlos. Sie können ihre Filterfunktion im extrem beschleunigten Online-Nachrichtenfluss nur noch schwer oder verzögert ausüben. Gleichzeitig muss der Mensch weiterhin die Einschätzung der Glaubwürdigkeit von Informationen teilweise delegieren, da die Masse an Information ansonsten nicht zu bewältigen wäre und der einzelne Mensch nicht so viele sorgfältige Einzelbeurteilungen und -entscheidungen treffen kann.

### Wenn sich Wahrheit und Lüge für den Einzelnen aber nicht mehr mit

#### Technische Lösungen zu Vertrauen im Spiegel der Theoriegeschichte

Technische Lösungen, die Vertrauenswürdigkeit messen und zwischenmenschliche Vertrauensprozesse nachbauen und ersetzen sollen, lassen sich vor dem Hintergrund der allgemeinen Theoriegeschichte von Vertrauen deuten. Es gibt aktuell zwei große Theorielinien, die einander widersprechende Antworten geben auf die Frage, wie begründet werden kann, dass man einer Person glaubt. Die eine Antwort lautet: Vertrauen *hat* epistemische Gründe. Die andere geht dagegen davon aus, dass Vertrauen nicht auf Gründe zurückgeführt werden kann, sondern selbst ein Orientierung gebender Grund *ist*, etwas für wahr zu halten. Die erste Linie versucht die Zeugenschaft also ganz auf ein *erkenntnistheoretisches* Problem zurückzuführen. Damit geht die Hoffnung einher, das vermeintliche Wissen, das wir durch *Andere* gewinnen, auf unser *eigenes* Wissen zurückzuführen. Wir beurteilen und bewerten die Vertrauenswürdigkeit aufgrund von ‘track records’ oder Indikatoren, die wir aufgrund unserer eigenen Erfahrung gewonnen haben. Diese Linie beginnt bei David Hume, führt über Psychologen wie William Stern oder Hugo Münsterberg bis zu den technischen Entwicklungen der jüngsten Zeit. Der Grundgedanke der so genannten *Evidential View* ist: Man muss Evidenz gewinnen, auf deren Grundlage man sich dann für oder gegen Vertrauen entscheidet. Die Evidenz zeigt an, ob eine Person oder Information vertrauenswürdig ist. Der Grad der Evidenz führt zu einem Grad an Wahrscheinlichkeit der Vertrauenswürdigkeit.



**realistischem Aufwand trennen lassen, dann ist das Konstrukt des „mündigen Bürgers“ bzw. einer „mündigen Gesellschaft“ als Ausgangspunkt aller staatlichen Gewalt** (nach Art. 20, Abs. 2 des deutschen Grundgesetzes) **faktisch unmöglich geworden**. Das demokratische Gemeinwesen verliert seine Grundlage. Die Auseinandersetzung mit der Vertrauenswürdigkeit von Informationen im digitalen Raum wird so zur dringlichen Notwendigkeit.

Neben der Gefährdung demokratischer Systeme durch die informationelle Entmündigung des Bürgers sind weitere schwerwiegende Folgen realistisch: Die Zuverlässigkeit polizeilicher Untersuchungen oder die rechtsstaatliche Gültigkeit von Gerichtsprozessen kann möglicherweise durch die ausgeweitete Fälschbarkeit von Informationen nur noch eingeschränkt gewährleistet werden. Die Öffentlichkeit vertraut dem staatliche Rechtssystem nicht mehr. Finanzmärkte können aufgrund Fehlinformationen, deren Falschheit sich nur schwer und zeitlich deutlich verspätet nachweisen lässt, bisweilen so manipuliert werden, dass die Auswirkungen auf die Weltwirtschaft massiv sind etc.

Verschärft wird die Situation durch **kommerzielle Anreize**, denn für Onlineinhalte dominiert als Geschäftsmodell die Werbefinanzierung. Dies führt dazu, dass die Auswahl und Darstellung von Inhalten allein auf eine hohe Aufmerksamkeit der Nutzer (z.B. gemessen durch Klicken oder Teilen) gerichtet ist.<sup>15</sup> Da extreme und/oder emotional erschütternde Nachrichten eine besonders hohe Aufmerksamkeit er-

regen, besteht ein hoher Anreiz für Plattformen, diese besonders prominent anzuzeigen. Gefälschte Informationen können genau dieses Muster besonders einfach bedienen und verbreiten sich daher besonders leicht.

Die skizzierten unerwünschten Entwicklungen lassen sich unter folgenden Überschriften diskutieren:

- Desorientierung
- Gefährdung des sozialen Zusammenhalts und demokratischen Systems
- Machtverschiebungen im öffentlichen Diskurs
- Gefährdung von Beweisverfahren (polizeiliche Ermittlungen, Gerichtsentscheidungen)

Es wäre im VR6 außerdem zu klären, **welche positiven Entwicklungen** den nicht-intendierten und unbeabsichtigten Nebenfolgen möglicherweise gegenüber stehen, differenziert nach verschiedenen Nutzergruppen (Entwickler, Medien, Rezipienten). Eine explizite Beschreibung der intendierten Folgen (*benefits*) sollte den *unseens* gegenüber gestellt werden. *Benefits* könnten etwa sein: eine neue Form von Öffentlichkeit und Kommunikationsintensivierungen mit Folgen für sozialen Zusammenhalt, usw; auf technischer Seite die möglicherweise günstigere Film-, Bild und Audioproduktion; in künstlerischer Hinsicht die Entwicklung von neuen Ausdrucksformen, Stilmitteln oder gar ganzen Kunstgattungen in den darstellenden wie bildenden Künsten.<sup>16</sup>

Bei der Frage wie wir den hier skizzierten Vulnerabilitäten begegnen, muss das besondere Verhältnis zwischen Vertrauen, Verletzlichkeit und Technik be-

achtet werden. Wenn wir Vertrauen als akzeptierte Vulnerabilität verstehen (siehe *Begriffliche und epistemologische Einordnung*), stellt sich die Frage, ob ein auf den ersten Blick eventuell naheliegendes Ziel, Vulnerabilität komplett zu eliminieren, tatsächlich sinnvoll und/oder wünschenswert wäre.

Für viele Akteure in der digitalen Welt scheinen Vertrauensprobleme durch technische Anwendungen lösbar. Aktuell werden diverse Technologien entwickelt oder angeboten, die von dieser Annahme auszugehen scheinen:

Facebook will die Vertrauenswürdigkeit seiner Nutzer messen. Journalistische Recherchenetzwerke wollen algorithmenbasiert die Glaubwürdigkeit von Nachrichten prüfen. Im Rahmen eines von der EU geförderten Forschungsprojekts (*iBorderctrl*) wird unter anderem ein KI-basierter Lügendetektor entwickelt, der in drei Staaten (Griechenland, Lettland, Ungarn) getestet wird. Er soll die Vertrauenswürdigkeit von Personen, die in die EU einreisen wollen, messen und prüfen. Lange schon werden im Rahmen von Computersimulationen so genannte *trust games* durchgeführt, die die Erfolgsträchtigkeit unterschiedlicher Vertrauensstrategien messen. Entsprechendes zeichnet sich bereits auch für das hier beschriebene Phänomen ab.<sup>17</sup>

Allerdings muss genau geprüft werden, in welcher Hinsicht solche technischen Werkzeuge beim Aufbau von klugen Netzen des Vertrauens hilfreich sind.

### 3 Welche Stakeholder sind für ein Verständnis und ein Management der *Unseens* von besonderer Bedeutung? Welche wissenschaftlichen Wissensbereiche sind relevant?

Die Auswahl der Akteure kann sich am Raster der unter Punkt 4) skizzierten Perspektiven orientieren und entspricht der im Gesamtprojekt DiDaT angelegten und transdisziplinären Paarung von Wissenschaft und Praxis. Experten zur Bearbeitung der Fragestellungen könnten aus folgenden Bereichen kommen (wobei angesichts der Gruppengröße von 12 Personen Akteure mit übergreifendem Fachwissen bzw. mehrfachen Rollen bevorzugt sind):

- Schwerpunkt technische Perspektive: Experten für Künstliche Intelligenz, Deep Fakes, IT-Sicherheit, Zertifikatswesen, auch Datenwissenschaftler/ Bibliothekare
- Schwerpunkt gesellschaftliche Perspektive: Publizisten wie Journalisten, Verlagsinhaber, Blogger, Medienwissenschaftler; Juristen für Internetrecht, Datenschützer, Normungsspezialisten; auch Politiker
- Schwerpunkt philosophische Perspektive: Wissenschafts- und Technikphilosophen; Sozialtheoretiker
- Schwerpunkt ökonomische Perspektive: Wirtschaftswissenschaftler; Akteure im Online Marketing, Product Information Management, Content-Plattformen, Anbieter von IT-Hardware und Software, die Verifizierungsprobleme lösen

Darüber hinaus ist die Einbindung von Stakeholdern denkbar, die Auskunft zu teilweise analogen Problem- und/oder Lösungsfeldern geben können. Die Parallelen zwischen einem „verschmutzten Informationsökosystem“ und der Verschmutzung unserer natürlichen Umwelt liegen nahe. Interessante Impulse könnten aber auch aus den Wirtschaftswissenschaften oder aus Buchhaltungsregularien (Basel III etc.) kommen, wo Systeme von jeher auf Resilienz gegenüber nicht gutwilligen Akteuren ausgerichtet werden.

Auch eine historische Betrachtung insbesondere von politischer Propaganda und Gegenmaßnahmen in unterschiedlichen Ländern und Gesellschaftssystemen wäre hilfreich und sollte idealerweise als Kompetenz im Kreis der Akteure vertreten sein.

Es ergibt sich folgende Matrix von Vulnerabilitäten und Stakeholdern. Dabei wird unterschieden zwischen Stakeholdern, die von einer Vulnerabilität besonders betroffen sind (B), und Stakeholdern, die diese Vulnerabilität lösen können (L), wobei diese Einordnung nur eine grobe Orientierung sein kann:

| <b>Vulnerabilität</b> →<br><br><b>Stakeholder</b> ↓                  | Technische Möglichkeit zur überzeugenden Fälschung von digitaler Realität | Herkunft und Vertrauenswürdigkeit von Information nicht mehr klar →<br>Vertrauensverlust, Reality Apathy, Zynismus | Verlust der informationellen Grundlage für funktionierende ökonomische, soziale und politische Systeme | Gefährdung der Beweissführung in Polizei und Justiz, Vertrauensverlust in den Rechtsstaat |
|--|---|--|--|---|
| <b>Sicherheits- und Prüfinstitutionen<br/>Technikanbieter</b>        |   |  |  |   |
| Polizei, Verfassungsschutz,<br>Rechtspflege                          | B   | B  | B,L  | B,L   |
| Datenschutzaufsicht  | L   | L  |  | B,L   |
| Vertrauenswürdige neutrale<br>Instanzen (z.B. Prüfer, Zertifizierer) | L   | B  | L  | L   |
| Anbieter IT-Industrie Technik &<br>Dienstleistungen                  | L   | L  |  |   |

| <b>Vulnerabilität</b> →<br><b>Stakeholder</b> ↓              | Technische Möglichkeit zur überzeugenden Fälschung von digitaler Realität | Herkunft und Vertrauenswürdigkeit von Information → nicht mehr klar → Vertrauensverlust, Reality Apathy, Zynismus | Verlust der informationellen Grundlage für funktionierende ökonomische, soziale und politische Systeme | Gefährdung der Beweisführung in Polizei und Justiz, Vertrauensverlust in den Rechtsstaat |
|--|---|---|--|--|
| <b>Politik und Gesellschaft</b>                              |   |   |  |  |
| Politikentwickler und -entscheider, Verwaltungsentscheider   | B   | B,L   | B,L  | L  |
| Einzelne Wähler / Bürger                                     | B   | B   | B  | B  |
| Netzaktivisten, NGOs   | L   | L   | B,L  |  |
| <b>Medien</b>  |   |   |  |  |
| Journalisten, Blogger, Influencer etc.                       | B   | B,L   | B,L  |  |
| Medieninhaber, Chefredakteure                                | B   | L   | L  |  |
| Contentkuratoren und -aggregatoren<br>Suchmaschinenbetreiber | L   | B   | L  |  |
| Betreiber sozialer Netzwerke                                 | L   | B,L   | L  | L  |
| <b>Wissenschaft</b>  |   |   |  |  |
| IT-Spezialisten, Kryptografen, Lösungsarchitekten            | L   | L   | L  | L  |
| Medien- und Kommunikationswissenschaftler, Psychologen       | L   | L   | L  |  |
| Technikphilosophen, Systemanalytiker, Sozialtheoretiker      | L   | L   | L  | L  |

## Die Aporien der epistemisch-technischen Lösung

Der den technischen Lösungen des Vertrauensproblems zugrundeliegende Vertrauensbegriff geht von folgenden Annahmen aus. Vertrauen (wie Vertrauenswürdigkeit) wird demnach begriffen als:

- (1) **quantifizierbar**: Es handelt sich um eine durch einen Wahrscheinlichkeitsbegriff artikulierbare Größe.
- (2) **begründet**: Vertrauen beruht auf Gründen, die für die Vertrauenden als evident gelten. Diese Evidenz geht auf Informationen zurück.
- (3) **kognitive Erwartung**: Vertrauen geht daher mit einer kognitiven Erwartung einher. Kognitive Erwartungen unterscheiden sich von normativen u.a. dadurch, dass sie im Enttäuschungsfall aufgegeben werden, mit anderen Worten: Es wird gelernt.
- (4) **epistemische Beziehung**: Die andere Person kommt darin nicht als Person vor. Denn eigentlich vertraue ich nicht ihr, sondern meinem Erkenntnisvermögen.

Diese Annahmen sind zwar grundsätzlich nicht falsch, führen jedoch in der vorliegenden Form zu Aporien. Wenn Evidenz, wie Annahme (2) besagt, unabhängig von Vertrauen ist, kann man auf der Grundlage der richtigen Evidenz (Information) eine vernünftige Entscheidung treffen zu vertrauen oder zu misstrauen. Genau diese Unabhängigkeit ist in den meisten Fällen jedoch nicht gegeben. Wir betrachten nämlich umgekehrt eine uns präsentierte Evidenz (Information) im Licht von Vertrauen und Misstrauen. Es stellt sich daher ein zirkuläres Verhältnis zwischen Evidenz/Informationen und Vertrauen ein. Diesen Zirkel kann man nicht vermeiden, man kann nur mit ihm besser oder schlechter umgehen.

Die Annahme, das Vertrauensproblem sei technisch lösbar, ignoriert diesen Zirkel. Gehen wir dazu zu der Annahme zurück, die Vertrauenswürdigkeit einer Information lasse sich evidentiell durch eine Technik messen. Selbst wenn dies der Fall wäre, bestünde das Problem, dass man dazu dieser Technik bzw. ihren Entwicklern vertrauen müsste. Man kann also nicht aus der Abhängigkeit von anderen und damit aus dem Zirkel aus Evidenz und Vertrauen herauspringen. **Daher ist auch die Vorstellung, Fake News ließen sich einfach als solche technisch identifizieren und aus der Welt schaffen, fraglich.** Keine Person, die glaubt, dass eine Regierung Informationen unterdrückt und ihre Bürgerinnen und Bürger gezielt desinformiert, wird sich durch eine von diesem Staat (oder von staatlich abhängigen Institutionen) entwickelten Technologie darüber 'aufklären' lassen, dass dies nicht der Fall ist.

## Ein soziotechnisches Problem - eine soziotechnische Lösung

Das Problem muss also differenzierter begriffen und angegangen werden. Es trifft zwar zu, dass es sich *auch* um ein technisches Problem handelt; jedoch lässt es sich darauf nicht reduzieren. Daher ist auch die Lösung nicht rein technisch zu erzielen. Für die Erarbeitung von Lösungsansätzen muss davon ausgegangen werden, dass es sich um eine ebenso gesellschaftliche wie technische Aufgabe handelt. Techniken müssen demgemäß zwar eingesetzt oder gar neu entwickelt werden; diese können (in der Regel) das Vertrauensproblem jedoch nicht lösen, sondern uns lediglich dabei unterstützen, eine **Urteilkraft** auszubilden, um angemessen zwischen vertrauenswürdigen und nichtvertrauenswürdigen Personen und Aussagen zu unterscheiden. Dazu muss Vertrauen jedoch in einem Praxiszusammenhang begriffen werden, der sich daran zeigt, dass man in der Regel einer Person oder Institution geneigt ist zu vertrauen, weil man anderen Personen oder Institutionen vertraut, die ihr vertrauen. Es gibt mit anderen Worten **Netze des Vertrauens**.

**Die Frage, wie wir die Vertrauenswürdigkeit und Zuverlässigkeit digitaler Daten und Informationen sichern, ist also im Kern eine Frage danach, wie sich kluge Netze des Vertrauens knüpfen lassen und wie dies auf angemessene Weise technisch unterstützt werden kann.**

## 4 Methodische Überlegungen zur Unterstützung von Kernaussagen

Zur Erarbeitung möglicher Antworten durchdringen wir gleichzeitig vier Perspektiven und führen sie zusammen:

### 1. Technische Perspektive:

Was ist technisch machbar (jetzt oder in naher Zukunft)?

### 2. Gesellschaftliche, politische und rechtliche Perspektive:

Was findet Akzeptanz? Was ist national/international wünschenswert, regulierbar und durchsetzbar? An welche Institutionen kann dies geknüpft werden?

### 3. Philosophische Perspektive:

Welche Antworten sind hinsichtlich des von ihnen vorausgesetzten Verständnisses von Vertrauen/Vertrauenswürdigkeit, Wahrheit/Unwahrheit, Realität/Fiktion, Gewissheit, Bezeugung/Zeugenschaft, Wahrnehmung, Geltung, Legitimation und Beweis sowohl kohärent als auch anschlussfähig an die etablierte Sprachverwendung der medialen und politischen Öffentlichkeit?

### 4. Ökonomische Perspektive:

Was ist finanzierbar, disseminierbar oder langfristig lohnend?

Es gilt die Hypothese, dass die grundsätzlichen IT-Voraussetzungen zur Beantwortung der Leitfragen im Wesentlichen bereits vorhanden sind und nicht im Rahmen des Projekts entwickelt oder gefordert werden müssen (siehe dazu Abschnitt 1).

Die Kombination eines ungewöhnlich breiten Spektrums von Akteuren aus Gesellschaft, Medien, Wissenschaft und

Wirtschaft (vgl. vorherigen Abschnitt) soll von Anfang an einen lebendigen Austausch von Wissen und die Generierung möglicher Lösungselemente ermöglichen. Gleichzeitig wird es dadurch möglich, die Folgen der aktuellen technischen und gesellschaftlichen Entwicklungen breitgefächert und konsequent zu Ende zu denken, ggf. mithilfe von Szenariotechniken. Dadurch wird zusätzlicher Handlungsdruck aufgebaut und ein späterer Ergebnistransfer vorbereitet.

Die Wirksamkeit und Sinnhaftigkeit von Lösungsansätzen soll anhand einer Reihe frühzeitig definierter **Test Cases** geprüft werden. Jeder Test Case beschreibt eine - bereits beobachtete oder auch konstruierte - problematische Situation, für die der Effekt eines Lösungsansatzes durchgespielt werden kann. - Beispiele: „Der gewählte Präsident eines einflussreichen Landes bestreitet offensichtliche Fakten und ermutigt Gewalt gegen kritische Journalisten“ oder „Ein Massenmedium stellt reißerische ‚Nachrichten‘ ohne Rücksicht auf deren Wahrheitsgehalt in den Vordergrund, um Aufmerksamkeit und Werbeeinnahmen zu generieren“ oder „Ein bestechlicher Politiker bestreitet die Echtheit von Videodokumenten und lässt gleichzeitig ein falsches Video seines politischen Gegners herstellen und streuen, in dem diesem abstoßende Aussagen in den Mund gelegt werden“ oder „Ein repressives Regime unterdrückt eine unabhängige Presse mit der Behauptung, sie würde ‚fake news‘ verbreiten“.

Mit der Formulierung von Test Cases wird frühzeitig ein Rahmen für die ge-

meinsame Arbeit und Diskussion gesetzt und ein Konsens zum Zielkorridor hergestellt.

### **Bedarf für Vertiefungsforschung**

Es ist unklar, wie eine Infrastruktur für eine breit anerkannte, nicht staatlich beeinflusste Zertifizierung von Informationsquellen technisch aussehen könnte, insbesondere durch Reputationsträger und Autoritäten, die nicht bereits als „Marken“ aus dem Offline-Bereich bekannt sind. Die Vergabe von SSL-Zertifikaten für elektronische Signaturen kann ein Ausgangspunkt der Überlegungen sein, ist aber nur begrenzt übertragbar und hat zahlreiche Schwächen. Wichtig ist auch, dass eine anonyme (bzw. irreversibel pseudonyme) Kommunikation möglich bleibt. Zur Erarbeitung und prototypischen Demonstration technisch realisierbarer Vorschläge sollte Vertiefungsforschung im Umfang von **1 Personenjahr** eingeplant werden.

Dies kann auf insg. **1,5 Personenjahre** aufgestockt werden, um mehrere Informationsökosysteme parallel zu betrachten und Infrastrukturlösungen zu erarbeiten. Infrage kommen Informationsökosysteme wie die klassischen Nachrichtenmedien, Plattformen für nutzergenerierte Inhalte (Twitter, Facebook, YK, Wordpress, Medium etc.) sowie Fachpublikationen (auch in der Wissenschaft). Es können überdies weitere Ökosysteme herangezogen werden, um technologische Eigenschaften und Regulierungsmechanismen zu verstehen und ggf. durch Analogieschlüsse Handlungsmöglichkeiten für die Informationsökosysteme zu identifizieren. Beispiele für dieses Umfeld sind App-

Ökosysteme (Google, Apple), Datenökosysteme im Industrie-4.0-Bereich (z.B. Bosch IOTA), Zertifikatökosysteme (SSL klassisch bzw. mit CaCERT), Digitale Währungen oder Peer-to-Peer-Dateiaustausch-Ökosysteme (z.B. BitTorrent). Sinnvolle Unterstützung können auch Forschungen zur gesellschaftlichen Wirkung von Deep Fake-Videos sowie zur Evolution von Nachrichten im Internet liefern.

### **5. Erwartete Ergebnisse und Folgeinitiativen**

**Im Rahmen des Vulnerabilitätsraums „Vertrauenswürdigkeit von Informationen im digitalen Raum“ erarbeiten wir im inter- und transdisziplinären Dialog konsensfähige und praktikable Antworten auf die oben beschriebenen Herausforderungen.**

**Mögliche Fragestellungen:** Wie können wir Informationsökosysteme so gestalten, dass ein faktenbasierter gesellschaftlicher, wissenschaftlicher und politischer Diskurs möglich bleibt? Wie sorgen wir dafür, dass ein Dialog und eine ggf. auch mühsame Konsensfindung attraktiver bleiben als das Verharren auf extremen Positionen? Welche Anreize für die Wahrheitsfindung und -verbreitung können wir schaffen? Wie kann auch künftig mündige politische Meinungsbildung ablaufen?

**Erste Arbeitshypothesen:** Als Gerüst für erwartete Ergebnisse, zur Planung der Akteursauswahl sowie für die ersten Dialogschritte im Vulnerabilitätsraum dienen die folgenden Arbeitshypothesen zur Gestaltung des künftigen digitalen Raums:

1. Bisherige primär technologische Gegenmaßnahmen gegen Fake News wie die KI-gestützte Untersuchung von Videos oder die internen Netzwerkaktivitätsanalysen großer Internetplattformen sind letztlich nur ein Wettrüsten mit immer besseren Fälschungswerkzeugen und -methoden und daher bestenfalls eine partielle Lösung.
  2. Das Vertrauen in Informationen fußt fast immer auf dem Vertrauen in die Person/Institution, die sie verbreitet. Die Mechanismen zur Genese dieses Vertrauens (z.B. Andocken an den Augenschein in der realen Welt, Transfer, Crowdansätze etc.) sollten daher einen Schwerpunkt bilden. Die Frage ist, wie eine institutionelle Infrastruktur des Vertrauens aussehen könnte.
  3. Es ist nicht ausgeschlossen, dass es in einzelnen sozialen Handlungsbereichen zukünftig mehr Sinn macht, etablierte Beweisarten (z.B. Videobeweise) komplett zu ersetzen, d.h. alternative Strategien zu entwickeln, die helfen, Fakten von Fiktionen zu unterscheiden.
  4. Auf technischer und regulatorischer Seite erscheinen Maßnahmen wie bspw. Mechanismen und Standards für die Rückverfolgbarkeit von Informationen (u.U. mit Offenlegungspflichten für große Internetplattformen), eine Zertifizierung von Quellen und Kuratoren in Anlehnung an die etablierte Vergabe von SSL-Zertifikaten etc. überlegenswert. Blockchain und andere Technologien mit Notariatsfunktion können eine unterstützende Rolle spielen (z.B. für fälschungssichere Fingerabdrücke und Zeitstempel). Grundlegende Werkzeuge zur elektronischen Verschlüsselung und Signierung sind seit vielen Jahren verfügbar und mathematisch abgesichert.
  5. Die Auseinandersetzung „Anonymität vs. Pseudonymität vs. Klarnamen“ im Netz ist ein künstlich konstruierter Konflikt. Jeder der drei Ansätze ist in bestimmten Kontexten sinnvoll und muss für Menschen zugänglich sein. Die Verantwortlichkeit für eigene Inhalte ebenso wie für das Teilen von Fremdinhalten muss neu gedacht werden.
  6. Der Nachweis einer Lüge genügt nicht. Gesellschaftliche Konventionen und andere Faktoren bestimmen den Umgang mit ertappten Lügneren (vgl. Trump vs. Relotius). Vgl. auch das Phänomen „Reality Apathy“. Wir benötigen eine pragmatische Auseinandersetzung zur Existenz „objektiver“ Fakten oder einer objektiven Wahrheit<sup>18</sup> sowie der Frage, inwieweit Wahrheit tatsächlich gewollt ist, auch mit Blick auf psychologische Mechanismen.
  7. Gängige Geschäftsmodelle für Onlineinhalte – vor allem die Werbefinanzierung – stehen im Zielkonflikt mit Vertrauenswürdigkeit und müssen vermutlich weiterentwickelt bzw. ersetzt werden; gleichzeitig ist zu erwarten, dass nicht alle Lösungsvorschläge kommerziell tragfähig und stattdessen bspw. staatlich zu finanzieren sind. Letzteres wirft wiederum die Frage auf, inwieweit diese im Kontext repressiver Regime funktionieren würden.
- Diese Liste ist naturgemäß unvollständig und wird laufend verfeinert und ergänzt.**

## Endnoten:

[1] Melanie Ehrenkranz (2018): Researchers Come Out With Yet Another Unnerving, New Deepfake Method; Gizmodo; 09.11. 2018; unter: <https://gizmodo.com/researchers-come-out-with-yet-another-unnerving-new-de-1828977488> (abgerufen am 20.11.2019). Katyanna Quach (2018): The eyes don't have it! AI's 'deep-fake' vids surge ahead in realism; The Register; 11.09.2018; unter: [www.theregister.co.uk/2018/09/11/ai\\_fake\\_videos/](http://www.theregister.co.uk/2018/09/11/ai_fake_videos/) (abgerufen am 20.11.2019). Bloomberg LP (2016): It's Getting Harder to Spot a Deep Fake Video; Youtube; 27.09.2018; unter: [www.youtube.com/watch?v=gLoI9hAX9dw](http://www.youtube.com/watch?v=gLoI9hAX9dw); (abgerufen am 20.11.2019)

[2] Der US-Bundesstaat Kalifornien verbietet seit Oktober 2019 60 Tage vor einer Wahl die Verbreitung von manipulierten Videos, Tonspuren und Bildern eines Wahlkandidaten, die in böswilliger Absicht den Rufschädigung oder Wählerbeeinflussung betreiben, es sei denn, sie wurden als manipuliert gekennzeichnet. Siehe [https://leginfo.ca.gov/faces/billTextClient.xhtml?bill\\_id=20190200AB730](https://leginfo.ca.gov/faces/billTextClient.xhtml?bill_id=20190200AB730) (abgerufen am 20.11.2019). Facebook hat angekündigt, solche Videos in seinen Streams zu entfernen, die Ergebnis einer KI-Manipulation sind (so dass sie authentisch erscheinen) und zugleich in einer Weise bearbeitet wurden, die für einen Durchschnittsrezipienten nicht erkennbar ist und die wahrscheinlich zu der Überzeugung führt, dass eine dargestellte Person etwas sagte, was sie in Wirklichkeit nicht sagte. Siehe <https://about.fb.com/news/2020/01/enforcing-against-manipulated-media> (abgerufen am 06.01.2020)

[3] Von 14,698 Deepfake Videos, die im Sommer 2019 online gefunden wurden, waren 96% pornographischer Natur. Siehe für diese und folgende Aussagen: Henry Ajder, Giorgio Patrini, Francesco Cavalli, and Laurence Cullen (2019): The State of Deepfakes: Landscape, Threats, and Impact; Amsterdam: Deeptrace; September 2019

[4] Drew Harwell (2018): Fake-porn videos are being weaponized to harass and humiliate women: 'Everybody is a potential target'; The Washington Post; 30.12.2018; [www.washingtonpost.com/technology/2018/12/30/fake-porn-videos-are-being-weaponized-harass-humiliate-women-everybody-is-potential-target/](http://www.washingtonpost.com/technology/2018/12/30/fake-porn-videos-are-being-weaponized-harass-humiliate-women-everybody-is-potential-target/) (abgerufen am 20.11.2019)

[5] Siehe die bewusste Streuung von Videos, die die angeblich betrunkene Sprecherin des US-Repräsentantenhauses Nancy Pelosi und die angebliche körperliche Aggressivität des US-Reporters Jim Acosta belegen sollten. Dazu Henry Ajder, Giorgio Patrini, Francesco Cavalli, and Laurence Cullen (2019): The State of Deepfakes: Landscape, Threats, and Impact; Amsterdam:

Deeptrace; September 2019; S. 11 f.

[6] Katyanna Quach (2018): An AI system has just created the most realistic looking photos ever; The Register; 14.12. 2018; unter: [https://www.theregister.co.uk/2018/12/14/ai\\_created\\_photos/](https://www.theregister.co.uk/2018/12/14/ai_created_photos/) (abgerufen am 20.11.2019) Man beachte auch das eingebettete Video im Artikel.

[7] Siehe die Recherchen der britischen Consumers' Association im Verbrauchermagazin *Which?*: Hannah Walsh (2019): Thousands of 'fake' customer reviews found on popular tech categories on Amazon; *Which?*; 16.04.2019; [www.which.co.uk/news/2019/04/thousands-of-fake-customer-reviews-found-on-popular-tech-categories-on-amazon](http://www.which.co.uk/news/2019/04/thousands-of-fake-customer-reviews-found-on-popular-tech-categories-on-amazon) (abgerufen am 20.11.2019). Shefalee Loth (2018): The facts about fake reviews *Which?* investigators reveal tricks that sellers use to mislead online shoppers; *Which?*; 25.10.2018; [www.which.co.uk/news/2018/10/the-facts-about-fake-reviews](http://www.which.co.uk/news/2018/10/the-facts-about-fake-reviews) (abgerufen am 20.11.2019)

[8] Lion Gu, Vladimir Kropotov, and Fyodor Yarochkin: The Fake News Machine. How Propagandists Abuse the Internet and Manipulate the Public (2017); TrendLabs Research Paper; Hrsg: Trend Micro, Incorporated; unter: [https://documents.trendmicro.com/assets/white\\_papers/wp-fake-news-machine-how-propagandists-abuse-the-internet.pdf](https://documents.trendmicro.com/assets/white_papers/wp-fake-news-machine-how-propagandists-abuse-the-internet.pdf) (abgerufen am 20.11.2019).

[9] Adrian Chen (2015): The Agency. From a nondescript office building in St. Petersburg, Russia, an army of well-paid "trolls" has tried to wreak havoc all around the Internet — and in real-life American communities; New York Times; 02.06.2015; unter [www.nytimes.com/2015/06/07/magazine/the-agency.html](http://www.nytimes.com/2015/06/07/magazine/the-agency.html) (abgerufen am 20.11.2019). The Economist (2018): Inside the Internet Research Agency's lie machine (Briefing); The Economist; 22.02.2018; [www.economist.com/briefing/2018/02/22/inside-the-internet-research-agencys-lie-machine](http://www.economist.com/briefing/2018/02/22/inside-the-internet-research-agencys-lie-machine) (abgerufen am 20.11.2019)

[10] British Broadcasting Corporation (2018): Cambridge Analytica planted fake news; BBC; 20.03.2018; unter [www.bbc.com/news/av/world-43472347/cambridge-analytica-planted-fake-news](http://www.bbc.com/news/av/world-43472347/cambridge-analytica-planted-fake-news) (abgerufen am 20.11.2019)

[11] Henry Ajder, Giorgio Patrini, Francesco Cavalli, and Laurence Cullen (2019): The State of Deepfakes: Landscape, Threats, and Impact; Amsterdam: Deeptrace; September 2019; S. 9

[12] Elizabeth F Loftus (2005): Searching for the neurobiology of the misinformation effect: Planting misinformation in the human mind: A 30-year investigation of the malleability of memory; in: *Learning & Memory*; July 1, 2005 Vol: 12; S. 361-366; DOI: 10.1101/lm.94705; unter: <https://www.researchgate.net/publication/80457>

---

38\_Searching\_for\_the\_neurobiology\_of\_the\_misinformation\_effect (abgerufen am 20.11.2019)  
[13] Charlie Warzel (2018): He predicted the 2016 fake news crisis. Now he's worried about an information apocalypse; Buzz Feed News; 11.02.2018; unter:  
[www.buzzfeed.com/charliewarzel/the-terrifying-future-of-fake-news](http://www.buzzfeed.com/charliewarzel/the-terrifying-future-of-fake-news) (abgerufen am 20.11.2019)  
[14] Miriam Meckel (2018): Eine neue Schicht Endfrustrierter treibt uns in die Infocalypse. Was dagegen helfen könnte; Wirtschaftswoche; 23.02.2018; unter:  
[www.wiwo.de/politik/deutschland/schlusswort-laesst-sich-die-infocalypse-noch-abwenden/20989742.html](http://www.wiwo.de/politik/deutschland/schlusswort-laesst-sich-die-infocalypse-noch-abwenden/20989742.html) (abgerufen am 20.11.2019). Oscar Schwartz (2018): You thought fake news was bad? Deep fakes are where truth goes to die; The Guardian; 12.11.2018;  
[www.theguardian.com/technology/2018/nov/12/deep-fakes-fake-news-truth](http://www.theguardian.com/technology/2018/nov/12/deep-fakes-fake-news-truth) (abgerufen am 20.11.2019). Elmar Theveßen (2019): Manipulierte Videos - Wie Deep Fakes die Welt gefährden ; in zdf - heute Nachrichten; 29.05.2019;  
[www.zdf.de/nachrichten/heute/infokalypse-wie-deep-fakes-die-welt-gefaehrden-100.html](http://www.zdf.de/nachrichten/heute/infokalypse-wie-deep-fakes-die-welt-gefaehrden-100.html) (abgerufen am 20.11.2019)  
[15]vgl. z.B. diese Tipps einer Werbeagentur:  
<https://blog.hootsuite.com/de/facebook-algorithmus-organische-reichweite/> (abgerufen am 20.11.2019)  
[16] Siehe zumindest für den Einsatz von KI bei Multimedia-Produktionen in den Bildenden Künste die Werke von Gillian Wearing, etwa in ihrer Ausstellung im Cincinnati Art Museum: [www.cincinnatiartmuseum.org/wearing](http://www.cincinnatiartmuseum.org/wearing) . Ein Videoausschnitt demonstriert den Einsatz besonders anschaulich:  
<https://player.vimeo.com/video/295991070> (abgerufen 20.11.2019)  
[17] Für die Enttarnung von DeepFakes siehe z.B.: TU München (2019): Software FaceForensics erkennt Fake-Videos am zuverlässigsten. Künstliche Intelligenz enttarnt Fake-Videos; Pressemitteilung der Technischen Universität München; 09.06.2019;  
[www.tum.de/nc/die-tum/aktuelles/pressemitteilungen/details/35501](http://www.tum.de/nc/die-tum/aktuelles/pressemitteilungen/details/35501) (abgerufen am 20.11.2019);  
[18] unter Berücksichtigung der bereits vorhandenen philosophischen Erkenntnisse und Traditionen