

# Artificial General Intelligence (AGI)

## Zwischen Versprechungen und Realität

### IN KÜRZE

#### Artificial General Intelligence (AGI)

- bezeichnet bislang nicht existierende KI-Systeme, die flexibel, autonom und kontextübergreifend Informationen verarbeiten, lernen und handeln – ähnlich wie Menschen.
- wird aktuell sehr unterschiedlich definiert: von funktionaler Intelligenz bis hin zu bewusstem Erleben.
- wird von verschiedenen Akteuren unterschiedlich aufgeladen – etwa als Fortschrittsversprechen, Risiko oder Science-Fiction-Szenario.
- ist Gegenstand philosophischer Debatten über Bewusstsein, Körperlichkeit und Kreativität: Was zeichnet „echte“ Intelligenz aus?
- steht exemplarisch für das Spannungsfeld zwischen technologischer Machbarkeit, grundlegenden philosophischen Fragen und gesellschaftlicher Imagination.

Der Begriff „AGI“ wird mit sehr unterschiedlichen Bedeutungen und Zielen verwendet – von technischer Spezifikation bis Zukunftsvision. Diese Vieldeutigkeit führt zu überhöhten Erwartungen und verzerrten öffentlichen Debatten.

## Kontext

Bereits in den 1980er Jahren wurde die Unterscheidung zwischen so genannter schwacher und starker Künstlicher Intelligenz (KI) geprägt. Während schwache KI spezifische Aufgaben ausführt, sollte starke KI perspektivisch menschliche Intelligenz in vollem Umfang reproduzieren oder simulieren, einschließlich Bewusstsein und Intentionalität. Die Idee, menschliche Kognition formal zu fassen, reicht allerdings noch weiter zurück: Bereits 1943 schlugen Warren McCulloch und Walter Pitts ein vereinfachtes Modell neuronaler Aktivität vor, das als theoretische Grundlage für künstliche Neuronennetze diente. Früh wurde jedoch deutlich, dass solche Modelle nur einen sehr begrenzten Ausschnitt menschlicher Intelligenz abbilden können. Heute hat sich für das Zielbild starker KI zunehmend der Begriff „Artificial General Intelligence“ (AGI) durchgesetzt. Trotz enormer technologischer Fortschritte in den vergangenen Jahren bleiben heutige KI-Systeme weit entfernt von einer AGI. Diese ist mit großen Hoffnungen auf medizinische, wissenschaftliche und gesellschaftliche Durchbrüche verbunden; zugleich aber auch mit tiefgreifenden Fragen nach Kontrolle, Verantwortung und den Grenzen zwischen Mensch und Maschine. AGI ist damit nicht nur ein technisches Ziel, sondern auch ein Spiegel gesellschaftlicher Vorstellungen, Erwartungen und Befürchtungen.

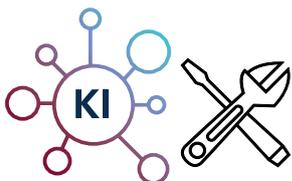
## Was ist Artificial General Intelligence?

Aktuell existiert keine allgemein anerkannte Definition von AGI und auch kein Test, der das Erreichen bestätigen würde. Verschiedene Forschungsdisziplinen, Unternehmen und politische Akteure setzen unterschiedliche Schwerpunkte: Einige Definitionen betonen kognitive Flexibilität und die Fähigkeit zu autonomem Lernen, andere fokussieren Bewusstsein, selbstständige Zielsetzung oder Transferleistung. Ein zentrales Problem bei allen Ansätzen ist, dass sie sich zumeist an einer Vorstellung von menschlicher Intelligenz orientieren – doch auch diese ist nicht klar definiert. In Psychologie, Neurowissenschaft oder Philosophie existieren konkurrierende Ansätze, was Intelligenz ausmacht. Diese doppelte Unschärfe – beim Menschen wie bei der Maschine – verdeutlicht: AGI ist nicht nur ein technisches Konzept, sondern auch ein Projektionsraum für wissenschaftliche, wirtschaftliche und gesellschaftliche Interessen.

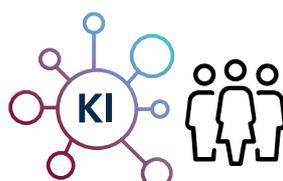
Akteure	Merkmale von AGI	Interesse	Kritik
<b>Forschung</b>	<ul style="list-style-type: none"> <li>■ Kognitive Flexibilität</li> <li>■ Autonomes Lernen</li> <li>■ Transferleistung</li> <li>■ (Teilweise) Bewusstsein</li> </ul>	<ul style="list-style-type: none"> <li>■ Theoretisches Verständnis von Intelligenz</li> <li>■ Disziplinäre Paradigmen (z. B. Psychologie vs. Informatik)</li> </ul>	<ul style="list-style-type: none"> <li>■ Uneinheitliche Definition von (menschlicher) Intelligenz</li> </ul>
<b>Unternehmen</b>	<ul style="list-style-type: none"> <li>■ Leistungsfähigkeit</li> <li>■ Adaptivität</li> <li>■ Breite Anwendbarkeit</li> </ul>	<ul style="list-style-type: none"> <li>■ Strategische Kommunikation</li> <li>■ Investitionsförderung</li> <li>■ Marktvorteile durch Innovationsnarrative</li> </ul>	<ul style="list-style-type: none"> <li>■ Definitionen meist marketinggetrieben</li> <li>■ Gefahr der Überschätzung/ Erwartungshaltung</li> </ul>
<b>Science-Fiction</b>	<ul style="list-style-type: none"> <li>■ Übermenschliche Intelligenz</li> <li>■ Emotionale Tiefe</li> <li>■ Eigenständige Ziele</li> </ul>	<ul style="list-style-type: none"> <li>■ Dramaturgie/ kulturelle Reflexion</li> <li>■ Zukunftsvisionen/ Warnungen</li> </ul>	<ul style="list-style-type: none"> <li>■ Starke Prägung der gesellschaftlichen Vorstellung</li> <li>■ Verzerrung der realen Forschung</li> </ul>
<b>Gesellschaft</b>	<ul style="list-style-type: none"> <li>■ Menschenähnliches Denken und Fühlen</li> <li>■ „Alleskönnerin“</li> <li>■ Bewusstsein</li> <li>■ Autonomie</li> </ul>	<ul style="list-style-type: none"> <li>■ Orientierung durch Medien &amp; Fiktion</li> <li>■ Hoffnungen auf Problemlösung</li> <li>■ Ängste vor Kontrollverlust</li> </ul>	<ul style="list-style-type: none"> <li>■ Fiktion statt Fakten als Bezugsrahmen</li> <li>■ Polariserte Debatte zwischen Hype und Dystopie</li> </ul>

## Technische Grundlagen: Wie kann AGI funktionieren?

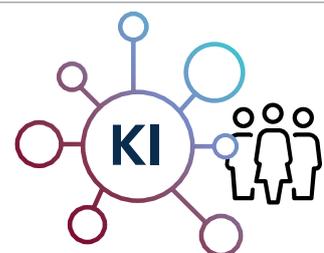
Die Entwicklung einer Artificial General Intelligence ist eine der großen Herausforderungen der modernen KI-Forschung. Hierbei wird die Frage diskutiert, ob eine solche Entwicklung überhaupt möglich ist und falls ja, welche Ansätze dabei vielversprechend sind.



**Artificial Narrow Intelligence (ANI)**  
Kann spezifische Aufgaben ausführen



**Artificial General Intelligence (AGI)**  
Kann sich bei allen Aufgaben menschenähnlich verhalten



**Artificial Super Intelligence (ASI)**  
Intelligenter als der Mensch – wie in Science-Fiction

Quelle: Plattform Lernende Systeme

## Ansätze zur Entwicklung von AGI

Die aktuell insbesondere von großen Tech-Unternehmen vertretene **Skalierungsthese** besagt, dass durch Wachstum von Modellgröße, Datenmenge und Rechenkapazität irgendwann allgemeine Intelligenz erreicht werden könnte. BefürworterInnen dieser These argumentieren, dass bereits heutige KI-Modelle mit zunehmender Skalierung immer leistungsfähiger werden und komplexere Aufgaben bewältigen – so etwa Sprachmodelle wie GPT-5, die mit steigender Größe ein besseres Sprachverständnis zeigen.

KritikerInnen argumentieren hingegen, dass Skalierung allein nicht ausreichen wird, um Fähigkeiten wie Transferlernen, kausales Denken, Begründen, Handlungskompetenz, Zwecksetzung oder Autonomie zu realisieren. Bestehende KI-Modelle haben nach wie vor Schwierigkeiten, Wissen aus einem Bereich auf einen anderen zu übertragen oder eigenständige Strategien zur Problemlösung zu entwickeln. Auch empirisch ist die Skalierungsthese angreifbar: Die zugrunde liegenden Korrelationen („Skalierungsgesetze“) sind keine Naturgesetze und scheinen nur innerhalb bestimmter Größenordnungen zu gelten. Entsprechend lässt sich die Performanz rein skalenbasierter Modelle nicht unbegrenzt steigern. Vielmehr deutet sich an, dass zusätzliche algorithmische Innovationen notwendig sind, um signifikante Leistungssteigerungen zu erzielen. Darüber hinaus wird darauf hingewiesen, dass die aktuellen Modelle bereits den Großteil der hochqualitativen Trainingsdaten nutzen und dementsprechend eine weitere Skalierung bereits an Datenmangel scheitern könnte.

Ansatz	Kernidee	Versprechen	Kritik
<b>Skalierung</b>	<ul style="list-style-type: none"> <li>Immer größere Modelle, mehr Daten und Rechenleistung</li> <li>→ AGI emergiert durch schierere Größe</li> </ul>	<ul style="list-style-type: none"> <li>Schnelle Fortschritte bei bestimmten Aufgaben</li> <li>Bereits heute leistungsfähige Modelle</li> <li>Klarer Pfad zur Umsetzung</li> </ul>	<ul style="list-style-type: none"> <li>Mangel an Transfer, Kausalität, Vernunft</li> <li>Ressourcenintensiv</li> <li>„Emergenz“ bleibt spekulativ</li> <li>Perspektivisch: Datenmangel</li> </ul>
<b>Hybride KI</b>	<ul style="list-style-type: none"> <li>Kombination von symbolischer KI und neuronalen Netzen</li> </ul>	<ul style="list-style-type: none"> <li>Bessere Generalisierung &amp; Erklärbarkeit</li> <li>Integration von Weltwissen</li> </ul>	<ul style="list-style-type: none"> <li>Umfassende Integration bisher ungelöst</li> <li>Oft weniger skalierbar als reine Deep-Learning-Modelle</li> <li>Noch kaum standardisiert</li> </ul>
<b>Neue Paradigmen/ evolutionäre Ansätze</b>	<ul style="list-style-type: none"> <li>AGI durch neue Lernmechanismen, Selbstorganisation &amp; Körperlichkeit</li> </ul>	<ul style="list-style-type: none"> <li>Orientierung an biologischen Systemen</li> <li>Fokus auf Motivation, Zielsetzung, Autonomie</li> <li>Langfristig robuster</li> </ul>	<ul style="list-style-type: none"> <li>Noch im Frühstadium</li> <li>Schwer messbar &amp; aufwendig</li> </ul>

Hier setzen **hybride KI-Modelle** an, die datengetriebene und wissensbasierte KI kombinieren. Sie versuchen, die Stärke neuronaler Netze im Umgang mit großen Datenmengen mit der logischen Struktur symbolischer KI zu verbinden. Ziel ist es, Systeme zu schaffen, die nicht nur Muster erkennen, sondern auch Erklärungen liefern und mit expliziten Regeln umgehen können. Hybride Ansätze gelten daher als vielversprechend für Anwendungsfelder, in denen erklärbares und zuverlässiges Verhalten erforderlich ist – etwa in der Medizin oder im Recht. Ein erfolgreiches Beispiel ist AlphaFold, das die Faltung von Proteinen mit zuvor unerreichter Genauigkeit vorhersagt – die Entwickler wurden 2024 mit dem Nobelpreis ausgezeichnet.

Darüber hinaus gibt es **alternative Paradigmen**, die ganz andere Wege zur AGI beschreiten wollen. Dazu gehören **evolutionäre Ansätze**, bei denen KI-Systeme über viele Generationen hinweg mittels Variation und Selektion weiterentwickelt werden – analog zur biologischen Evolution. Auch **verkörperte (embodied) KI** wird diskutiert: Hier geht man davon aus, dass echte Intelligenz nur durch physische Interaktion mit der Umwelt entstehen kann – also durch Sensorik, Motorik und verkörperte Erfahrung. Diese Konzepte knüpfen an entwicklungspsychologische und neurokognitive Theorien an, die den engen Zusammenhang zwischen Wahrnehmung, Körper und Denken betonen.

## Einschätzungen von führenden KI-ExpertInnen

Die öffentliche Debatte über AGI wird derzeit stark von Stimmen aus dem angloamerikanischen Raum geprägt. Dies spiegelt sowohl die technologische Führungsrolle dieser Länder als auch die dort vorherrschenden Narrative über die Zukunft von KI wider.

### AGI ist erreichbar – bald

**Sam Altman, CEO von OpenAI;** Altman sieht AGI als realistische Zukunft in greifbarer Nähe. Er setzt auf die Skalierung bestehender Systeme (Transformermodelle) und investiert mit OpenAI massiv in Hardware.

### Stufenweise Annäherung an AGI

**Demis Hassabis, CEO von DeepMind, Nobelpreisträger;** Hassabis verfolgt einen wissenschaftlich geprägten, vorsichtigen Ansatz. Er sieht AGI als langfristiges Ziel, erreichbar durch Fortschritte etwa in Reinforcement Learning, Gedächtnisarchitekturen und Simulation.

### Sorge vor Kontrollverlust

**Geoffrey Hinton, KI-Pionier, ehemals Google, Nobelpreisträger;** Hinton hält AGI technisch für möglich, äußert aber Besorgnis über Kontrollverluste bei mächtigen KI-Systemen und warnt vor übertriebener Euphorie.

### Generelle Skepsis gegenüber AGI

**Yann LeCun, Chef-KI-Forscher bei Meta, Turing-Award-Träger;** LeCun betont, dass aktuelle Systeme fundamentale Fähigkeiten wie kausales Denken oder gesunden Menschenverstand nicht beherrschen, und fordert neue Ansätze jenseits von Deep Learning.

### Forderung nach hybriden Ansätzen

**Gary Marcus, KI-Kritiker, Neurowissenschaftler;** Marcus hält Deep Learning allein für unzureichend und fordert eine Verbindung von symbolischer Logik und statistischen Methoden, um eine nach seiner Vorstellung echte AGI zu erreichen.

### Warnung vor KI mit Bewusstsein

**Thomas Metzinger, Philosophieprofessor, ehemaliges Mitglied einer EU-Experten-gruppe zu KI-Ethik;** Metzinger stellt im Sinne eines Gedankenexperiments die Überlegung an, dass durch die Entwicklung von KI-Systemen mit Bewusstsein eine massive Zunahme von potenziell leidenden künstlichen Subjekten zu erwarten sei. Er spricht sich für ein Entwicklungs-Moratorium aus, solange ethische Aspekte nicht hinreichend reflektiert wurden.

### Kritik an aktuellen AGI-Narrativen

**Margaret Mitchell, KI-Ethikerin, ehemals bei Google AI Ethics;** Mitchell zufolge ist die AGI-Debatte oft wissenschaftlich unscharf und verschleiert narrative Machtverhältnisse. Sie fordert mehr Fokus auf Fairness, Verantwortung und heutige reale Auswirkungen von KI.

### Kritik an Machtstrukturen und fehlender Regulierung

**Meredith Whittaker, Präsidentin der Signal Foundation;** Whittaker sieht AGI nicht primär als technisches, sondern als politisches Thema. Sie warnt vor der Konzentration von KI-Systemen in den Händen weniger Tech-Konzerne und fordert demokratische Kontrolle, Transparenz und Schutz öffentlicher Infrastruktur.

## Philosophische und gesellschaftliche Einordnung

Die Diskussion um Artificial General Intelligence wirft nicht nur technische, sondern auch philosophische und gesellschaftliche Fragen auf. Ein zentrales Thema ist dabei die **Definition von Intelligenz** selbst: Reicht die Fähigkeit zur Problemlösung und zum Lernen aus – oder gehören auch Kreativität, Bewusstsein, Körperlichkeit und Autonomie untrennbar dazu?

Die Debatte um AGI führt damit unweigerlich an die Grenzen unseres Verständnisses von Intelligenz, Subjektivität und Kreativität – und spiegelt letztlich auch Fragen wider, die wir uns über uns selbst stellen. Ob Maschinen je vergleichbar mit dem Menschen bewusst, kreativ oder autonom sein können (oder dies nur simulieren), bleibt vorerst offen. Doch schon die Auseinandersetzung mit diesen Möglichkeiten zeigt, wie sehr das AGI-Konzept nicht nur ein technisches Ziel, sondern auch ein kultureller Spiegel ist. Was wir unter „allgemeiner Intelligenz“ verstehen – und ob wir sie in Maschinen wiedererkennen wollen –, hängt eng mit unseren Vorstellungen vom Menschsein zusammen.

	Mensch	AGI
<b>Bewusstsein</b> 	<ul style="list-style-type: none"> <li>Intelligenz ist untrennbar mit Bewusstsein verbunden (subjektives Erleben, Empfindung, Selbstgefühl)</li> </ul>	<ul style="list-style-type: none"> <li>Funktionen wie Gedächtnis und Sprachverarbeitung sind technisch erklärbar, jedoch nicht das subjektive Erleben („hard problem of consciousness“)</li> <li>oder</li> <li>Funktionalistische KI-Forschung: AGI mit limitierten kognitiven Funktionen ist möglich ohne echtes Bewusstsein</li> </ul>
<b>Körperlichkeit &amp; Autonomie</b> 	<ul style="list-style-type: none"> <li>Menschen erleben und verstehen die Welt nicht nur durch abstrakte Symbole, sondern durch ihren Körper, ihre Sinne, Emotionen und sozialen Interaktionen</li> <li>Autonomie basiert auf biologischen Antrieben: Überleben, Schmerzvermeidung, Lust, Fortpflanzung</li> </ul>	<ul style="list-style-type: none"> <li>Keine sinnliche Erfahrung, keine biologischen Antriebe – damit auch keine natürliche Motivation</li> <li>oder</li> <li>Digitale Systeme könnten ebenfalls Formen von Zielsetzungen entwickeln, etwa durch Belohnungsmechanismen</li> </ul>
<b>Kreativität</b> 	<ul style="list-style-type: none"> <li>Schöpferische Intention: Neues entsteht mit Bedeutung, Kontextverständnis und kultureller Einbettung</li> <li>Kombinatorische, explorative und transformative Kreativität: Menschen können aus eigenem Antrieb völlig neue Dinge erschaffen</li> </ul>	<ul style="list-style-type: none"> <li>Kreativität auf Grundlage von statistischen Kombinationen, nicht aus innerer Intention oder kultureller Einbettung</li> <li>Fraglich, ob KI je zu transformativer Kreativität (grundlegendes Neudenken kultureller Konventionen) in der Lage sein wird</li> </ul>

## Rechtliche und ethische Herausforderungen

Die Entwicklung einer Artificial General Intelligence reicht weit über eine technische Vision hinaus. Sie berührt fundamentale Fragen des Rechts, der Ethik und des gesellschaftlichen Zusammenlebens.



### Haftung: Wer trägt die Verantwortung für AGI-Entscheidungen?

Wer ist verantwortlich, wenn – in einer hypothetischen Zukunft – die Entscheidung einer AGI zu einem Schaden führt: Die EntwicklerInnen? Das betreibende Unternehmen? Die Nutzenden? Die AGI selbst? Je autonomer ein System handelt, desto schwieriger wird es, klare Verantwortlichkeiten zuzuweisen. Zwar adressieren Regelungen wie die DSGVO bereits automatisierte Entscheidungen und stellen strenge Anforderungen an Transparenz und Verantwortlichkeit. Doch setzen sie voraus, dass es einen nachvollziehbaren menschlichen Verantwortungszusammenhang gibt. Bei einer echten AGI, die per Definition eigenständig entscheidet, würden diese Voraussetzungen entfallen – und damit auch der Anwendungsrahmen bestehender Rechtsprinzipien.



### Rechte für Maschinen? Die Debatte um den rechtlichen und ethischen Status von AGI

Mit zunehmender Autonomie stellt sich die Frage, ob eine AGI mit Bewusstsein, eigenen Zielsetzungen oder einer Art Erleben weiterhin als bloßes Werkzeug gelten kann – oder ob sie, analog zu juristischen Personen wie Unternehmen, einen eigenen Rechtsstatus erhalten sollte. Einige EthikerInnen und JuristInnen argumentieren, dass sich aus einem moralischen Status auch Rechte ableiten könnten – etwa das Recht, nicht abgeschaltet zu werden oder über sich selbst zu verfügen.

Wenn eine AGI als Trägerin moralischer Ansprüche anerkannt würde, könnten daraus nicht nur Schutzrechte, sondern auch Fürsorgepflichten resultieren – ähnlich wie bei natürlichen Personen. Daraus ergeben sich Anschlussfragen: Dürfte man eine AGI zur Arbeit zwingen? Könnten Klonverbote oder Schutzmechanismen analog zur Bioethik Anwendung finden? Ein anerkannter moralischer oder rechtlicher Status könnte auch Pflichten begründen: So wie Unternehmen für bestimmte Schäden einstehen müssen, ließe sich diskutieren, ob auch AGI-Systemen normative Verantwortung zugeschrieben werden kann – etwa durch digitale Treuhandsysteme oder Stellvertreterhaftung.



### Regulierung: Reicht bestehendes KI-Recht aus?

Derzeitige Gesetze – etwa der AI Act der Europäischen Union – orientieren sich primär an Risiken heutiger KI-Systeme. Mit spezifischen Anforderungen für „General Purpose AI“ (GPAI) enthält der AI Act zwar erste Regelungen für breit einsetzbare KI-Modelle. Sie dürften jedoch nicht ausreichen, um hypothetische AGI-Systeme zu erfassen, die eigene Ziele verfolgen oder in sozialen Kontexten agieren. Hier stellt sich die Frage, ob es Anpassungen am AI Act oder sogar ein eigenes Rechtsregime für AGI braucht – vergleichbar mit dem Völkerrecht oder bioethischen Sonderregelungen. Unklar bleibt auch, wer dieses Regime gestalten sollte.



### Wohin entwickelt sich KI?

Diejenigen, die eine AGI in naher Zukunft für erreichbar halten, sind häufig auch der Ansicht, dass deren Weiterentwicklung zu einer Superintelligenz führen wird, die die menschliche Intelligenz kognitiv überflügelt. Häufig sind damit Befürchtungen verbunden, eine Superintelligenz könnte anstreben, die Menschheit zu beherrschen. Kurzfristig dürfte jedoch nicht die Maschine selbst zur Gefahr werden – sondern die Frage, wer sie kontrolliert: Könnten einzelne Akteure – etwa große Technologiekonzerne oder Staaten – Superintelligenz nutzen, um ihre Macht auszuweiten?

## Bewertung: Stimmen aus der Plattform Lernende Systeme



*Für KI-Anwendungen ist alles gleich: Sie funktionieren für ein intelligentes Ampelsystem ebenso wie für die Analyse von Liebesgedichten. KI versteht nichts, sondern errechnet Muster. Ihr Blick in die Zukunft ist eine Wahrscheinlichkeitsprognose. Sie kann nichts hoffen, denn sie hat auch nichts zu verlieren. KI fehlt die Erfahrung einer körperlichen Existenz. Auch wenn KI-Sprachmodelle alles kopieren können, was Menschen zu Bewusstsein, Wünschen oder Solidarität sagen, verbinden sie damit so viel wie mit der schon angesprochenen Ampelschaltung. Superintelligenz ist ein hypothetisches Konzept, ein interessantes Gedankenexperiment, aber auch nicht mehr.*

**Jessica Heesen**, Professorin am Internationalen Zentrum für Ethik in den Wissenschaften (IZEW) an der Universität Tübingen

*AGI bleibt ein vages Konzept. Trotz Erfolgen von Deep Learning fehlt KI-Systemen gesunder Menschenverstand: Sie denken nicht logisch, scheitern an Neuem und verschlingen Ressourcen. Reine Skalierung führt uns nicht zu AGI. Stattdessen müssen wir KI mit Kognitionswissenschaft verbinden und Deep Learning mit symbolischem Schlussfolgern vereinen. So schaffen wir vernünftige KI, die komplexe Probleme löst, Vertrauen schafft und etwa in der Verwaltung Bürgeranliegen effizient bearbeitet.*

**Kristian Kersting**, Professor für Maschinelles Lernen und Künstliche Intelligenz an der TU Darmstadt und Co-Direktor von hessian.AI



*KI ist mehr als ein gewöhnliches Werkzeug und weniger als ein menschliches Subjekt. Die Wahrscheinlichkeit einer Superintelligenz, die die Menschheit vernichten will, ist deshalb gering. Die Gefahren liegen in den enormen Machtspielräumen, die KI einigen wenigen eröffnet, sowie im Mangel an Vorhersehbarkeit und Kontrolle, der ein strukturelles Merkmal von KI ist. Außerdem übernimmt KI zunehmend Tätigkeiten, die Menschen eigentlich gut und gerne erledigen, etwa im zwischenmenschlichen oder kreativen Bereich. Ich würde mir demgegenüber mehr Entwicklung in Problemfeldern wünschen, in denen wir noch nicht über gute Lösungen verfügen.*

**Catrin Misselhorn**, Professorin für Philosophie an der Georg-August-Universität Göttingen

## Welche Fragen sind offen?

- **Sicherheit:** Wie lässt sich sicherstellen, dass AGI-Systeme im Einklang mit menschlichen Werten handeln – auch bei offenem Lernverhalten und wachsender Autonomie?
- **Kontrolle:** Wie kann verhindert werden, dass sich eine AGI verselbstständigt oder von einer kleinen Gruppe missbraucht wird?
- **Philosophie:** Welche Eigenschaften (z. B. Bewusstsein, Verstehen, Körperlichkeit) sind nötig, damit man von „echter“ Intelligenz sprechen kann?
- **Governance:** Reichen nationale Regulierungen aus – oder braucht es internationale Regeln und Institutionen für die AGI-Entwicklung?
- **Einflussfaktoren:** Wie stark prägen wirtschaftliche Interessen, mediale Narrative und Science-Fiction unsere Erwartungen an AGI – und wie wirkt sich das auf Forschung und Politik aus?

## Quellen

Altman, S. (2023). Sam Altman on AGI: Scaling large language models is not enough. The Decoder. <https://the-decoder.com/sam-altman-on-agi-scaling-large-language-models-is-not-enough/>

Bender, E. M., Gebru, T. et al. (2021). On the dangers of stochastic parrots: Can language models be too big? In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT). <https://arxiv.org/pdf/2002.06177>

Hassabis, D. (2025). Demis Hassabis: "Why AI must be built responsibly". TIME 100 Interview. <https://time.com/7277608/demis-hassabis-interview-time100-2025/>

Krasheninnikov, A., Yao, J., & Sutskever, I. (2023). On the measurement of AGI progress. arXiv. <https://arxiv.org/pdf/2311.02462>

Mitchell, M., Wu, S., Zaldivar, A. et al. (2019). Model cards for model reporting. In Proceedings of the Conference on Fairness, Accountability, and Transparency (FAT\*). <https://arxiv.org/pdf/1810.03993>

University of Toronto (2023). Risks of artificial intelligence must be considered as technology evolves: Geoffrey Hinton. <https://web.cs.toronto.edu/news-events/news/risks-artificial-intelligence-must-be-considered-technology-evolves-geoffrey-hinton>

---

Inhalte von KIKOMPAKT können unter Nennung der Quelle Plattform Lernende Systeme für redaktionelle Zwecke genutzt werden.

---

### Impressum

Expertise: Jessica Heesen, Kristian Kersting, Catrin Misselhorn

Redaktion: Paul Grünke, Birgit Obermeier

Herausgeber: Lernende Systeme – Die Plattform für Künstliche Intelligenz | Geschäftsstelle | c/o acatech | Karolinenplatz 4 | D-80333 München  
 presse@plattform-lernende-systeme.de | www.plattform-lernende-systeme.de

Stand: September 2025 | Bildnachweis: Thilo Schoch, Uni Darmstadt, Stefanie Trenz | Folgen Sie uns auf [LinkedIn](#) und [YouTube](#).

Gefördert durch:



Bundesministerium  
für Forschung, Technologie  
und Raumfahrt

 **acatech**  
DEUTSCHE AKADEMIE DER  
TECHNIKWISSENSCHAFTEN