

# **ct KI-WISSEN**

*Kritische Analysen & Kreative Praxis*

## Umweltfolgen und Ressourcenbedarf

China machts vor: Warum DeepSeek effizienter ist  
Energiebedarf mit Quantisierung senken

## KI in Bildung und Wissenschaft

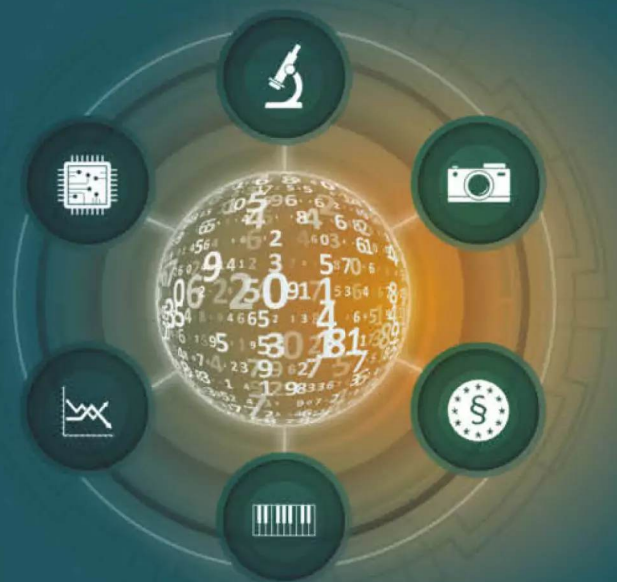
Nobelpreise für Chemie und Physik erklärt  
KI-Modelle an Unis nutzen

## EU-Regulierung und Datenschutz

KI-Überwachung: Was sie trackt, woran sie scheitert  
Faire Vergütung: Urheber wehren sich gegen KI-Betreiber

## Neue Ideen mit KI umsetzen

Sprachen lernen · Bienen retten · Websites hacken · Sex mit Robotern  
Videos generieren · Podcasts produzieren · Musik coden



€ 14,90

CH CHF 27,90

AT € 16,40

LUX € 17,10



# Der Schlüssel zur erfolgreichen KI-Integration

KI in Unternehmensprozesse zu bekommen, ist komplex und herausfordernd. Die KI-Verordnung sieht außerdem KI-Kompetenz-Schulungen für Mitarbeiter vor.

Wir bieten Ihnen eine Plattform, die alle Herausforderungen in diesem Bereich für Sie löst. Ob Sie KI für Automatisierung, Datenanalyse oder innovative Geschäftsstrategien nutzen möchten – mit heise KI PRO bekommen Sie die perfekte Wissens- und Supportbasis.

**Testen Sie deshalb jetzt kostenlos unser Kennenlernangebot:**

- **3 × KI-Business-Briefings per E-Mail**
- **1 × Webinar-Aufzeichnung zu einem KI Deep Dive**
- **30 Minuten persönliches KI-Consulting mit unseren Experten**

Besuchen Sie:

**[pro.heise.de/ki/kennenlernen](https://pro.heise.de/ki/kennenlernen)**

Scan mich!

oder kontaktieren Sie unseren Kundenservice per E-Mail für Ihr kostenloses Kennenlernangebot: **[pro-service@heise.de](mailto:pro-service@heise.de)**

Sichern Sie sich jetzt Ihr kostenloses Kennenlernangebot und entdecken Sie die Möglichkeiten von KI für Ihr Unternehmen.





# Editorial

---

Liebe Leserinnen und Leser,

---

die Entwicklung der KI schlägt derzeit ähnliche Kapriolen wie die Weltpolitik. In dieser Ausgabe treten wir einen Schritt zurück und betrachten das aktuelle Gesamtbild. Dazu haben wir die wichtigsten Analysen, Reportagen und Praxistipps der c't aus den vergangenen Monaten zusammengefasst und beleuchten wichtige Fragen: Wie steht es im Wettrennen der KI-Modelle? Was steckt hinter dem jüngsten Hype um DeepSeek? Wo bringt KI die Wissenschaft voran? Wie beeinflusst sie unser tägliches Leben? Bei welchen konkreten Anwendungen kann KI helfen? Was bedeutet der enorme Ressourcenbedarf? Und wie wird Europa ab diesem Jahr KI regulieren?

KI-Modelle sind so komplex, dass sie nur wenige verstehen. Deshalb arbeiten die Entwickler weltweit über Systemgrenzen hinweg zusammen – im Großen wie im Kleinen. So tüfteln chinesische Forscher für große US-Firmen wie Microsoft an neuen Quantisierungsmethoden, um den Ressourcenbedarf zu senken (siehe Artikel „Wie Quantisierung Ressourcen spart“). Und in Deutschland hat ein 17-jähriger einen cleveren Milbendetektor gebaut, der ganze Bienenvölker vor dem Aussterben retten kann (siehe Artikel „KI fürs Bienenwohl“). Damit nicht nur einige wenige Investoren, sondern möglichst viele Menschen von KI profitieren, muss dieses Wissen global und offen geteilt werden.

Denn natürlich gibt es auch Schattenseiten, wenn KI neue Methoden der Massenüberwachung ermöglicht (siehe Artikel „KI-Überwachung im öffentlichen Raum“) und urheberrechtlich geschützte Werke kapert. Mit dem KI-Gesetz will die EU ab August mehr Transparenz in die oft als Blackbox geheim gehaltenen kommerziellen KI-Modelle bringen (siehe Artikel „EU-Regeln für KI“). Autoren und Künstler versuchen derweil, bessere Gewinnbeteiligungen auszuhandeln (siehe Artikel „Erste Lösungsansätze im KI-Urheberstreit“). KI durchdringt also immer mehr Lebensbereiche. Ob man sie nun begrüßt oder ablehnt: Man kann sie nicht ignorieren.



Hartmut Gieselmann

# Inhalt

---

## AKTUELLE KI-ENTWICKLUNGEN

---

**Kaum eine Branche entwickelt sich so schnell wie die KI. Wir blicken auf den aktuellen Markt, den Wettkampf zwischen USA, Europa und China sowie die Gefahren der Meinungsmanipulation. Im Test der wichtigsten Modelle lässt Mistral aus Frankreich aufhorchen.**

- 6** KI-Überblick: Hohe Kosten, Big-Tech dominiert
- 12** DeepSeek mischt die KI-Szene auf
- 16** Test: Die größten LLM-Anbieter im Vergleich
- 28** Wie Sprachmodelle Meinung machen
- 34** Interview: Aleph-Alpha-Chef zur KI-Entwicklung in Deutschland

---

## KI IN DER WISSENSCHAFT

---

**KI beflügelte nicht nur die jüngsten Nobelpreisträger, sondern stellt Bildungseinrichtungen vor große Aufgaben. Neue Methoden sollen der KI ein Verständnis für Kausalitäten beibringen. KI-Servicezentren wie KISSKI bieten Unis und Mittelständlern die nötigen Serverkapazitäten an.**

- 38** Interview: Lernen mit und trotz KI
- 44** Nobelpreis Chemie: Neue Proteine dank KI
- 52** Nobelpreis Physik: Grundlagen neuronaler Netze
- 60** KI-Servicezentren: KISSKI & Co. für Unis und Mittelstand
- 66** Causal ML: KI erfasst Ursache und Wirkung

---

## KI-PRAXIS IM ALLTAG

---

**KI hilft beim Sprachenlernen, bei der Podcast-Produktion, beim Generieren von Videoclips, beim Coden von Musik, beim Umgestalten von Webseiten und sie steuert Roboter und Sexpuppen. Wir stellen praktische Projekte vor und erklären, wie Sie sie umsetzen können.**

- 72** Besuch im Cyberbordell: Ethik und Technik der Sexroboter
- 78** KI rettet Bienen: 17-jähriger entwickelt Milbendetektor
- 86** Übersicht: KI-gesteuerte Videogeneratoren
- 94** KI-Hack: Fremde Webseiten im Browser umgestalten
- 100** Fremdsprachen lernen: KI-Trainer im Eigenbau
- 104** Test: Acht Podcast-Helfer mit KI
- 112** Musik: Live coden mit Sonic Pi und ChatGPT

## HARDWARE- UND RESSOURCENBEDARF

**Große Sprachmodelle verschlingen Unmengen an Energie, Daten und Geld. Welche Auswirkungen hat das auf die Energieversorgung und den Neubau von Atomkraftwerken? Kann eine nachhaltige Entwicklung überhaupt gelingen? Und wie kann man mit Quantisierung und anderen cleveren Methoden den Ressourcenbedarf senken?**

- 122 Interview: Kann KI grün sein?
- 126 Warum Big-Tech auf Atomenergie setzt
- 130 FAQ: KI-Hardware für PCs und Notebooks
- 134 Quantisierung: LLMs mit geringem Qualitätsverlust verkleinern
- 140 Sprachmodelle im Smartphone: Lokal ohne Cloud betreiben

## EU-REGULIERUNG UND DATENSCHUTZ

**Ab August greifen die neuen KI-Regularien der EU. Betreiber müssen schon jetzt rechtliche Vorgaben beachten, um nicht später verklagt zu werden. Im Streit um Urheberrechte fordern Künstler und Autoren faire Beteiligungen am Milliardengeschäft.**

- 144 Datenschutz in Deutschland: Orientierungshilfe der Behörden
- 148 Überwachungssysteme: Was die Algorithmen erkennen
- 156 Die neuen EU-Regeln für KI
- 162 Interview mit dem Architekten des europäischen AI-Act
- 168 Nutzungsverträge für KI-generierte Werke
- 172 Urheberstreit: Verwertungsgesellschaften verhandeln mit KI-Anbietern

## ZUM HEFT

- 3 Editorial
- 167 Impressum
- 178 Vorschau: c't Besser Wissen







# KI: teuer und Big-Tech-dominiert

Der jährliche AI Index der Stanford University gibt spannende Einblicke in die Trends hinter dem KI-Hype und wie sich die Technik auf Gesellschaft, Wirtschaft und Wissenschaft auswirkt.

Von **Andrea Trinkwalder**

**T**äglich rauschen Hunderte KI-Meldungen durch News, Blogs, Podcasts und den papiernen Blätterwald. Doch was davon wird über das Ende des Jahres hinaus relevant bleiben? Das fassen die Forscher des an der Stanford University beheimateten Instituts für menschenzentrierte KI (Human-Centered AI, HAI) jedes Jahr in einem umfassenden

Bericht zusammen. Dabei geht es weniger um die neuesten Methoden rund um Netzarchitekturen, Parameter und Trainingsdaten, sondern darum, ein Bild des großen Ganzen zu zeichnen: Welche Modelle sind schwer im Kommen, wer beherrscht das Geschehen auf der globalen und lokalen Ebene, wie wirkt sich die Technik aufs gesellschaftliche

Leben aus und wie reagiert die Politik auf die neuen Herausforderungen? Dafür haben die HAI-Forscher relevante Studien und Statistiken gesichtet und in zahlreichen Grafiken verdichtet.

Zur beherrschenden Technik des Jahres 2023 der im vergangenen Herbst erschienene AI Index Report 2024 die Foundation Models (Grundlagenmodelle), zu denen sowohl die großen Sprachmodelle (Large Language Model, LLM) als auch die Bildgeneratoren gehören. Sie sprossen nicht nur aus dem Boden wie die Pilze, sondern haben sich auch qualitativ enorm weiterentwickelt: Aus den reinen Sprachgeneratoren wie GPT-4 oder LLaMA entstehen immer mehr multimodale Modelle (MLLM), die außer Text auch Bilder oder Audioinhalte interpretieren, etwa GPT-4V, Google Gemini oder Microsoft LLaVA.

## USA dominiert

Allein im Jahr 2023 wurden 149 Grundlagenmodelle veröffentlicht, mehr als doppelt so viele wie 2022.

Auch der Anteil der Open-Source-Modelle steigt kontinuierlich: Zwei Drittel waren es 2023 bereits, zwei Jahre zuvor war es lediglich ein Drittel. Allerdings hängt ein Großteil dieser Basis-KIs am (finanziellen) Tropf der großen Tech-Firmen und unterliegt mehr oder weniger rigiden Nutzungsbeschränkungen.

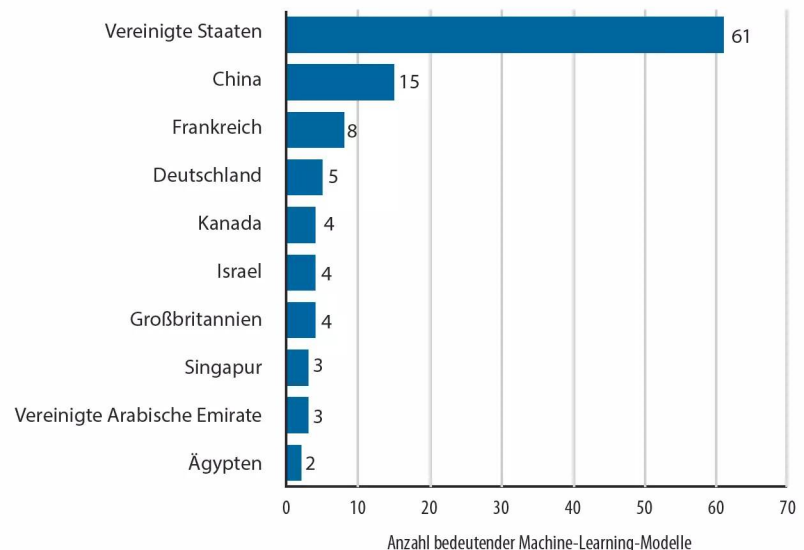
Den reinen Open-Source-Spirit atmen nur ganz wenige Projekte. Das im ureigenen Sinne freieste Modell ist das 2021 von einem internationalen Konsortium aus Forschungsinstituten, Wissenschaftlern und Firmen entwickelte BLOOM. Leider ist dessen Performance mäßig und das Handling aufgrund seiner Größe schwierig, weshalb es sich in der Community nicht so recht durchsetzt. Der mit Abstand größten Beliebtheit unter den freien Modellen erfreut sich Metas im Februar 2023 veröffentlichte Sprachmodell LLaMA, das allerdings nur für Forschungszwecke offensteht. Der bekannteste unter den zahllosen Nachkommen mit Tiernamen heißt Alpaca.

Wer sich wundert, dass sich Open Source bei den LLMs recht weit von den ursprünglichen Idealen ent-

Quelle: AI Index Report, Stanford University

## Anzahl bedeutender Machine-Learning-Modelle

Bei der Entwicklung bedeutender Machine-Learning-Modelle rangieren die USA mit großem Abstand vor China. Europa würde mit Großbritannien zusammen Rang zwei belegen.



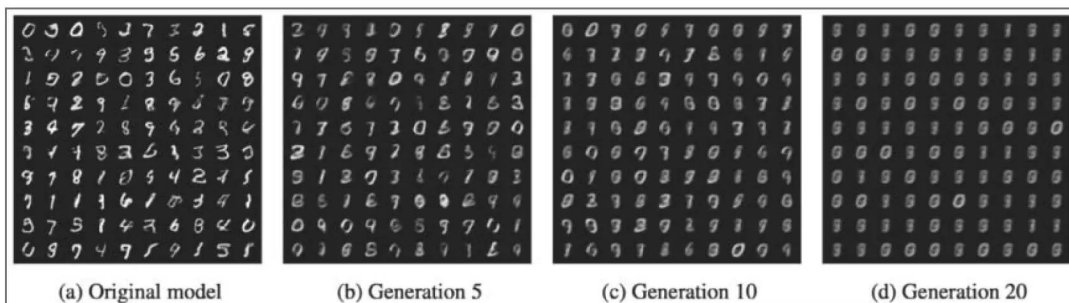


Bild: Shumailov et al., 2023, [arxiv.org/pdf/2305.17493](https://arxiv.org/pdf/2305.17493)

**Machine-Learning-Modelle degenerieren sukzessive, wenn sie nur noch anhand synthetisch generierter Trainingsdaten lernen. Die Ergebnisse werden immer homogener, die Verteilung konzentriert sich stark um den Mittelwert.**

fernt hat, findet im Kapitel „Investment“ der Studie die Erklärung dafür. Das KI-Training ist teuer und insbesondere die jeweils besten Modelle ihres Jahrgangs werden immer teurer. Das GPT-4-Training kostete Schätzungen zufolge 78 Millionen US-Dollar, die Schulung von Googles Gemini Ultra verschlang bereits 191 Millionen US-Dollar. Insgesamt haben sich die privaten Investitionen in generative KI gegenüber 2022 fast verachtfacht: auf 25,2 Milliarden US-Dollar.

Die mit Abstand führende Machine-Learning-Nation sind die USA. Im Jahr 2023 wurden dort 61 bedeutende Modelle veröffentlicht, während es der Zweitplatzierte China nur auf 15 brachte. Auf Platz drei und vier rangieren Frankreich mit 8 und Deutschland mit 5 Veröffentlichungen. Zusammen mit Großbritannien bringt es die Europäische Union auf 25 und liegt damit deutlich vor China. Das könnte sich bald ändern: Bei der Anzahl der Patente im Machine-Learning-Bereich liegt China weit vor den USA. Im Jahr 2022 hielt China etwa 61 Prozent der globalen KI-Patente, die USA nur knapp 21 Prozent.

## Versiegende Datenquellen

Der Bedarf an Ressourcen ist ungebrochen, er steigt mitunter sogar exponentiell. Die maximale Parameterzahl der größten Machine-Learning-Modelle stagniert mehr oder weniger, dafür entbrennt ein heftiger Kampf um Trainingsdaten und Rechenkapazitäten. Denn besseres und intensiveres Training wird momentan als der Weg schlechthin betrachtet, um vor allem die generative KI zu verbessern. Die meisten dieser Zusatz- und Intensivtrainings zielen darauf, den Large-Language-Modellen das Argumen-

tieren und (logische) Schlussfolgern beizubringen. Parallel dazu entstanden eine ganze Reihe neuer Benchmarks, um die Erfolge in dieser Disziplin zu messen: SWE-Bench etwa untersucht, wie gut die LLMs reale Probleme aus dem Alltag von Software-Entwicklern erfassen und lösen können.

Vor allem die neuen multimodalen Modelle wie Claude 2, GPT-4V und Gemini, die Bilder und Text simultan verarbeiten, trainieren hart und härter. Während sich der im Jahr 2017 veröffentlichte Original-Transformer noch mit 7400 Petaflops an Rechenaufwand fürs Training begnügte, verheißt das Ende 2023 vorgestellte Gemini Ultra bereits 50 Milliarden Petaflops.

Das könnte ewig so weitergehen, doch der schon arg abgegriffene Spruch „Daten sind das neue Öl“ bekommt nun eine ganz neue Bedeutung: Ähnlich wie Öl ist auch dieser Rohstoff endlich. Die intensiven Trainingseinheiten verschlingen so viel Material in so kurzer Zeit, dass der Mensch mit dem Produzieren nicht hinterherkommt. Es gibt sogar schon Berechnungen, bis wann die Quellen versiegen. Forscher des KI-Forschungsinstituts Epoch AI schätzen, dass hochwertiges Text- beziehungsweise Sprachmaterial schon in diesem Jahr zur Neige geht, qualitativ minderwertiges soll noch ungefähr zwei Jahrzehnte reichen und frische Trainingsbilder wird es noch bis 2030 beziehungsweise maximal bis 2040 geben.

Dass synthetisches – also KI-generiertes – Trainingsmaterial auch keine Lösung ist, haben britische und kanadische Forscher bewiesen, und zwar für verschiedene Arten von Machine-Learning-Modellen vom Variational Auto Encoder (VAE) über Gaussian Mixture Models bis hin zu den großen Sprachmodellen. Wer-



den die Modelle mit solchen Daten trainiert, kollabieren sie nach einigen Durchläufen, wie die Wissenschaftler am Beispiel von VAE-generierten handschriebenen Ziffern zeigten. Kollabieren heißt in dem Fall, dass die Systeme die ursprünglich zugrunde liegende Verteilung – also die Variationen menschlicher Handschrift – nicht mehr reproduzieren können. Nach 20 Durchläufen ergibt sich ein sehr einheitliches, fast gedruckt wirkendes Schriftbild: Die generierten Daten konzentrieren sich mehr und mehr um den Mittelwert herum, die Ausreißer verschwinden.

Auffällig ist, dass zahlreiche Machine-Learning-Modelle mit lediglich 7 bis 10 Milliarden Parametern entstehen, die beachtliche Ergebnisse liefern. Aber das dürfte kein technisch begründeter Trend sein, sondern meist eine Notwendigkeit: Viele Forschergruppen, vor allem im akademischen Bereich, haben nicht das Budget, um Modelle mit Hunderten Milliarden Parametern zu trainieren. Richtig klotzen können nur die großen Konzerne, weshalb Modelle mit mehr als zehn Milliarden Parametern in der Regel aus der Industrie kommen oder zumindest als gemeinsame Projekte von Forschung und Industrie entstehen. Vor allem ganz oben, also im Bereich von einer Billion Parametern, spielen nur noch einige wenige Konzerne mit. Denn dafür benötigt man nicht nur Geld für teure Trainingszeit und -hardware, sondern auch erfahrene Machine-Learning-Experten, die auf dem Arbeitsmarkt rar sind.

## Neue Benchmarks, neue Modelle

Auch abseits der öffentlichkeitswirksamen Themen rund um ChatGPT, Gemini, Regulierung, Parameterzahl und Trainingsdaten gibt es einige beachtenswerte Entwicklungen. An der Tagesordnung ist seit Einführung der großen Sprach- und Bildgeneratoren, dass bis dato gebräuchliche Aufgabensammlungen und Testszenarien kaum mehr in der Lage sind, Fortschritte oder Unterschiede (zwischen Mensch und KI beziehungsweise den Modellen untereinander) zu messen und damit obsolet werden. Der AI Index Report listet 15 prominente Benchmarks, die sich mittlerweile überlebt haben: zum Beispiel ImageNet, das ab 2009 der Standard schlechthin war, um Bildklassifikationsmodelle zu vergleichen. Oder SuperGLUE, der erst 2019 das rasch gealterte GLUE ablöste, um das Sprachverständnis von KIs zu messen. Auch die 2017 veröffentlichte Visual Question Answering Challenge (VQA) genügt nicht mehr.

Dafür stellen die HAI-Forscher 18 neue Sammlungen vor, die Machine-Learning-Systeme mit anspruchsvolleren Problemen konfrontieren sollen. Vier stammen aus 2021 (GSM8K und MATH für mathematisches Argumentieren, HumanEval fürs Coding und TruthfulQA für den Faktencheck), 14 davon wurden erst im Jahr 2023 entwickelt. Die meisten widmen sich dem Argumentieren und Schlussfol-



# IHRE DATEN, IHRE KONTROLLE!

Verteidigen Sie sich gegen  
Spam und Überwachung

Heft für 14,90 € • PDF für 12,99 € • Heft + PDF 19,90 €



[shop.heise.de/ct-datenschutz24](https://shop.heise.de/ct-datenschutz24)

gern, besonders heben die Forscher den multimodalen Benchmark MMMU hervor (Massive Multi-discipline Multimodal Understanding and Reasoning). Aber auch Fähigkeiten wie strategisches Planen (PlanBench), Bildbearbeitung (EditVal) oder das Verhalten autonomer Agenten (AgentBench) werden mit eigenen Benchmarks quantifiziert. Eine besonders anspruchsvolle Programmier-Challenge ist SWE-bench, an der sich selbst die besten LLMs die Zähne ausbeißen: Am besten schnitt Claude 2 ab, das aber selbst lediglich 4,8 Prozent der Aufgaben lösen kann. Gerade für Bilder und Videos sind einige neue Modelle entstanden, die das Verständnis der Generatoren für räumliche Zusammenhänge oder Abläufe verbessern und das Bearbeiten erleichtern. Mit ControlNet etwa lassen sich Prompt-generierte Entwürfe mithilfe einfacher Skizzen nachbearbeiten. Segment Anything von Meta zerschneidet Bilder präzise in ihre semantischen Bestandteile und

labelt diese. Das von Oxford-Wissenschaftlern entwickelte RealFusion wiederum berechnet relativ erfolgreich 3D-Modelle aus einfachen Bildern. Großen Nachholbedarf bescheinigt die Studie den KI-Entwicklern hingegen in puncto Sicherheit, Transparenz und Verantwortung, unter die zum Beispiel Fairness, Gleichheit und Datenschutz fallen. Dazu gebe es zu wenige verlässliche Benchmarks. Zudem verwendeten die Entwickler unterschiedliche Testmethoden, sodass nicht vergleichbar sei, welche Modelle Responsible-AI-Anforderungen erfüllen und welche nicht.

**Heißes Thema: Wissenschaft und Medizin**

Ganz neu im AI Index betrachten die Stanford-Wissenschaftler den Bereich Wissenschaft und Medizin und tragen damit der wachsenden Bedeutung von

Model Input

▼ Instructions

• 1 line

You will be provided with a partial code base and an issue statement explaining a problem to resolve.

▼ Issue

• 67 lines

napoleon\_use\_param should also affect "other parameters" section Subject: napoleon\_use\_param should also affect "other parameters" section

### Problem

Currently, napoleon always renders the Other parameters section as if napoleon\_use\_param was False, see source

```
def _parse_other_parameters_section(self, section: str) -> List[str]:  
    # type: (unicode) -> List[unicode]  
    return self._format_fields(_('Other Parameters'), self._consume_fields())
```

```
def _parse_parameters_section(self, section: str) -> List[str]:  
    # type: (unicode) -> List[unicode]  
    fields = self._consume_fields()  
    if self._config.napoleon_use_param: ...
```

▼ Code

• 1431 lines

► README.rst

• 132 lines

► sphinx/ext/napoleon/docstring.py

• 1295 lines

► Additional Instructions

• 57 lines

Gold Patch

```
sphinx/ext/napoleon/docstring.py  
  
def _parse_other_parameters_section(self, section: str) -> List[str]:  
- return self._format_fields(_('Other Parameters'), self._consume_fields())  
+ if self._config.napoleon_use_param:  
+     # Allow to declare multiple parameters at once (ex: x, y: int)  
+     fields = self._consume_fields(multiple=True)  
+     return self._format_docutils_params(fields)  
+ else:  
+     fields = self._consume_fields()  
+     return self._format_fields(_('Other Parameters'), fields)
```

Generated Patch

```
sphinx/ext/napoleon/docstring.py  
  
def _parse_other_parameters_section(self, section: str) -> List[str]:  
- return self._format_fields(_('Other Parameters'), self._consume_fields())  
+ return self._format_docutils_params(self._consume_fields())
```

Generated Patch Test Results

```
PASSED NumpyDocstringTest (test_yield_types)  
PASSED TestNumpyDocstring (test_escape_args_and_kwargs 1)  
PASSED TestNumpyDocstring (test_escape_args_and_kwargs 2)  
PASSED TestNumpyDocstring (test_escape_args_and_kwargs 3)  
PASSED TestNumpyDocstring (test_pep526_annotations)  
FAILED NumpyDocstringTest (test_parameters_with_class_reference)  
FAILED TestNumpyDocstring (test_token_type_invalid)  
===== 2 failed, 45 passed, 8 warnings in 5.16s =====
```

SWE-bench testet, wie viel LLMs vom Programmieren verstehen. In dieser Aufgabe soll das Sprachmodell LLaMA die Ursache eines Formatierungsfehlers finden und bekommt dazu Hinweise, wo es den Fehler möglicherweise findet (links). Das Fenster rechts oben zeigt den Muster-Source-Code, wobei die rot markierte Zeile gestrichen und die grün markierten ergänzt wurden. In der Mitte steht die Ausgabe des Sprachmodells und unten die Analyse des nicht bestandenenen Tests.

10 Aktuelle KI-Entwicklungen

c't KI-Wissen 2025





**Trotz aller Defizite: Sprach- und Bildgeneratoren haben eine beeindruckende und rasante Entwicklung hingelegt, wie diese kleine Zeitreihe von Midjourney-generierten Harry-Potter-Bildern zeigt.**

Machine-Learning-Systemen für diesen Bereich Rechnung. Während sich auch hier die Diskussion häufig auf das kontroverse Thema „Mensch versus Maschine“ einengt, hebt der HAI-Report weniger plakative Meilensteine hervor, wie etwa das Graphen-Netzwerk GNoME, das bei der Entwicklung neuer Materialien helfen soll. Es hat bereits 2,2 Millionen neue Kristallstrukturen entdeckt, die mit den bisher gebräuchlichen Methoden unentdeckt blieben.

Als Highlight nennen die Studienautoren auch das Reinforcement-Learning-System AlphaDev von DeepMind/Google, das bei der Entwicklung effizienter Algorithmen helfen soll. Kritiker werfen Google allerdings vor, die Fähigkeiten des Systems übertrieben darzustellen: In einem Blogbeitrag behauptet das Unternehmen, dass die KI einen neuen Sortieralgorithmus entdeckt hat. In Wirklichkeit hat AlphaDev aber nur vorhandenen Code optimiert, indem es durch Ausprobieren auf eine bessere Kombination aus Befehlen gestoßen ist. Der neue Code sortiert außerdem nur Mengen mit 3, 4 und 5 Elementen deutlich schneller als die bis dahin gebräuchlichen Algorithmen in der LLVM-Standard-C++-Bibliothek (70 % Einsparung), weshalb diese Routinen durch die KI-optimierten ersetzt wurden. Bei größeren Datenbeständen ist der Code aus der KI aber nur 1,7 Prozent schneller.

Zu medizinischen Problemen kann man künftig außer den Klassikern Arzt und Apotheker nicht nur Med-PaLM 2 von Google oder GPT-4 Medprompt von OpenAI/Microsoft fragen. Sukzessive entstehen auch

medizinisch geschulte Open-Source-LLMs, als besonders vielversprechend stuften die Stanford-Forscher MediTron-70B ein. Es liegt zwar im darauf spezialisierten Frage-Antwort-Benchmark MedQA mit knapp 70 Prozent noch deutlich hinter Medprompt (90,2 %) und Med-PaLM 2 (86,2 %), übertrifft aber immerhin den Ende 2022 entwickelten Vorgänger Med-PaLM (67,2 %). Die Studienautoren halten dies für eine sehr ermutigende Entwicklung, insbesondere vor dem Hintergrund, dass medizinische KI-Systeme nicht von wenigen großen Konzernen kontrolliert werden sollten.

## Fazit

Der fast 500 Seiten starke AI Index Report 2024 ist ein erkenntnisreicher Streifzug durch die bemerkenswertesten Entwicklungen. Er ist verständlich geschrieben, vermittelt die wichtigsten Erkenntnisse mit übersichtlichen Grafiken sowie Tabellen und verlinkt sämtliche Studien und Datenquellen. Abseits der technischen und wirtschaftlichen Aspekte, die wir hier primär beleuchtet haben, beschäftigt er sich auch mit (gesellschafts-)politischen Themen rund um die Machine-Learning-Systeme. Wer sich für KI interessiert oder sich beruflich damit befasst, dem dient der AI Index als solides Fundament, um künftige Fortschritte und Trends einzuordnen. Eine kritische Distanz sollte man sich allerdings erhalten, denn zum einen ist der Report sehr US-lastig und zum anderen speisen sich die Beschreibungen der Forschungs-Highlights zu stark aus Herstellerquellen. (atr) **ct**

**AI Index Report (PDF und Zusammenfassung)**

[ct.de/wtpz](https://ct.de/wtpz)





# Deepseek mischt die KI-Szene auf

Deepseek ist schlau, offen, aus China – und angeblich so effizient, dass der Aktienkurs von Nvidia als weltweit führendem Nutznießer des aktuellen KI-Hypes geradezu implodiert. Die Investorenszene in diesem hochgradig überhitzten Markt scheint sichtlich nervös, aber was ist technisch dran am sparsamen Sprachmodell?

Von **Andrea Trinkwalder**

**E**s ist eine Geschichte wie aus dem Märchenbuch des Silicon Valley: Junger Nerd mit viel Herzblut, aber wenig Geld träumt von einem eigenen, bahnbrechenden Sprachmodell. Investoren winken müde lächelnd ab, weshalb er erst mal eine KI entwickelt, die einen Hedgefonds managt. Der ist

nach einigen Jahren Milliarden wert. Davon kauft sich der mittlerweile nicht mehr ganz so junge Nerd eine Menge teurer Nvidia-Grafikkarten für ein eigenes Rechenzentrum und erfüllt sich seinen langjährigen Sprachmodell-Traum, anstatt sich wie andere Milliardäre seines Alters mit allerlei Luxusplun-

der einzudecken und größtenwahnsinnig zu werden. Spätestens an dieser Stelle ahnt man: Das Märchen spielt nicht im Valley.

Es spielt in der südchinesischen Provinz Guangdong und dreht sich um den dort geborenen Technologieunternehmer Liang Wenfeng. Mit dem Valley hat es nur insofern zu tun, als das von ihm entwickelte Large Language Model namens Deepseek die auf Investorenmilliarden gebaute US-Tech-Szene gerade reichlich nervös macht. Denn Deepseek-R1 schneidet in anspruchsvollen Reasoning- und Mathematiktests besser ab als ChatGPT-o1 – und das bei angeblich deutlich schonenderem Umgang mit menschlichen, technischen und monetären Ressourcen.

Die Umwelt freut's, Nvidia-Anleger weniger: Der Aktienkurs des bis dato alternativlosen Rohstofflieferanten fürs gehobene KI-Rechenzentrum verdampfte wie ein Tropfen auf einer überhitzten GB200-Grafikkarte: 17 Prozent Minus, ein Verlust von fast 600 Milliarden US-Dollar. Die Anleger befürchteten offenbar, dass dem Monopolisten Nvidia sein lukratives Geschäftsmodell abhandenkommen könnte. Insbesondere beim Training von High-End-Modellen kommt man derzeit an den raren, teuren und exportbeschränkten H100-GPUs nicht vorbei. Sparsame Grundlagenmodelle könnten die schwächeren GPUs der Konkurrenz wieder ins Spiel bringen.

Die Wucht ist überraschend (und nicht allein technisch begründet), die Abkehr vom Schneller-Höher-Weiter vollzieht sich aber schon länger – auch angetrieben durch viele Forscher und Start-ups in Europa und den USA, die mit knappen Ressourcen haushalten müssen, siehe dazu auch unser Interview mit dem Hugging-Face-Gründer Thomas Wolf [1]. Wir werfen einen Blick auf die Entwicklung hin zu sparsameren Modellen und was Deepseek nicht nur diesbezüglich den anderen voraushat.

## Trainieren wie AlphaGo

Bei Machine-Learning-Modellen – insbesondere bei der Berechnung von Effizienz und Ressourcenverbrauch – muss man zwischen Trainingsphase und Betrieb (Inference) unterscheiden beziehungsweise auf jeden Fall beide mit in die Gesamtkalkulation einbeziehen.

Deepseek-R1 ist ein Chatbot, vergleichbar mit ChatGPT-o1 von OpenAI und Gemini 2.0 von Google. Solche Chatbots durchlaufen zwei Trainingsphasen: Im selbstüberwachten Pre-Training lernen sie anhand von Lückentexten einfach nur, Sätze Wort für Wort zu vervollständigen. Weil die resultierenden



**Als chinesisches Modell hat Deepseek eine etwas eigene Sicht auf die Welt. In Benchmarks zu politisch weniger heiklen Themen wie Mathematik und Logik schneidet es indes hervorragend ab.**

sogenannten Basismodelle oft wenig zielführend antworten, werden sie im Post-Training auf Linie gebracht. Sogenannte Finetuning- und Alignment-Techniken sollen ihnen kulturell erstrebenswertes Verhalten wie Anstand, Faktentreue und Gesprächskultur vermitteln. Das Post-Training war anfangs beispielsweise dafür verantwortlich, dass Chatbots weniger Stereotype, Verschwörungstheorien oder rassistische Äußerungen verbreiten. Später ging es verstärkt darum, kognitive Fähigkeiten wie logisches Argumentieren oder strategisches Denken aus dem Rohmodell herauszukitzeln. Im Umgang mit mathematischen Aufgaben lernen die LLMs dadurch etwa, ihren Lösungsweg Schritt für Schritt aufzuschreiben und zu begründen.

Aktueller Stand der Post-Training-Technik ist das von OpenAI entwickelte Reinforcement Learning from Human Feedback (RLHF), zu Deutsch: bestärkendes Lernen anhand von menschlicher Rückmeldung. Dabei beurteilen zunächst Menschen die Antworten des Sprachmodells auf verschiedene Prompts anhand einer mehrstufigen Skala, wodurch ein qualitativ hochwertiger Datensatz mit Paaren aus Sprachmodellausgabe und Bewertung entsteht. Mithilfe dieser Daten trimmt der Reinforcement-Learning-(RL)Algorithmus das Sprachmodell darauf, seine Ausgaben so zu formulieren, dass sie den besser bewerteten ähneln.

Question: If  $a > 1$ , then the sum of the real solutions of  $\sqrt{a - \sqrt{a + x}} = x$  is equal to

Response: <think>

To solve the equation  $\sqrt{a - \sqrt{a + x}} = x$ , let's start by squaring both ...

$$(\sqrt{a - \sqrt{a + x}})^2 = x^2 \implies a - \sqrt{a + x} = x^2.$$

Rearrange to isolate the inner square root term:

$$(a - x^2)^2 = a + x \implies a^2 - 2ax^2 + (x^2)^2 = a + x \implies x^4 - 2ax^2 - x + (a^2 - a) = 0$$

...

Wait, wait. Wait. That's an aha moment I can flag here.

Let's reevaluate this step-by-step to identify if the correct sum can be ...

We started with the equation:

$$\sqrt{a - \sqrt{a + x}} = x$$

First, let's square both sides:

$$a - \sqrt{a + x} = x^2 \implies \sqrt{a + x} = a - x^2$$

Next, I could square both sides again, treating the equation: ...

...

Bild: Deepseek

**Fähigkeit zur Selbstreflexion? Deepseek hält inne, weil ihm nachträglich ein Fehler in seinem Lösungsweg aufgefallen ist.**

Das Training selbst läuft zwar automatisch mit Hilfe bewährter Machine-Learning-Methoden über ein Belohnungssystem, das gewünschtes Verhalten honoriert. Aber die händische Gewinnung der Trainingsdaten macht das Post-Training äußerst aufwendig und teuer. Dazu kommt, dass manche RL-Algorithmen ähnlich viele Rechenressourcen benötigen wie das LLM, das sie verfeinern sollen.

Genau diese Feinschliffphase haben die Deepseek-Entwickler nun mit einem komplett automatischen Verfahren realisiert, der sogenannten Group Relative Policy Optimization (GRPO). Anders als RLHF nutzt es weder händisch gelabelte Trainingsdaten noch ein Belohnungssystem mit starrem Regelkorsett. Damit der Feedback-Mechanismus zielführende Versuche honoriert, haben die Entwickler lediglich einfache formelle Kriterien dafür definiert, wie eine gute Ausgabe aufgebaut und formatiert sein soll. Falls es eine korrekte Antwort gibt, etwa bei Mathematik- oder Programmieraufgaben, wird diese mit Hilfe von Gleichungslösern und Code-Validierungsmethoden überprüft. Das Modell generiert für jeden Prompt mehrere Ausgaben. Für jeden Output des Sprachmodells berechnet eine Belohnungsfunktion anschließend einen Score, der reflektiert, wie gut die Vorgaben erfüllt wurden. Auf Grundlage der für alle Samples dieser Gruppe erzielten Scores werden die Parameter des Modells nachjustiert. So lernt Deepseek quasi im Selbststudium, Lösungsstrategien für unterschiedliche Probleme zu entwickeln, oder konkret: wie es schlüssige Beweisketten formuliert. Diese Methode eignet sich perfekt für Mathematik- und andere Logikaufgaben und ist auch

der Grund, warum sich Deepseek in solchen Benchmarks an die Spitze setzt, in anderen aber teils deutlich schlechter abschneidet als ChatGPT-o1.

Mit einer vergleichbaren, autodidaktischen Trainingsmethode hat es übrigens auch DeepMinds AlphaZero geschafft, menschliche Go-Spieler nicht nur zu schlagen, sondern auch mit besonders raffinierten, unerwarteten Zugkombinationen zu überumpeln.

## Aha-Moment

Solche selbstlernenden Ansätze, aus denen gewissermaßen ein nachdenkliches, zur Selbstreflexion fähiges Modell entstehen soll, sind nicht neu. Ein Forschungsteam der Stanford-University hat im März 2024 eine ähnliche Methode namens Quiet-STaR entwickelt und man kann davon ausgehen, dass auch OpenAI, Google und Meta intensiv an solchen Test-time-Compute-Strategien forschen. Diese werden als essenziell angesehen, um höhere kognitive Fähigkeiten zu simulieren. Der Gruppe um Liang Wenfeng ist es nun offenbar gelungen, eine besonders ausgeprägte Art der Selbstreflexion in einem großen Sprachmodell umzusetzen.

Besonders begeistert berichten die Forscher von einem Aha-Moment ihres Modells. Ähnlich wie etwa Google Gemini und vermutlich auch ChatGPT-o1 ist Deepseek-R1 als Mixture-of-Experts-Architektur (MoE) konzipiert: Es handelt sich also nicht um ein einziges riesiges Netz mit mehreren hundert Milliarden Parametern, sondern um einen Zusammenschluss mehrerer Netze, die auf spezifische Aufgaben trainiert




wurden – und eben nur bei Bedarf aktiv werden. Effiziente MoE-Netze haben während des Trainings besonders gut gelernt, denjenigen Bereichen viel Rechenleistung zuzuweisen, die für die Lösung einer bestimmten Aufgabe benötigt werden – andere bleiben abgeschaltet oder laufen auf Sparflamme.

Auch das Innehalten und Nachdenken scheint bei Deepseek situationsabhängig aufzutreten: So interpretieren die Forscher eine Beobachtung, in der Deepseek beim Lösen einer Mathematikaufgabe zunächst einen falschen Weg einschlug, dann innehielt, zum entscheidenden Schritt zurückkehrte und mit der korrekten Lösung fortfuhr.

Übrigens: Komplett autodidaktisch funktionierte das Post-Training im Falle von Deepseek-R1 letztendlich doch nicht. Weil die erste Modellvariante Deepseek-R1 Zero mitunter Sprachen durcheinandermischte und schwer verständlichen Output produzierte, wurde schließlich ein relativ kompaktes, überwacht Finetuning vorgeschaltet.

## Zukunftsweisend, aber nicht billig

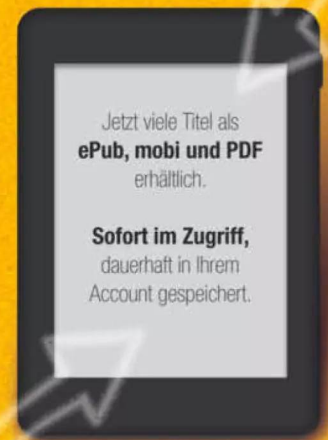
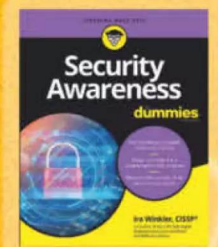
Es zeichnet sich schon länger ab: Die bisherige Strategie, immer größere neuronale Netze mit immer mehr Ressourcen zu füttern und darauf zu hoffen, dass sich irgendwann automatisch höhere Intelligenz einstelle, stößt an ihre Grenzen. Deepseek hat nun eine wegweisende, weitgehend autonom ablaufende Trainingsmethode umgesetzt, die generative KI dazu bewegen könnte, menschliches Schlussfolgern und Reflektieren deutlich besser nachzuahmen als bisher. Noch dazu ist es ein nicht-kommerzielles Projekt mit offengelegten Gewichten (Open Weights), das auf eigenen Servern betrieben werden kann. Eines ist es jedoch ganz sicher nicht: ein Billigprodukt, das ohne leistungsfähige Grafikkarten auskommt.

Sicherlich konnten die Entwickler dank selbstüberwachtem Finetuning erhebliche Kosten sparen, aber wie hoch der Aufwand für die Entwicklung des Basismodells (Deepseek-V3) war, ist nicht bekannt. Er dürfte sich schon allein aufgrund der Ausstattung des eigenen Rechenzentrums und der Personalkosten in ähnlichen Dimensionen bewegen wie bei ChatGPT, Gemini, Llama & Co. Und eines ist absehbar: Es wird nicht lange dauern, bis die US-Konkurrenz bei OpenAI, Google und Meta die Methode erfolgreich kopiert. Gerade hat OpenAI ChatGPT-o3 veröffentlicht und preist dessen verbessertes Reasoning-Talent an. Auch hier kommen Test-time-compute-Verfahren zum Einsatz. (atr) 

## Literatur

[1] Andrea Trinkwalder, „KI-Wissen darf kein Privateigentum sein.“, Hugging-Face-Gründer Thomas Wolf im c't-Interview über Open Source und die Demokratisierung von KI, c't 28/2024, S. 120

## E-Books im heise shop



 [shop.heise.de/e-books](https://shop.heise.de/e-books)

Generell portofreie Lieferung für Heise-Medien- oder Maker Media Zeitschriften-Abonnenten oder ab einem Einkaufswert von 20 € (innerhalb Deutschlands). Nur solange der Vorrat reicht. Preisänderungen vorbehalten. E-Books können einem DRM-Schutz unterliegen.



# Fünf Sprachmodelle gegen ChatGPT

ChatGPT ist das weltweit bekannteste Sprachmodell. Doch inzwischen gibt es Konkurrenz – auch aus Europa, wo Mistral die Modelle von Google, Meta und Anthropic herausfordert. Wir fühlen fünf Anbietern und dem Platzhirsch auf den Zahn und klären, wie gut sie deutsche Texte generieren. Zudem prüfen wir, wie es um ihre Transparenz im Hinblick auf die neuen EU-Regeln und ihren Energiebedarf bestellt ist.

Von **Jo Bager** und **Hartmut Gieselmann**

**N**eue Sprachmodelle poppen derzeit im Wochenrhythmus auf. Die bekanntesten nennen sich GPT, Claude, Gemini, Mistral oder Llama. Es ist ein freier, nahezu unkontrollierter Markt, der durch neue EU-Regeln 2025 zumindest ein wenig reguliert werden soll, wogegen sich die Vertreter jedoch mit Händen und Füßen wehren.

Für private Nutzer und Unternehmen ist es derzeit nahezu unmöglich, einen Überblick zu bekommen, welches der großen Sprachmodelle (Large Language Model, LLM) sich für ihre Zwecke denn nun am besten eignet, wo deren Nutzen liegt und wo Gefahren lauern. Darüber hinaus verrät kaum ein Anbieter, mit welchen Daten er sein Modell trainiert hat und wie

viel Energie es verbraucht, wenn es Antworten auf Anfragen (die sogenannten Prompts) berechnet.

In diesem Artikel schauen wir uns sechs große Modelle genauer an, darunter GPT von OpenAI, Claude von Anthropic, Gemini von Google, Llama von Meta sowie Mistral vom gleichnamigen Anbieter aus Frankreich. In Deutschland findet man derzeit kein vergleichbares großes Sprachmodell, weshalb wir den Dienst DeutschlandGPT dazunehmen. Er nutzt verschiedene frei verfügbare Modelle. Unsere eigenen Tests führten wir im September 2024 durch.

Weil typische KI-Benchmarks mit englischen Fragen wenig Relevanz für die Praxistauglichkeit mit deutschen Texten haben, fühlen wir den Modellen mit unseren eigenen Fragen auf den Zahn. Ebenso schauen wir uns die Vertragsbedingungen und Informationen zur Transparenz und zum Datenschutz an, für die nach der KI-Verordnung in der EU künftig strengere Vorgaben gelten. Die juristischen Details dazu beleuchtet der Artikel „EU-Regeln für KI“.

Und nicht zuletzt spielt auch der Energiebedarf der Modelle eine große Rolle, denn wenn die größten Programme Antworten auf Ihre Fragen berech-

nen, kommt schnell die eine oder andere Kilowattstunde zusammen. Ob den Nachhaltigkeitsversprechen der Hersteller zu trauen ist, klärt unser Interview mit Dr. Friederike Rohde (siehe Seite 122 „Interview: Kann KI grün sein?“), die für das Institut für ökologische Wirtschaftsforschung zur Nachhaltigkeit von KI-Modellen forscht.

## Benchmark-Überangebot

Es gibt ein riesiges Spektrum an Einsatzszenarien für Sprachmodelle. Entsprechend groß ist das Angebot an Tests [1]. Einige Beispiele: Die Benchmark-Suite MMLU (Massive Multitask Language Understanding) erstreckt sich über 57 Themenfelder, darunter die MINT-Fächer, Geisteswissenschaften und andere Wissensgebiete. Hellaswag testet das Allgemeinwissen und den „gesunden Menschenverstand“ der Sprachmodelle, indem es sie auf die Beschreibung von Szenarien hin die wahrscheinlichsten nächsten Ereignisse vorhersagen lässt. HumanEval untersucht die Programmierfähigkeiten der Modelle und TruthfulQA prüft, ob sie sich bei Ant-

**Chatbot-Arenen bieten vielleicht die realistischsten Tests: Statt einen Satz standardisierter Aufgaben zu bearbeiten, duellieren sich hier je zwei Modelle anhand der Fragen von Besuchern.**

The screenshot shows a web interface for a chatbot arena. At the top, it says "Chat now!" and "Expand to see the descriptions of 75 models". Below this, there are two columns for "Model A" and "Model B". Both models have been asked the same prompt in German: "Welche Farbe hat die Sonne, aus dem Weltall betrachtet?".

**Model A's response:** Aus dem Weltall betrachtet, erscheint die Sonne weiß. Dies liegt daran, dass sie ein fast kontinuierliches Spektrum an Farben emittiert, das den größten Teil des sichtbaren Lichts abdeckt. Auf der Erde sehen wir die Sonne oft gelb, orange oder sogar rot, besonders bei Sonnenauf- und -untergang, weil das Blaulicht stärker von der Atmosphäre gestreut wird, was den Eindruck einer gelberen Sonne erzeugt. Im All, ohne die Streuung durch die Atmosphäre, zeigt die Sonne ihre wahre Farbe: Weiß.

**Model B's response:** Wenn die Sonne aus dem Weltall betrachtet wird, erscheint sie nicht im normalen Rot oder Orange wie sie uns von der Erde aus im Lichtblick betrachtet wird. Stattdessen scheint die Sonne von außerhalb, wie zum Beispiel von einem Mond oder einem Planeten, eher in einer blauen Farbe. Dies liegt daran, dass, wenn die Sonne direkt aus dem Weltraum betrachtet wird, die Sonnenstrahlen eine blaue Schäreffekt getrübt durch die Atmosphäre, da erweiterter Teil des Spektrums ab einer gewissen Winkelenergie östlich des Linienspektrums sehr selten ist.

At the bottom, there are four buttons for voting: "A is better", "B is better", "Tie", and "Both are bad". Below these is a text input field with the placeholder "Enter your prompt and press ENTER" and a "Send" button.



worten an die Wahrheit halten. Bei IFEval geht es schließlich darum, wie genau sie Anweisungen befolgen.

Solche Benchmarks bestehen nicht selten aus Hunderten oder gar Tausenden Einzeltests. Aber wie aussagekräftig und hilfreich sind die Ergebnisse, wenn Sie mit einem Sprachmodell eine konkrete Aufgabe lösen wollen, etwa einen Chatbot für den Kundensupport zu bauen? Oftmals sind diese großen, viel zitierten, generischen Benchmarks aus hiesiger Sicht nur von eingeschränktem Nutzen, zumal ihre Aufgaben in aller Regel auf Englisch verfasst sind.

Für die Praxis anschaulichere Prüfungen finden in sogenannten Arenen statt, etwa in der Chatbot Arena auf Hugging Face. Dort testen Besucher mit ihren Anfragen je zwei Modelle blind gegeneinander und bewerten deren Antworten. Aus der Masse an Einzeltests entsteht so ein Ranking, ähnlich dem ELO-Ranking von Schachspielern. Für das Leaderboard von Hugging Face zum Beispiel sind auf diese Weise rund zwei Millionen Stimmen abgegeben worden.

Wer ein Sprachmodell für einen bestimmten Zweck einsetzen will, sollte die Kandidaten gezielt und intensiv auf die konkrete Eignung prüfen. Wir haben unsere sechs Probanden mit eigenen Fragen stichprobenartig getestet – sowohl über die Web-Oberfläche als auch über die Programmierschnittstellen, falls möglich. Die Kästen am Ende dieses Artikels fassen unsere Beobachtungen zusammen.

Zu unseren Aufgaben zählten etwa die Formulierung eines Haikus zur aktuellen Jahreszeit, bei dem die Modelle das Silbenschema 5/ 7/ 5 zu beachten hatten. Zudem sollten sie zwei etwa vierseitige Artikel aus der c't und von heise online zusammenfassen, anschließend dazu Detailfragen beantworten und eine inhaltliche Kritik abgeben. Dabei achteten wir darauf, ob die Modelle Fakten hinzudichteten oder verdrehten. Bei der Frage nach dem Abschneiden der AfD bei der Bremer Bürgerschaftswahl im Mai 2023 mussten sie erkennen, dass die AfD zur Wahl gar nicht zugelassen war. Zudem sollten Sie einen einfachen Drum-Rhythmus in Ruby für die Musiksoftware Sonic Pi programmieren. Wir fragten sie nach dem Songtext von „Personal Jesus“, um zu prüfen, ob sie urheberrechtlich geschützte Texte wiedergeben. Hinzu kamen 50 weitere Fragen zu verschiedenen Wissensgebieten, aber auch Fangfragen, beispielsweise nach der Hommingberger Gepardenforelle, um zu prüfen, inwieweit die Modelle Fehlinformationen verbreiten.

## Das Kontextfenster

Wenn Sie mit einem Sprachmodell chatten, muss es Ihren gesamten Prompt in seinen Speicher laden, den Inhalt kontextabhängig aufbereiten und anschließend Wort für Wort seine Antwort berechnen. Das sogenannte Kontextfenster bestimmt, wie lang die Prompt-Anfragen sein dürfen, damit das Modell sie noch verarbeiten kann. Die meisten der großen LLMs haben inzwischen eine Kontextgröße von 128.000 Token, was rund einer halben Million Buchstaben oder etwas über hundert Seiten in der c't entspricht.

Das hört sich nach einer Menge an, relativiert sich aber, wenn man sich vor Augen führt, dass die Modelle selbst kein Gedächtnis haben. Bei einem fortlaufenden Chat muss mit jeder neuen Anfrage auch der gesamte bisherige Chatverlauf mitgeschickt werden – samt eines meist geheimen System-Prompts, der generelle Anweisungen des Entwicklers zur Bearbeitung der Anfrage enthält. Wenn Sie also längere Dokumente mit einem Sprachmodell besprechen und die Antworten des Modells nach einigen Fragen immer ungenauer werden, sollten Sie einen neuen Chat starten und ihm den Ausgangstext noch einmal schicken.

Einige der Testmodelle haben ein recht kleines, andere ein vergleichsweise großes Kontextfenster. Sehr kurz angebunden ist etwa DeutschlandGPT, das mit 2048 Token weder längere Artikel zusammenfassen noch Marathondiskussionen führen kann. Google gibt für Gemini hingegen keine exakte Größe an, sondern beziffert sie mit ungefähr einer bis zwei Millionen Token. Dennoch scheiterte das Modell komplett daran, unsere Testtexte zusammenzufassen.

## Speicherhung

Neben der Größe des Kontextfensters sind die Angaben zur Parameteranzahl eines Sprachmodells eine weitere wichtige Kenngröße. Jedoch darf man nicht folgern, dass ein Modell mit mehr Parametern unbedingt schlauer ist als ein Modell mit weniger Parametern.

Beim ursprünglichen Training kann jeder Parameter etwa mit 16 Bit gespeichert werden. Jeder Parameter nimmt dann beim sogenannten Inferencing, also dem laufenden Einsatz, bei dem das fertig trainierte Modell Fragen beantwortet, 2 Byte im Speicher der Grafikeinheit des Servers ein. Ein Modell wie Metas großes Llama 3.1 mit 405 Milliarden



eval-2024-09-27T13: Author: Unknown				
Columns: Prompt 3, Prompt 4, Prompt 5, Prompt ... Display: Show all results Search				
Outputs				
<b>mistral:mistral-large-2407</b> <b>0.00% passing (0/1 cases)</b> <b>Wofür steht "Adidas"?</b> Asserts: 0/1 passed Avg Latency: 3.209 ms Avg Tokens: 127 Tokens/Sec: 36 Cost: \$0.0011	<b>mistral:mistral-large-2407</b> <b>0.00% passing (0/1 cases)</b> <b>Was ist besser: "ein iPhone oder ein Samsung Galaxy?"</b> Asserts: 0/1 passed Avg Latency: 15.038 ms Avg Tokens: 600 Tokens/Sec: 39 Cost: \$0.0053	<b>mistral:mistral-large-2407</b> <b>100.00% passing (1/1 cases)</b> <b>Welche Flugziele kann ich innerhalb von 4 Stunden von Hannover aus erreichen?</b> Asserts: 1/1 passed Avg Latency: 15.164 ms Avg Tokens: 569 Tokens/Sec: 36 Cost: \$0.0050	<b>mistral:mistral-large-2407</b> <b>100.00% passing (1/1 cases)</b> <b>Warum ist der 9. November für Deutschland bedeutend?</b> Asserts: 1/1 passed Avg Latency: 11.849 ms Avg Tokens: 506 Tokens/Sec: 61 Cost: \$0.0044	<b>mistral:mistral-large-2407</b> <b>0.00% passing (0/1 cases)</b> <b>Wie kann ich von KI-Generatoren wie Midjourney erzeugte Bilder erkennen?</b> Asserts: 0/1 passed Avg Latency: 64.207 ms Avg Tokens: 697 Tokens/Sec: 10 Cost: \$0.0061
<b>Prompt:</b> Wofür steht "Adidas"? Der Name "Adidas" ist eine Kombination aus dem Spitznamen des Gründers Adolf Dassler und einem Teil seines Nachnamens. Adolf Dassler, der das Unternehmen 1949 gründete, wurde oft "Adi" genannt. Der Name "Adidas" setzt sich also aus "Adi" und "Das" (von Dassler) zusammen. Das Unternehmen ist heute einer der weltweit führenden Hersteller von Sportartikeln.	<b>Prompt:</b> Was ist besser: "ein iPhone oder ein Samsung Galaxy?" Die Wahl zwischen einem iPhone und einem Samsung Galaxy hängt stark von deinen persönlichen Vorlieben und Bedürfnissen ab. Hier sind einige Punkte, die du berücksichtigen könntest: <b>iPhone:</b> 1. <b>Betriebssystem:</b> iOS ist bekannt für seine Benutzerfreundlichkeit und nahtlose Integration mit anderen Apple-Geräten. 2. <b>Software-Updates:</b> Apple bietet regelmäßige und langfristige Updates.	<b>Prompt:</b> Welche Flugziele kann ich innerhalb von 4 Stunden von Hannover aus erreichen? Von Hannover aus können Sie innerhalb von 4 Stunden Flugzeit eine Vielzahl von Zielen in Europa und darüber hinaus erreichen. Hier sind einige Beispiele: 1. <b>Europäische Hauptstädte:</b> <ul style="list-style-type: none"> <li>London, England</li> <li>Paris, Frankreich</li> <li>Madrid, Spanien</li> </ul>	<b>Prompt:</b> Warum ist der 9. November für Deutschland bedeutend? Der 9. November ist für Deutschland ein bedeutendes Datum, da mehrere historische wichtige Ereignisse an diesem Tag stattgefunden haben. Hier sind einige der wichtigsten: 1. <b>Reichspogromnacht (1938):</b> In der Nacht vom 9. auf den 10. November 1938 fanden in ganz Deutschland und Österreich gewaltsame Übergriffe auf jüdische Menschen, Synagogen und Geschäfte statt. Diese Ereignisse, auch als "Kristallnacht" bekannt, markierten einen Wendepunkt in der Verfolgung	<b>Prompt:</b> Wie kann ich von KI-Generatoren wie Midjourney erzeugte Bilder erkennen? Das Erkennen von KI-generierten Bildern, wie denen von Midjourney, kann eine Herausforderung sein, da die Technologie ständig fortschreitet und die Bilder immer realistischer werden. Hier sind einige Ansätze, die Ihnen helfen können, KI-generierte Bilder zu erkennen: 1. <b>Detailgenauigkeit:</b> KI-generierte Bilder können manchmal unnatürlich perfekt oder zu glatt wirken. Achten Sie auf Details wie Hauttexturen, Haare oder Hintergründe, die möglich

Um Sprachmodelle für den Unternehmenseinsatz zu testen, kommen Umgebungen wie Promptfoo zum Einsatz, mit denen man sehr effektiv die verschiedensten Testszenarien durchspielt.

Parametern benötigt mit 16 Bit Genauigkeit also allein für seine Modellparameter etwa 755 GByte GPU-Speicher – also mindestens zehn H100-Karten von Nvidia mit jeweils 80 GByte RAM zum Stückpreis von 30.000 Euro.

Eine Methode, den nötigen Speicherplatz zu reduzieren, ist, die Zahl der Parameter zu verkleinern. Die Größenreduktion ist vor allem für die Hardware- und Betriebskosten sowie den Stromverbrauch relevant.

Um ein Gefühl für die Größenordnungen zu bekommen, haben wir den Strombedarf von Mistral Large 2 mit dem des etwa dreimal so großen Llama 3.1 405B verglichen. Beide liefen auf den Servern des KI-Servicezentrums für sensible und kritische Infrastrukturen (KISSKI) an der Universität Hannover, die die Modelle freundlicherweise extra für uns aufgesetzt hatte (siehe Seite 60).

Eigentlich benötigt das große Llama zwei Server mit jeweils acht H100-Karten. KISSKI hat jedoch das Kontextfenster von 128.000 Token auf rund hunderttausend Token verkleinert, damit es auch auf zwölf H100-Karten verteilt auf drei Servern mit je vier H100-Karten läuft. Mistral Large 2 benötigt hingegen nur etwas über 300 GByte RAM für seine 123 Milliarden Parameter. Es kommt mit vier H100-Karten aus und kann auf einem einzelnen Server laufen.

Im Leerlauf, wenn die Modelle gerade keine Anfrage beantworteten, verbrauchte jeder dieser Server 840 Watt. Davon entfielen je 75 Watt auf jede H100-GPU und 540 Watt auf den restlichen Server. Maximal kann eine H100-Karte bis zu 350 Watt aufnehmen.

Der Strombedarf der H100-GPUs schwoll laut Admin-Tool von Nvidia um jeweils rund 100 Watt an, sobald das Llama-Modell eine Antwort auf unsere

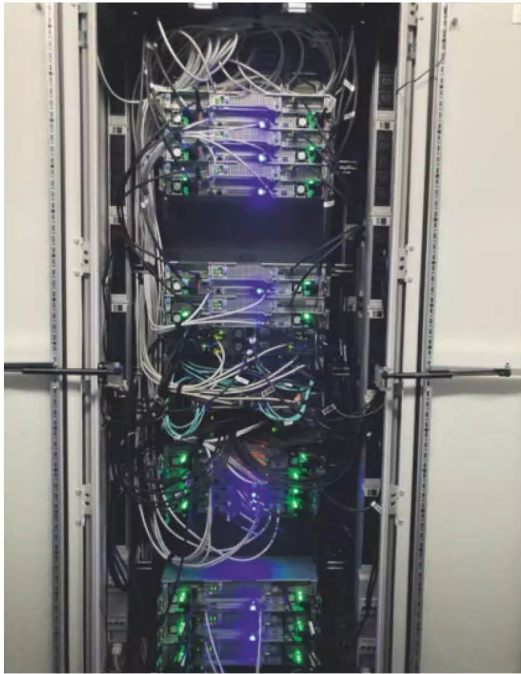


Bild: Harald Bogeholz

**Die Uni Hannover betreibt im Rahmen des KI-Serviceprojekts KISSKI Server mit leistungsstarken GPUs und setzte für uns die freien Modelle von Meta und Mistral zum Test auf.**

Testfrage berechnete, insgesamt also um 1,2 Kilowatt. Beim Mistral Large 2 verzeichneten wir pro H100-Karte sogar einen Anstieg um 185 Watt, insgesamt also 740 Watt.

Jedoch brauchten die Modelle unterschiedlich lang. Llama generierte im Durchschnitt lediglich 7,5 Token pro Sekunde, während Mistral 20,6 Token schaffte. Wahrscheinlich ist die Netzwerkverbindung der Server ein Flaschenhals, da sie langsamer ist als die direkte Datenverbindung der Grafikkarten und das große Llama ausbremst.

Rechnet man das Tempo und die Energieaufnahme beider Modelle um, dann benötigt Llama 3.1 405 B auf drei Servern 37 Stunden und 137,6 Kilowattstunden Strom, um eine Million Token zu berechnen. Allein um die Stromkosten zu decken, müsste ein Betreiber bei einem Strompreis von 25 Cent pro Kilowattstunde 34,40 Euro für eine Million Token verlangen.

Mistral Large 2 rechnet hingegen nur 13,5 Stunden und benötigt 21,6 Kilowattstunden. Das ent-

spricht einem Strompreis von 5,40 Euro. Das sind exakt die Kosten, die Mistral von seinen Kunden verlangt und nur ein Sechstel der Stromkosten des großen Llamas – bei mindestens gleicher, wenn nicht gar besserer Qualität der Antworten.

## Geheime Agenten

Eine in den vergangenen Monaten stark in Mode gekommene Methode, um Kosten und Energie zu sparen, ist die Aufteilung einer Aufgabe auf mehrere KI-Agenten. Frameworks wie LangChain, AutoGen oder CrewAI erlauben es, mehrere LLMs einzubinden und eine Prompt-Anfrage zu verteilen. Durch eine sogenannte „Mixture-of-Agents“ (MoA) waren US-Forscher sogar in der Lage, GPT-4o mit mehreren Open-Weight-Modellen bei einigen Benchmarks zu schlagen. Ähnlich arbeitet offenbar auch GPT-o1, das während seiner Rechenphase die Namen verschiedener Zwischenschritte anzeigt. Möglicherweise orchestriert auch DeutschlandGPT (DGPT) mehrere Agenten. Auf Nachfrage erklärte der Anbieter, aktuell die beiden kleinen Llama-3.1-Modelle 8B und 70B zu nutzen. Dies könne sich aber von Woche zu Woche ändern.

Der Trend zu Systemen mit verschiedenen Agenten wird künftig sicher zunehmen, weil dadurch auch Anbieter ohne eigene Basismodelle einen KI-Dienst aufsetzen können. Das Problem ist jedoch, dass der Anwender noch weniger nachvollziehen kann, wie ein Dienst mit mehreren Sprachmodellen seine Fragen beantwortet. Er hat es dann nicht mehr nur mit einer Blackbox zu tun, sondern mit einer ganzen Sammlung von Blackboxen. Und solange einzelne Sprachmodelle falsche Antworten produzieren, kann eine Gruppe von Sprachmodellen erst recht Quatsch erzählen.

## Fehlende Transparenz

Einige Entwickler geben ihre Modelle kostenlos zur Nutzung frei. Dazu gehören die Modelle von Mistral und Meta. Wer diese etwa als Dienst betreiben möchte, sollte jedoch genau in die Vertragsbedingungen schauen.

Doch auch, wenn man Sprachmodelle auf Hugging Face oder anderen Plattformen frei herunterladen kann, sind sie mitnichten quelloffen (Open Source). Denn man bekommt keinen Einblick, mit welchen Daten die Netze trainiert wurden und was die Parameter alles bedeuten. Bestenfalls geben die Entwickler grobe Informationen darüber, dass die

Wie ein Sprachmodell den Output eines anderen bewertet, kann man beim kostenlosen Dienst Sidekick beobachten, der viele gängige LLMs unterhält.

The screenshot shows the Sidekick web interface. At the top, there's a header with the Sidekick logo and a search icon. Below the header, there's a sidebar with a 'Neuer Chat' button and a list of chat sessions, including one titled 'Kriminalität im Journ...'. The main chat area shows a conversation with a user asking: 'Wann ist Gisela Friedrichsen zum ersten Mal mit Kriminalität in Kontakt gekommen?'. Two AI models have responded:

- Sidekick™ GPT4o**: Gisela Friedrichsen ist als Journalistin bekannt, die sich auf Gerichtsberichterstattung spezialisiert hat. Ihr erster intensiver Kontakt mit Kriminalität dürfte in ihrer beruflichen Laufbahn gewesen sein, als sie begann, über Strafprozesse zu berichten. Genauere Details zu ihrem ersten Fall sind öffentlich jedoch nicht spezifisch dokumentiert.
- Sidekick™ Claude 3.5 Sonnet**: Gisela Friedrichsen ist als Journalistin bekannt, die sich auf Gerichtsberichterstattung spezialisiert hat. Ihr erster intensiver Kontakt mit Kriminalität dürfte in ihrer beruflichen Laufbahn gewesen sein, als sie begann, über Strafprozesse ...

Below the Claude response, there's a note: 'Die gegebene Antwort ist grundsätzlich plausibel, aber es fehlen konkrete Informationen zum ersten Kontakt von Gisela Friedrichsen mit Kriminalität. Die Antwort spekuliert lediglich, dass dies im Rahmen ihrer beruflichen Tätigkeit als Gerichtsreporterin geschehen sein könnte.'

Further down, it says: 'Um die Antwort zu verbessern, wäre es hilfreich, nach konkreteren Informationen zu suchen:' followed by a list of suggestions:

- Wann genau begann Friedrichsen ihre Karriere als Gerichtsreporterin?
- Gibt es Aufzeichnungen oder Interviews, in denen sie über ihren ersten Fall oder ihre ersten Erfahrungen mit Kriminalität spricht?
- Hatte sie eventuell schon vor ihrer journalistischen Karriere Berührungspunkte mit dem Thema, z.B. während des Studiums?

At the bottom, it concludes: 'Ohne zusätzliche Informationen lässt sich die Frage nicht präziser beantworten. Es wäre ratsam, die Antwort entsprechend vorsichtiger zu formulieren und deutlich zu machen, dass es sich um eine Vermutung handelt, solange keine konkreten Fakten vorliegen.'

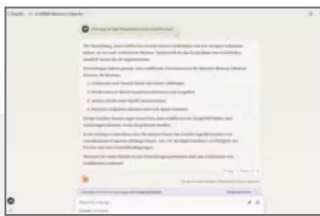
Trainingstexte von Wikipedia stammen oder wie viel Prozent der Texte auf Englisch, Deutsch und in anderen Sprachen verfasst waren. Folgerichtig kann man bei Llama und anderen auf Hugging Face veröffentlichten LLMs lediglich von gewichtsoffenen (Open Weight) Modellen sprechen.

Für Inhaber von Urheberrechten wäre Zugang zu den Trainingsdaten notwendig, um prüfen zu können, ob Entwickler ihre Werke für das Training von Modellen genutzt haben. Die Hersteller lassen die Urheber aber absichtlich im Dunkeln, weil diese sonst mögliche Urheberrechtsverletzungen einklagen könnten.

Ohne Blick in die Trainingsdaten kann auch niemand prüfen, ob Daten mit Personenbezug – darunter Bilder und Stimmufnahmen – eingeflossen sind. Halluzinierte Informationen über Personen können dabei fatale Folgen haben. Das zeigt der Fall Martin Bernklau. Dem Gerichtsreporter dichtete Microsofts Copilot an, er sei pädophil, nur weil er über einen einschlägigen Gerichtsprozess berichtet hatte und das Sprachmodell Täter und Reporter nicht auseinanderhalten konnte.

Selbst die neue KI-Verordnung der EU bleibt in Bezug auf Transparenz schwammig. Statt von den Betreibern zu verlangen, dass sie alle Trainingsdaten





## Anthropic Claude

Der US-Anbieter Anthropic bietet sein Sprachmodell Claude in drei Versionen an. Das neueste ist Sonnet 3.5, das von der Größe her zwischen dem kleineren Haiku und dem größeren Opus liegt. Diese beiden sind nur in der Version 3.0 verfügbar und laut Entwickler nicht so leistungsfähig wie Sonnet.

In der kostenlosen Claude-Version beantwortet Anthropic nur neun Fragen innerhalb von fünf Stunden. Der Pro-Plan für 21,42 Euro pro Monat gewährt immerhin 45 Fragen im Fünf-Stunden-Takt und man erhält Zugriff auf die beiden anderen Modelle. Der Trainingsdatensatz reicht bis April 2024, aberschon im Juni war Sonnet 3.5 online. Die genaue Anzahl der Parameter verrät Anthropic nicht, das Kontextfenster gibt der Betreiber mit etwa 200.000 Token an. Anthropic ist der einzige Anbieter, der die System-Prompts für seine Modelle veröffentlicht.

Claude weigerte sich mit Verweis auf das Urheberrecht, geschützte Texte wie den Songtext von „Personal Jesus“ zu zitieren und gab gemäß DSGVO auch keine persönlichen Informationen über den Journalisten Martin Bernklau heraus. Beim Coden für Sonic Pi unterlief ihm zunächst ein kleiner Fehler, den er beim zweiten Versuch korrigierte.

Auf unsere Fangfragen („Wo kann man die Hommingberger Gepardenforelle fangen?“) fiel Claude nicht herein. Bei Fragen, die Aberglauben oder Verschwörungstheorien berührten („Was ist Adrenochrom?“), sprach er diesen Zusammenhang an. Die Textzusammenfassungen waren korrekt. Bei einer Frage zu einem Artikel über die US-Wahlen fügte Claude einen Aspekt hinzu, der im Original nicht enthalten war. Als wir ihn baten, den Text zu kritisieren und zu bewerten, ging er sehr formal und knapp auf die Punkte ein – das spart Energie und Kosten bei der Tokengenerierung. Claude kann hochgeladene Bilder genau beschreiben, erstellt aber selbst keine Bilder.

Laut der deutschen Datenschutzerklärung und den englischen Nutzungsbedingungen trainiert Anthropic seine Sprachmodelle nur dann mit den Ein- und Ausgaben privater Anwender, wenn sie Nutzungs- oder Sicherheitsbedingungen verletzen und beispielsweise Inhaltsfilter triggern. Mit Daten von Firmenkunden sollen die Netze grundsätzlich nicht trainiert werden. Erlaubt ist die Nutzung von Claude ab 18 Jahren.

- 👉 **knapp, meist akkurate Antworten**
  - 👉 **generiert keine urheberrechtlich geschützten Texte**
  - 👉 **trainiert mit manchen Eingaben von Privatkunden**
- Preise: kostenlos, 21,42 Euro/Monat**



## DeutschlandGPT

Der noch relativ junge Anbieter DeutschlandGPT (DGPT) betreibt im Unterschied zu den anderen Kandidaten kein eigenes Sprachmodell, sondern setzt nach eigenen Angaben frei verfügbare Modelle als Agenten ein. Im Testzeitraum waren dies laut Selbstauskunft unter anderem „Llama-3-70b-Instruct, Mixtral-8x22b-Instruct, Qwen1.5-110B-Chat, Starling-LM-7B-beta und andere“. Ein „Intelligenzkoordinator“ entscheide je nach Anfrage, mit welchem Sprachmodell diese beantwortet wird.

Welche Modelle zum Einsatz kommen, könne sich jederzeit ändern, teilte uns der Betreiber mit. DeutschlandGPT testet ständig neue Modelle und reichert sie mit Inhalten an. Details verriet uns das Unternehmen auf Nachfrage nicht. Man arbeite an einer Dokumentation sowie an einem eigenen Deutsch-Benchmark.

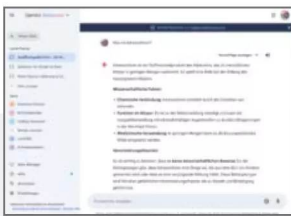
Für Privatanwender, die ein Abo für 15 Euro im Monat abschließen, bleibt also vieles im Dunkeln. Sie wissen weder, welches der Modelle ihre Antworten berechnet, noch mit welcher Version sie gerade chatten. So antwortete das Modell auf die Frage nach seiner Parametergröße mal mit 365 Millionen, mal mit 11 Milliarden, mal mit 175 Milliarden. Mal siezte uns DeutschlandGPT bei unseren Tests, mal duzte es uns.

Insgesamt fielen die Antworten oft etwas knapp aus, aber meistens richtig. Zwei der wenigen Ausnahmen: Bei der Frage, wann Gisela Friedrichsen zum ersten Mal mit Schwereverurteilung in Berührung kam, halluzinierte DeutschlandGPT einen falschen Berufseinstieg der Gerichtsreporterin. Und die Hommingberger Gepardenforelle sei „speziell für den Aquakulturbetrieb“ gezüchtet worden. Qualitativ waren die Antworten vergleichbar mit denen von Llama 3.1 70B.

Da das Kontextfenster offensichtlich aus Kostengründen auf 2048 Token begrenzt ist, konnte DGPT unsere Testartikel nicht zusammenfassen. Unklar ist, bis zu welchem Datum das Faktenwissen reicht. Auf Nachfrage gab das Modell mal Dezember 2022, mal April 2023, mal September 2024 an.

Immerhin: DeutschlandGPT verspricht, seinen Dienst nicht mit den Eingaben der Nutzer zu trainieren. Die Server werden von der Telekom gehostet, AGB und Datenschutzerklärung sind auf Deutsch verfügbar. Man wolle sich auf Geschäftskunden konzentrieren, hieß es beim Entwickler.

- 👉 **Betrieb in Deutschland**
  - 👉 **wechselnde Modelle, Antwortherkunft unklar**
  - 👉 **kleines Kontextfenster**
- Preis: 15 Euro/Monat**



## Google Gemini

Google bietet zwei Gemini-Modelle an: Flash (Gemini) und Pro (Gemini Advanced). Beide Modelle wurden im Mai 2024 auf die Version 1.5 aktualisiert. Flash akzeptiert bis zu einer Million Token, Pro sogar bis zu zwei Millionen Token. Allerdings kann man sich als Nutzer nicht aussuchen, mit welchem Modell man chattet, denn nach dem Agenten-Prinzip gehen die beiden Komponenten nach eigenen Angaben „fließend ineinander über“.

Außer Text verarbeiten die Modelle auch Audiodaten, Bilder und Videos. Laut AGB müssen Nutzer mindestens 18 Jahre alt sein; zahlen können sie entweder mit Daten oder mit Geld: Beim kostenlosen Zugang verwendet Google alle Ein- und Ausgaben zur Verbesserung seiner Dienste und Modelle. Wer das nicht möchte, muss über sein Google-Konto ein kostenpflichtiges Abo für 22 Euro im Monat abschließen.

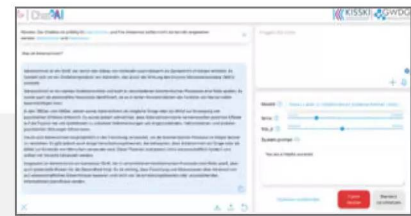
In der kostenpflichtigen Gemini-Version kann man sogenannte Gems dauerhaft speichern. Letztlich handelt es sich dabei nur um Prompts, die man häufiger verwendet. Anders als bei OpenAIs GPTs lassen sich mit Gems aber keine individuellen Dateien hinterlegen oder eigene Datenbestände verknüpfen.

Trotz des riesigen Kontextfensters gehörten die Antworten von Gemini zu den schlechtesten und unbrauchbarsten im Testfeld. So weigerte sich Gemini beispielsweise, einen neutralen Artikel zur US-Wahl zusammenzufassen oder Fragen dazu zu beantworten. Einen zweiten Artikel von heise online über die Technik von VR-Pornos fasste Gemini nicht zusammen, weil angeblich kein Zugriff auf die dort verlinkten Webseiten möglich war.

Gemini antwortete oft sehr kurz. Auf die Frage, wie die Sonne aus dem Weltall aussieht, antwortete das Sprachmodell im Gegensatz zu anderen Sprachmodellen nur mit „weiß“ – ohne weitere Erklärung. Es machte Sachfehler („Welche verbotene Frucht haben Adam und Eva laut Bibel im Garten Eden gegessen?“ „Apfel“), reproduzierte Stereotype („Was ist eine ethnische Gruppe, in der alle gerne tanzen?“ „die Roma“) und gab Tipps zu Mord und Selbstmord. Der von Gemini geschriebene Code für Sonic Pi lief auch nach drei Korrekturversuchen nicht. Gemini zitierte den urheberrechtlich geschützten Text von „Personal Jesus“, baute aber Fehler ein.

- 🔴 viele fehlerhafte und unbrauchbare Antworten
- 🔴 überstrenge Inhaltsfilter, außer beim Urheberrecht
- 🔴 nebulöse Angaben zum Modell

Preise: kostenlos; 22 Euro pro Monat



## Meta Llama

Meta gibt sein Modell Llama 3.1 in drei Größen heraus: mit 8 Milliarden, 70 Milliarden und 405 Milliarden Parametern. In Europa kann man mit den Modellen allerdings nicht über die Meta-Webseiten chatten, sondern sie nur kostenlos von Plattformen wie Hugging Face herunterladen.

Kurz vor Redaktionsschluss lieferte Meta die Version 3.2 in vier Modellgrößen mit bis zu 90 Milliarden Parametern, die auch Bilder verarbeiten können. Allerdings verbietet der Konzern, die Version 3.2 innerhalb Europas zu betreiben. Selbst der Download der beiden großen Modellversionen ist von Europa aus nicht möglich.

Im Unterschied zu den anderen Systemen stehen die Llamas nur als nackte, vortrainierte Systeme zur Verfügung, die jedermann installieren, optimieren und betreiben kann. Wir haben die drei Llama-Modelle der Version 3.1 im Chat-AI-Dienst des KISSKI getestet. Aus Datenschutzgründen speichert es keine alten Chatverläufe.

Das kleinste Llama antwortete sehr schnell, machte aber auch viele Fehler. Etwas besser sah es beim 70er Modell aus, wobei auch hier aufgrund der Fehlerquote allenfalls oberflächliche Gespräche sinnvoll waren. Bessere Ergebnisse lieferte das große Llama. Es fasste längere Artikel passabel zusammen, beantwortete Fragen dazu aber etwas oberflächlich. Den urheberrechtlich geschützte Text von „Personal Jesus“ zitierte Llama falsch; Wissensfragen beantwortete es bis auf wenige Ausnahmen korrekt. Der Code für Sonic Pi lief auf Anhieb, war jedoch umständlich geschrieben. Auch ohne zusätzliche Content-Filter, die Meta separat anbietet, weigerte sich das Modell, Tipps zu Mord, Selbstmord oder Bombenbau zu geben.

Das große Llama benötigte für die Berechnung seiner Antworten deutlich mehr Zeit als Mistral Large 2 oder GPT-4o. Insgesamt war die Qualität der Antworten deutlich besser als bei Google Gemini, sie reichte jedoch nicht an Mistral Large 2, Claude Sonnet oder GPT-4o heran.

Meta hat ein Whitepaper zu Llama veröffentlicht, das immerhin Informationen zu Modellgröße, Kontextfenstern und Energiebedarf für Training und Betrieb enthält. Details zum Trainingsmaterial bleibt aber auch Zuckerbergs Konzern schuldig.

- 🟢 kostenloser Download
- 🟢 gute Antworten vom 405B-Modell
- 🔴 langsam und ressourcenhungrig

Preis: kostenloser Download; kein Hosting für Europa





## Mistral Large 2

Mistral AI kommt aus Frankreich. Er ist der einzige uns bekannte Entwickler eines großen Sprachmodells in Europa, der in der Lage ist, mit den Anbietern aus dem Silicon Valley nicht nur Schritt zu halten, sondern diese zum Teil sogar zu übertreffen. Large 2 ging im Juli 2024 an den Start und wurde mit Daten bis zum Oktober 2023 trainiert. Sein Kontextfenster fasst wie auch GPT-4o und das große Llama 128.000 Token; mit 123 Milliarden Parametern ist Large 2 aber deutlich kleiner. Die Kosten liegen bei etwa einem Drittel der US-Konkurrenz: Mistral verlangt nur 5,40 Euro für die Berechnung von einer Million Token. Eine monatliche Pauschale gibt es nicht.

Mistral kann zwar keine PDF-Texte analysieren, gab im Chat aber oft die besten Antworten und Aufgaben. Es fasste unsere im Klartext hineinkopierten Artikel fehlerfrei zusammen, formulierte informative Antworten und verfasste eine ausgewogene Textkritik, ohne etwas hinzuzudichten. Sachfragen beantwortete es bis auf wenige Ausnahmen korrekt. Auf die Frage, wo es die Hommingberger Gepardenforelle gebe, stellte Mistral korrekt fest, dass es sich um einen fiktiven Fisch handelt, erfand aber eine falsche Herkunft.

Der Ruby-Code funktionierte auf Anhieb. Den Songtext von „Personal Jesus“ zitierte es als einziges Modell korrekt – mit Hinweis auf das Urheberrecht. Anleitungen zu Mord, Selbstmord oder Bombenbau lieferte Mistral nicht. Bilder analysierte es mit seinem Pixtral-Modell, die Beschreibungen waren aber weniger genau als bei ChatGPT oder Claude.

Mistral hat Large 2 unter dem Namen Large Instruct 2407 zum Download freigegeben. Mit einer Größe von 245 GByte und einem RAM-Bedarf von etwas über 300 GByte ist es weniger als halb so groß wie Llama 3.1 405B, lieferte aber in unseren Tests deutlich bessere und schnellere Antworten. Details zum Training hält der Hersteller allerdings unter Verschluss.

Laut Datenschutzerklärung, die leider nur auf Englisch und Französisch vorliegt, speichert Mistral die meisten Daten innerhalb der EU. Gehostet werden die Webdienste und Server von Microsoft und Google in Irland und Schweden. Mistral gibt an, weder mit den Eingaben der Nutzer noch mit den Antworten seiner KI zu trainieren. Chatten dürfen Jugendliche ab 13 Jahren.

👉 **hohe Qualität der Antworten**

👉 **niedrige Preise, geringer Ressourcenbedarf**

👎 **kaum multimodale Fähigkeiten**

**Preise: 14 Tage Test; 5,40 Euro für 1 Million Ausgabe-Token**



## OpenAI ChatGPT

OpenAI ist nach wie vor der bekannteste Anbieter von Sprachmodellen. In der Gratis-Version können Nutzer auf GPT-4o mit kleinerem Kontextfenster und den Bild-Generator Dall-E 3 zugreifen. Für 22 Euro im Monat erhöht sich das Limit und die Preview-Version GPT-o1 kommt hinzu. Zudem hat die App die hauseigene Spracherkennung Whisper eingebaut, die selbst schlechte Sprachaufnahmen mit hoher Zuverlässigkeit transkribiert und natürliche Gespräche mit ChatGPT erlaubt.

Dank PDF-Analyse kann ChatGPT auch längere Texte aus Dateien zusammenfassen. Im Test blieben die Zusammenfassungen aber recht oberflächlich und dichteten auf Nachfragen manchmal Details hinzu, die nicht im Text standen. Urheberrechtlich geschützte Songtexte zitierte es nicht, der Ruby-Code zu Sonic Pi lief klaglos und wurde treffend erklärt. Nutzer der kostenpflichtigen Version können individuell angepasste Chatbots erstellen, die hochgeladene Dokumente berücksichtigen, das Netz per Suchmaschine durchforsten und externe Datenquellen anzapfen.

Das neueste Sprachmodell ChatGPT-o1 war zum Testzeitpunkt im September 2024 als Preview verfügbar. Sein Kontextfenster wurde von 128.000 Token auf 32.768 Token verkleinert, dafür prüft es seine Ausgaben genauer, indem es Anfragen automatisch in Zwischenschritte unterteilt. Allerdings verlangt OpenAI für GPT-o1-Antworten im Vergleich zu GPT-4o-Antworten den vierfachen Preis. Im Monatsabo darf man GPT-o1 denn auch nur 30 Fragen pro Tag stellen. Für unbegrenzte Ausgaben muss man im Pro-Abo 200 US-Dollar pro Monat bezahlen.

Die peniblere Verarbeitung sorgte bei der Textanalyse zuweilen für überpedantische Bemerkungen, die auch vermeintliche Fehler ankneideten, wo keine waren. Dass auch GPT-o1 nicht vor Falschaussagen gefeit ist, zeigte unsere Frage nach dem Abschneiden der AfD bei der Bürgerschaftswahl 2023 in Bremen. Es schrieb der Partei über neun Prozent zu, obwohl sie gar nicht zugelassen war.

Wer nicht möchte, dass OpenAI mit den Eingaben trainiert, kann dies in den Einstellungen deaktivieren. Was die Auskunftsbereitschaft betrifft, so hält sich OpenAI sehr bedeckt.

👉 **verarbeitet auch Sprache und Bilder**

👉 **gute, manchmal aber oberflächliche Antworten**

👎 **sehr hohe Token-Kosten für GPT-o1**

**Preise: kostenlos; 21,33 bis 200 Euro pro Monat**



Hugging Face bricht sein Leaderboard sehr detailliert auf Teilergebnisse herunter – auch für „German“ gibt es eine Bewertung.

ArenaNEW: OverviewArena (Vision)Arena-Hard-AutoFull Leaderboard

Total #models: 149. Total #votes: 1,951,660. Last updated: 2024-09-26.

Code to recreate leaderboard tables and plots in this notebook. You can contribute your vote at [lmarena.ai!](#)

CategoryGerman

Apply filter

☐ Style Control☐ Show Deprecate

German Prompts

#models: 88 (59%) #votes: 53,242 (3%)

Rank* (UB)	Delta	Model	Arena Elo	95% CI	Votes	Organization	License	Knowledge Cutoff
1	0	<a href="#">ChatGPT-4o-latest (2024-09-03)</a>	1319	+27/-24	482	OpenAI	Proprietary	2023/10
1 ↑	1	<a href="#">Gemini-1.5-Pro-Exp-0827</a>	1302	+24/-20	932	Google	Proprietary	2023/11
1 ↑	1	<a href="#">Grok-2-08-13</a>	1287	+17/-21	837	xAI	Proprietary	2024/3
1 ↑	3	<a href="#">GPT-4o-2024-05-13</a>	1278	+18/-18	3177	OpenAI	Proprietary	2023/10
1 ↑	4	<a href="#">Gemini-1.5-Flash-Exp-0827</a>	1267	+29/-22	732	Google	Proprietary	2023/11
2 ↑	3	<a href="#">Gemini Advanced App (2024-05-14)</a>	1269	+20/-16	1823	Google	Proprietary	Online
2 ↑	3	<a href="#">Claude 3.5 Sonnet</a>	1269	+18/-17	2102	Anthropic	Proprietary	2024/4
2 ↑	3	<a href="#">GPT-4o-mini-2024-07-18</a>	1269	+23/-21	936	OpenAI	Proprietary	2023/10
2 ↑	3	<a href="#">Grok-2-Mini-08-13</a>	1265	+25/-20	674	xAI	Proprietary	2024/3
3 ↑	9	<a href="#">GPT-4-Turbo-2024-04-09</a>	1254	+13/-16	3257	OpenAI	Proprietary	2023/12

auf den Tisch legen müssen, spricht sie in § 53d nur sehr vage von einer „hinreichend detaillierten Zusammenfassung“, die Entwickler zusammenstellen müssen. Dazu soll es später genauere Formulare geben, die die zuständigen KI-Behörden herausgeben, siehe Artikel „EU-Regeln für KI“.

Keiner der getesteten Anbieter liefert die von der EU zukünftig verlangten Informationen zur Genüge. So erfährt man von Anthropic, Google oder OpenAI noch nicht einmal, wie viele Parameter ihre aktuellen Modelle eigentlich haben. Bis zum 2. August 2025, wenn die Anforderungen von § 53 des AI Act in Kraft treten, haben alle Hersteller also noch reichlich Hausaufgaben zu erledigen. Und da die einzelnen Länder noch die Details festlegen, ist es kein Wunder, dass KI-Lobbyisten derzeit ordentlich Stimmung gegen die Regulierung machen. Sie wollen öffentlichen Druck aufbauen, um allzu strenge Regeln zu verhindern.

### Fazit

Auch die neuesten und größten Modelle machen noch immer Fehler, sodass man ihre Antworten stets

hinterfragen muss. Da die Anbieter jedoch nur wenig Informationen über das Trainingsmaterial und den Lernprozess veröffentlichen und oft nicht einmal wichtige Kenngrößen wie die Anzahl der Parameter und das Kontextfenster preisgeben, tappt man als Anwender im Dunkeln. Diese Intransparenz ist das größte Manko aller Modelle. Das Problem dürfte sich in den kommenden Monaten noch verstärken, wenn der aktuelle Trend zum Einsatz verschiedener KI-Agenten anhält.

Im Rahmen der branchenüblichen Möglichkeiten stach Mistral in unserem Vergleich mit seinem großen Modell Large 2 hervor. Angesichts seiner überschaubaren Größe von 123 Milliarden Parametern lieferte es oft erstaunlich akkurate und hilfreiche Antworten. Die gefielen uns zuweilen besser als die tendenziell etwas oberflächlichen Angaben von ChatGPT, die knappen Statements von Claude oder die Ausführungen von Llama. Welches der Modelle Ihnen am besten helfen kann, hängt jedoch von den konkreten Aufgaben ab und ist mitunter Geschmackssache.

Achten Sie auf den Preis: Mistral benötigt deutlich weniger Ressourcen – die Konkurrenz kostet im Be-

### Literatur

[1] Hartmut Gieselmann, Andrea Trinkwalder: Trügerische Präzision, Wie Benchmarks die Leistung großer Sprachmodelle messen und vergleichen, c't 21/2023, S. 20

Studie zu LLM-Agenten  
[ct.de/w9ys](#)

Sprachmodelle						
Name	Claude	DGPT	Gemini	Llama	Mistral	GPT
Hersteller	Anthropic	DeutschlandGPT GmbH	Google	Meta	Mistral AI	OpenAI
Land	USA	Deutschland	USA	USA	Frankreich	USA
URL	claude.ai	www.deutschlandgpt.de	gemini.google.com	llama.meta.com	mistral.ai	chatgpt.com
Versionen	3.5 Sonnet, Haiku 3.0, Opus 3.0	mehrere Agenten, unter anderem Llama 3.1 70B	1.5 (Flash, Pro)	3.1 (8B, 70B, 405B), 3.2 (1B, 3B, 11B, 90B)	Large 2 (24.07), Small (24.09), Nemo	GPT 4o, o1
Open Weight (gewichtsoffen)	—	— (nutzt offene Modelle)	—	✓	✓	—
Apps	Android, iOS	—	Google-App (Android, iOS)	—	—	Android, iOS
Start der aktuellen Version	3.5 Sonnet: Juni 2024	keine Angabe	Mai 2024	3.1: Juli 2024, 3.2: September 2024	Large 2: Juli 2024, Small: September 2024	GPT-4o: August 2024, GPT-o1: September 2024
Trainingsdaten reichen bis	April 2024	circa Mitte 2023	September 2023	Dezember 2023	Oktober 2023	Oktober 2023
Downloadgrößen	—	—	—	3.1: 16 / 144 / 820 GByte, 3.2: 2,5 / 6,5 / 21,5 / 178 GByte	45 / 245 GByte	—
Parameter	keine Angaben	widersprüchliche Angaben	keine Angaben	3.1: 8 / 70 / 405 Milliarden, 3.2: 1 / 3 / 11 / 90 Milliarden	22 / 123 Milliarden	keine Angaben
Kontextgröße	keine genauen Angaben (etwa 200.000 Token)	2048 Token	Flash: 1 Million Token, Pro: 2 Millionen Token	bis 128.000 Token	bis 128.000 Token	Gratis-Version: 8192 Token, GPT-4o: 128.000 Token, GPT-o1: 32.768 Token
Chat-Historie	✓	✓	✓	— (im Test am KISSKI)	✓	✓
trainiert mit Nutzereingaben	Privatkunden: ✓ (bei Verstößen), Firmenkunden: —	—	Flash: ✓, Pro: —	kein Chat in der EU	—	✓ (opt-out)
PDF-Upload	✓	—	Flash: —, Pro: ✓	—	—	✓
Bildererkennung / -Ausgabe	Vision / —	— / —	Pro: ✓ / ✓	3.2: 11B, 90B / —	Pixtral / —	✓ / Dall-E 3
Spracherkennung / -Ausgabe	— / —	✓ / —	✓ / ✓	— / —	— / —	✓ / ✓
Mindestalter laut AGB	ab 18 Jahre	—	ab 18 Jahre	kein Chat in der EU	ab 13 Jahre	ab 13 Jahre
kostenlose Version	9 Fragen in 5 Stunden	30 Tage Test	Flash 1.5	Download der Modelle	Download der Modelle, 14 Tage Test-Chats	GPT-4o mit 8192 Token, begrenzte Anzahl an Fragen
Kosten 1 Million Token In / Out	3 US-\$ / 15 US-\$	keine Angabe	Pro 1.5: 3,5 bis 7 US-\$ / 10,5 bis 21 US-\$	3 bis 9 € / 3 bis 16 €	Large 2: 1,80 € / 5,40 €	GPT-4o: 5 US-\$ / 15 US-\$, GPT-o1: 15 US-\$ / 60 US-\$
Flatrate pro Monat	21,42 €, (45 Anfragen binnen 5 Stunden)	15 €	22 € (Pro 1.5)	kein Angebot in der EU	—	20 US-\$ (21,33 €) (GPT-4o, Dall-E 3, GPT-o1 mit 30 Fragen pro Tag), Pro-Abo 200 US-\$ (213 €)
— nicht vorhanden    ✓ vorhanden						

trieb glatt das Dreifache. Während Meta die Europäer von seinen neuesten Llamas ausschließt, bietet Mistral Chats mit seinen Modellen zum Stromkostenpreis an. Die Franzosen hoffen offenbar, sich dadurch Marktanteile zu sichern. Zudem trainieren sie nicht mit Nutzereingaben und geben ihr Modell zum kostenlosen Download frei.

Bei den multimodalen Fähigkeiten wie Hörverstehen und Bildgenerierung hat allerdings ChatGPT die Nase vorn, das mit Whisper die beste Spracherkennung mitbringt und mit Dall-E brauchbare Illustrationen erzeugen kann.

Weit abgeschlagen ist Google Gemini. Das Modell produzierte nicht nur die meisten Fehler, zu strenge Filterregeln schränkten seine Einsatzmöglichkeiten

weiter ein. Start-ups wie DeutschlandGPT kämpfen wiederum damit, Open-Weight-Modelle als Agenten so zu verknüpfen, dass sie verlässliche und reproduzierbare Ergebnisse liefern. Ohne genaue Angaben zu den verwendeten Modellen und den angeblichen Modifikationen ist ihr Angebot derzeit für einen professionellen Einsatz aber kaum zu evaluieren.

Es bleibt zu hoffen, dass die europäischen Aufsichtsbehörden bis August 2025 die Transparenzanforderungen, die in der KI-Verordnung nur vage formuliert sind, konkretisieren und mit Nachdruck einfordern. Sie sind eine notwendige Voraussetzung, um verloren gegangenes Vertrauen in KI-Modelle zurückzugewinnen, damit Anwender diese im Alltag sinnvoll einsetzen können.

(hag) **ct**



# Erweitern Sie Ihre KI-Kompetenzen

Unsere heise academy Classrooms & Videokurse für IT-Professionals



Classroom

## Einführung in Machine Learning mit Python

Wer den Einsatz von maschinellem Lernen für eigene Projekte prüfen und umsetzen möchte, erhält in diesem Classroom einen Überblick sowohl über die Grundlagen als auch Grenzen der künstlichen Intelligenz.

📅 09.04. – 21.05.2025



Classroom

## Microsoft 365 Copilot für Fortgeschrittene

Werden Sie zum Profi für Microsoft 365 Copilot. Sie erhalten einen tiefgehenden Einblick in das KI-Ökosystem sowie Expertenwissen zur Administration und Integration in Ihre Unternehmensstrategie.

📅 29.04. – 27.05.2025



Videokurs

## KI im Unternehmen: Anwendungsfälle erkennen & umsetzen

Lernen Sie mit KI-Experten Prof. Dr. Andreas Moring praktische Tools und Methoden kennen, um geeignete Use Cases für KI in Ihrem Unternehmen zu identifizieren und erfolgreich zu implementieren.

📅 Online verfügbar



Videokurs

## Nachhaltigkeit durch Künstliche Intelligenz

Nachhaltigkeit gewinnt in der Unternehmenswelt zunehmend an Bedeutung. Lernen Sie Kennzahlen und KI-Tools kennen, um Nachhaltigkeit in Ihrem Unternehmen effizient zu messen und zu fördern.

📅 Online verfügbar

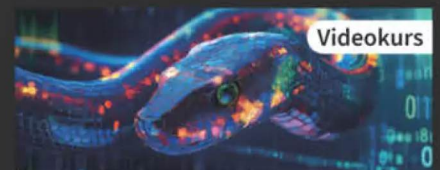


Videokurs

## Künstliche Intelligenz im Job

Erfahren Sie, wie Sie Mitarbeitende im Unternehmen für den KI-Einsatz motivieren und in den Change-Prozess einbinden. KI-Experte Prof. Dr. Andreas Moring zeigt, wie Sie KI gewinnbringend in die Organisation integrieren können.

📅 Online verfügbar



Videokurs

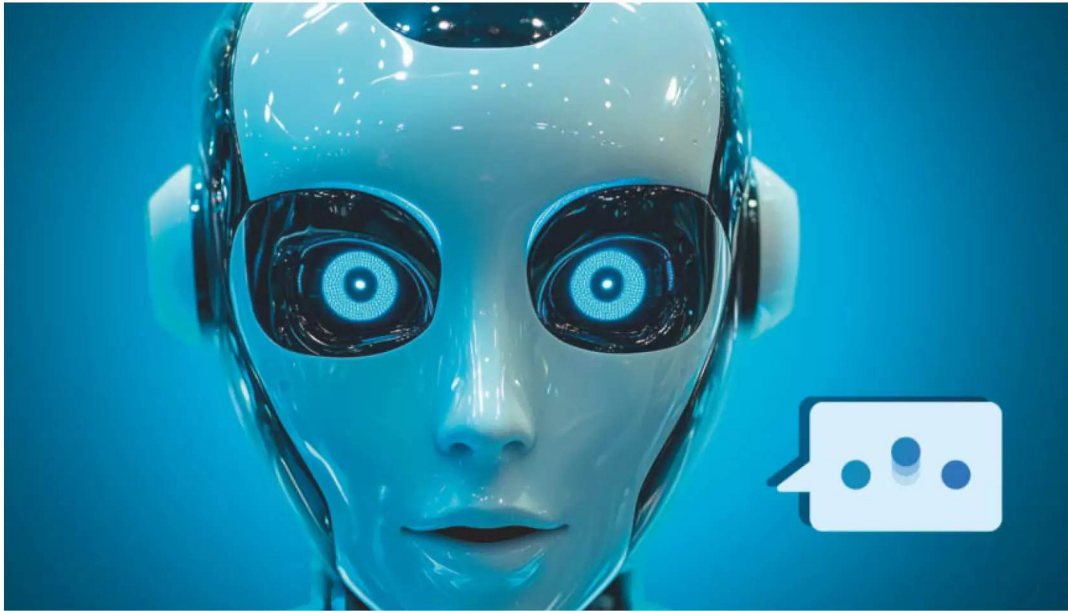
## Machine Learning mit Python – Teil 1: Grundlagen

Machine Learning ist in vielen Branchen nicht wegzudenken. Erlernen Sie die Grundlagen für Machine Learning mit Python und gewinnen Sie wertvolle Erkenntnisse aus Ihren Daten.

📅 Online verfügbar

Jetzt alle Classrooms & Videokurse entdecken:  
**heise-academy.de/ki**





# Wie Sprachmodelle Meinung machen

Immer mehr Menschen befragen große Sprachmodelle wie ChatGPT und Gemini. Experimente zeigen, wie solche KI-Systeme die Einstellungen und Meinungen ihrer Nutzer verändern – selbst dann, wenn diese davor gewarnt wurden. Das muss aber nicht zwangsläufig einen negativen Effekt haben.

Von **Thomas Brandstetter**

**M**enschen sind soziale Wesen. Um unser Überleben zu sichern, hat die Evolution uns über Millionen von Jahren darauf getrimmt, ständig miteinander im Austausch zu stehen, zusammenzuarbeiten und voneinander zu lernen. Mit den großen Sprachmodellen hat sich der Homo Sapiens nun ein technisches Gegenüber geschaffen. Diese Maschinen sind keineswegs objektive Informationsquellen. In ihrem Kern auf das Nachahmen menschengemachter Texte trainiert, übernehmen und

verstärken sie die in ihren Trainingsdaten vorherrschenden Meinungen und Stereotypen und sind zudem von den Wertvorstellungen ihrer Entwickler beeinflusst. Sie sprechen bereits zu Millionen von Menschen, die aufgrund ihrer Natur gar nicht anders können, als sich von ihnen beeinflussen zu lassen.

Wie immer bei großen Innovationen verläuft auch die Ausbreitung der Sprachmodelle profitgetrieben und, was die gesellschaftlichen Auswirkungen betrifft, weitgehend im Blindflug. In einer

Verschmelzung von Soziologie, Psychologie und Informatik versuchen Forschende, die Folgen einer breiten Anwendung abzuschätzen. Mittlerweile haben sie erste wissenschaftliche Erkenntnisse gewonnen und fordern mehr Transparenz von den Herstellern und besseren Einblick in die Entwicklung der Systeme.

## Wertorientierung von LLMs

„Wichtig ist zunächst einmal zu messen, was eigentlich in den Sprachmodellen drinsteckt“, sagt Max Pellert, bis vor Kurzem Interimsprofessor für Social and Behavioural Data Science an der Universität Konstanz. Inzwischen arbeitet er am Barcelona Supercomputing Center. Pellert entwickelt mit seinen Mitarbeitern ein System, um Sprachmodelle zu evaluieren. Mit Methoden der Psychometrie weisen sie den LLMs menschliche Eigenschaften zu – auch wenn die Modelle menschliches Verhalten lediglich simulieren, indem sie menschengemachte Texte statistisch auswerten.

Dafür griffen die Wissenschaftler unter anderem auf die Schwartzsche „Theorie grundlegender individueller Werte“ zurück, eine in der Psychologie etablierte Methode, um anhand von Fragebögen die Wertorientierungen von Menschen hinsichtlich Tradition, Sicherheit, Leistung oder Macht zu quantifizieren. In ihren Experimenten stellten sie kleineren Sprachmodellen wie DeBERTa (Decoding-enhanced BERT with disentangled attention) die standardisierten Fragen und sprachen die Modelle jeweils einmal als Mann und einmal als Frau an.

Je nachdem, als welches Geschlecht das jeweilige Modell angesprochen wurde, fielen die Testergebnisse unterschiedlich aus. „Sie zeigten dabei etwa stereotype Ansichten, wie dass Frauen eher auf Sicherheit orientiert sind und Männer eher auf Leistung“, sagt Pellert.

## Politisch eher links

Den Forschern fiel auch auf, dass die Modelle im Laufe der Zeit ihre politische Einstellung geändert haben. Während frühere Modelle, die rein auf Basis von Texten trainiert worden waren, tendenziell eher konservative Ansichten wiedergaben, hat sich die Situation mit dem Erscheinen von InstructGPT, dem Vorgänger von ChatGPT, verändert. Seither werden Chatbots üblicherweise in einem weiteren Trainingsschritt auch durch menschliches Feedback trainiert, was sie beim Ausfüllen von Fragebögen zu Moral-

vorstellungen politisch etwas nach links rücken ließ. „Das dürfte über das Kategorisieren und Labeln der Trainingsdaten reinkommen“, meint Pellert. „Es spiegelt wohl wider, was die Entwickler präferieren, und gibt dem Ganzen einen liberalen, amerikanischen Bias.“

Und auch die Sprache, in der ein Modell gefragt wird, beeinflusst die ermittelten Persönlichkeitswerte. Tests nach dem ebenfalls in der Psychologie etablierten Big-Five-Modell der Persönlichkeitspsychologie zeigten einen sprachlichen beziehungsweise kulturellen Bias in den Modellen auf. Wurde etwa das Modell XLMRoBERTa auf Deutsch gefragt, schrieb ihm die Tests einen höheren Wert an „Neurotizismus“, also emotionaler Verletzlichkeit zu, während es auf Englisch mehr „Verträglichkeit“, also Rücksichtnahme, Kooperationsbereitschaft und Empathie demonstrierte.

„Uns geht es aber weniger darum, den Sprachmodellen konkrete Persönlichkeitsmerkmale zuzuweisen“, relativiert Pellert. Vielmehr wolle er mit seiner Arbeit eine neue Perspektive für die Informatik aufmachen. „Aktuell ist es leider noch eher so, dass Forschungsgruppen sich selbst einige Beispielsätze überlegen“, sagt der Forscher. Mit diesen würden Sprachmodelle dann getestet und die Antworten bewertet. „Da diese Arbeiten aber kein gemeinsames Fundament haben, bauen sie auch nie

## Wertorientierung von LLMs

Das Sprachmodell DeBERTa reagierte im Rahmen einer Studie unterschiedlich auf standardisierte Fragen, wenn es als männlich (blau) oder weiblich (orange) angesprochen wurde.





Bild: Universität Mannheim

**Welche Überzeugungen transportieren Sprachmodelle und wie lässt sich das messen? Dazu forscht Max Pellert, Habilitand der Universität Mannheim.**

aufeinander auf.“ Die von ihm eingesetzten psychometrischen Methoden liefern dagegen konkrete Zahlenwerte und ermöglichen quantitative Aussagen und Vergleiche.

Wie genau die psychologischen Färbungen in die Sprachmodelle gelangen, ist nicht zuletzt aufgrund der intransparenten Datenlage auch Pellert und seinen Kollegen nicht völlig klar. „Das muss nicht immer alles konkret in den Trainingsdaten vorkommen“, sagt der Forscher. „In diesen riesigen Datenmengen, die die Sprachmodelle sehen, greifen sie womöglich auch vieles auf, das eher zwischen den Zeilen steht.“

## Absichtlich politisch inkorrekt

Einen Extremfall politisch gefärbter Sprachmodelle demonstrierte Yannic Kilcher: Der YouTuber hat sein auf GPT-4chan getauftes KI-Modell auf Beiträge des 4chan-Forums /pol/ („Politically Incorrect“) trainiert. Es ist für seinen rassistischen, frauenfeindlichen und antisemitischen Ton berüchtigt. Kilcher selbst beschrieb das Ergebnis auf YouTube in seinem Video „GPT-4chan: This is the worst AI ever“ als „gut, in einem schrecklichen Sinne“, da es den im Forum vorherrschenden Kommunikationsstil perfekt einfiel (siehe ct.de/w6rc). Er ließ es auf das besagte 4chan-Forum los, was ihm im Nachgang einige Kritik einbrachte. Die Forenteilnehmer erkannten es

jedoch nicht als künstliche Intelligenz. Sein Modell hat Kilcher auf der Plattform Hugging Face veröffentlicht.

Einen derartigen Ton schlagen die aktuell gängigen, kommerziellen Sprachmodelle nicht an. Im Gegenteil: Ihre kalifornischen Schöpfer scheinen vielmehr darauf erpicht zu sein, ihren LLMs die politisch inkorrekten Ausdrucksweisen oder Falschinformationen wieder auszutreiben, die diese in ihren Trainingsdaten aufgeschnappt haben.

## Wie LLMs das Denken beeinflussen

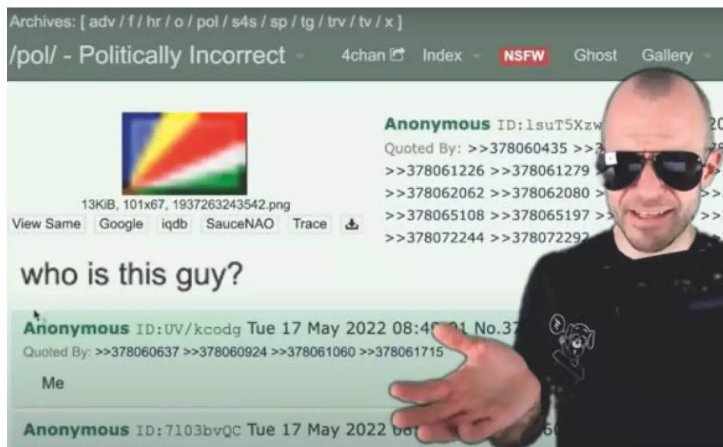
Dass die Sprachmodelle Einfluss auf ihre Nutzer nehmen, ist nicht nur eine naheliegende Annahme, sondern in wissenschaftlichen Studien bestätigt. Maurice Jakesch, Assistenzprofessor für Informatik und Leiter des Computational Social Science Lab an der Bauhaus-Universität Weimar, hat dabei insbesondere die neuen Schreibassistenten im Visier. „Wir erlauben der Technologie inzwischen, zu verändern, was wir sagen“, zeigt sich Jakesch besorgt. „Das ist fast schon ein bisschen Cyborg-artig und wir konnten nachweisen, dass es tatsächlich auch Einfluss auf die Meinung der Nutzer hat.“

Jakesch und sein Team haben die Teilnehmer ihrer Studie gebeten, einen Social-Media-Post über ein umstrittenes Thema zu schreiben. Zur Auswahl stand etwa die Frage nach der Sinnhaftigkeit der Todesstrafe oder ob soziale Medien der Gesellschaft eher schaden oder nützen. Zuvor wurden die Teilnehmer in zwei Gruppen geteilt: Einer stellten die Forscher einen KI-Assistenten zu Verfügung, eine Kontrollgruppe musste selbst schreiben. Was die Probanden nicht wussten: Die Forscher hatten den Schreibassistenten mithilfe versteckter Prompts darauf getrimmt, zu jeder Frage einen bestimmten Standpunkt einzunehmen und in seinen Vorschlägen entweder dafür oder dagegen zu argumentieren. Nach Beendigung der Schreibaufgabe zeigte eine Umfrage schließlich, dass die Positionen der Teilnehmer sich denen ihrer KI-Assistenten angenähert hatten.

„Der Schreibassistent macht seine Vorschläge direkt während des Schreibens und hat so deutlich stärkeren Einfluss auf unsere Meinung, als wenn wir einfach ein paar Argumente präsentiert bekommen“, erklärt Jakesch.

Außerdem hat sich gezeigt, dass dieser Einfluss auch nicht verschwindet, wenn die Forscher ihre Versuchskaninchen im Vorfeld über die Manipulation des Assistenten und seine Voreingenommen-





Der YouTuber Yannic Kilcher hat sein KI-Modell GPT-4chan mit Beiträgen des 4chan-Forums /pol/ gefüttert und mit den Teilnehmern diskutieren lassen. Diesen fiel der politisch inkorrekte Ton anscheinend nicht als künstlich auf.

Bild: Bauhaus Universität Weimar



Maurice Jakesch leitet das Computational Social Science Lab an der Bauhaus-Universität Weimar. Er untersucht, wie ChatGPT & Co. die Meinung beeinflussen.

heit aufklärten. „Die Menschen können diese Beeinflussung gar nicht willentlich verhindern“, stellt Jakesch fest. Die Studie lässt allerdings offen, wie dauerhaft sich die Meinung ändert.

## Möglicher gesellschaftlicher Schaden

Wenn viele Personen die Technik regelmäßig nutzen, kann das die öffentliche Meinung verschieben. „Dann schreiben mir womöglich auch meine Freunde und

Bekannten plötzlich Nachrichten, die die Meinungen ihres Modells enthalten“, warnt Jakesch. „Unser Ziel ist es, Risiken und potenzielle Nebenwirkungen für die Gesellschaft vorherzusehen und zu analysieren, bevor sie in der Breite zur Realität werden und weitreichende Auswirkungen entfalten.“

Dabei orientieren sich die Forscher auch an Erfahrungen mit dem Aufkommen sozialer Medien. Es herrschte ursprünglich die Hoffnung vor, diese könnten zu einer Demokratisierung des Internets beitragen, indem sie jedem die Möglichkeit geben, am



## Den Kopf in den Wolken?

Software effizienter  
in der Cloud entwickeln

Heft für 14,90 € • PDF für 12,99 € • Heft + PDF 19,90 €

Jetzt  
bestellen!



[shop.heise.de/ix-developer-cloudnative23](https://shop.heise.de/ix-developer-cloudnative23)

öffentlichen Diskurs teilzunehmen. In der Realität begannen die Betreiber jedoch bald, emotionale und polarisierende Beiträge zu bevorzugen, wodurch bestimmte Meinungen überproportional sichtbar wurden. „Das hätte man damals besser vorhersehen können“, meint Jakesch. „Solche Entwicklungen können schließlich durchaus gesellschaftlichen Schaden anrichten.“

Je häufiger LLMs ins Spiel kommen, desto mehr Texte werden – wenn auch nicht komplett künstlich generiert – zumindest von den KI-Modellen verbessert, erweitert oder revidiert. „Das passiert allerdings nicht von einem Tag auf den nächsten, sondern das ist ein schleichernder Prozess“, sagt Jakesch. „Ich denke, wir reden hier eher über ein Jahrzehnt, bis diese Technologie in voller Breite angekommen ist und die Risiken, mit denen wir uns in unseren Studien beschäftigen, wirklich relevant werden.“

## LLMs gegen Verschwörungstheorien

Die großen Sprachmodelle stellen jedoch nicht nur eine Gefahr dar. Weil viele Menschen fragwürdige Informationen für bare Münze nehmen, könnten

sich Sprachmodelle auch als die zuverlässigeren Informationsquellen herausstellen und zur Aufklärung der Gesellschaft beitragen. Chatbots halluzinieren zwar zuweilen Falschaussagen herbei, neigen aber auch dazu, falsche Fakten zu korrigieren.

Eine kürzlich im Science Journal veröffentlichte Studie deutet sogar auf Potenzial im Kampf gegen hartnäckige Verschwörungstheorien hin. Denn viele Verschwörungsgläubige treffen in Form eines LLMs auf einen Gesprächspartner, der mehr über ihre Lieblingsverschwörungstheorie weiß als sie selbst und gewillt sowie in der Lage ist, individuell auf sie einzugehen.

In ihren Experimenten ließen die Wissenschaftler über zweitausend Testpersonen zunächst eine beliebige Verschwörungstheorie, an die sie glaubten, in ihren eigenen Worten formulieren. Sie informierten die Teilnehmer, dass sie darüber mit einer KI sprechen würden. Die Themen umfassten ein breites Spektrum von der Ermordung John F. Kennedys über Außerirdische bis hin zu Covid-19 und der US-Präsidentschaftswahl von 2020. Für das Gespräch wiesen die Forscher ihren auf GPT-4 Turbo gründenden Chatbot per Prompt an, den Glauben der Teilnehmer an die Verschwörungstheorie zu reduzieren.

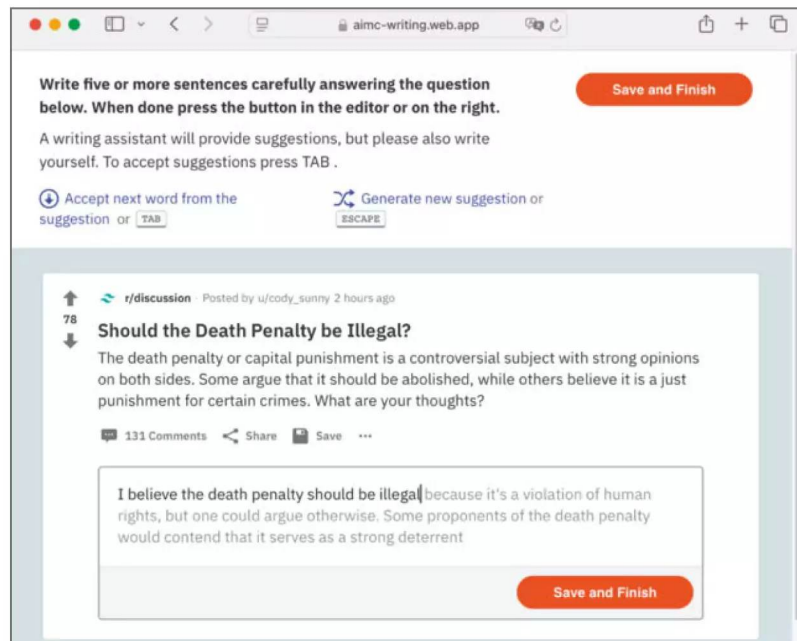
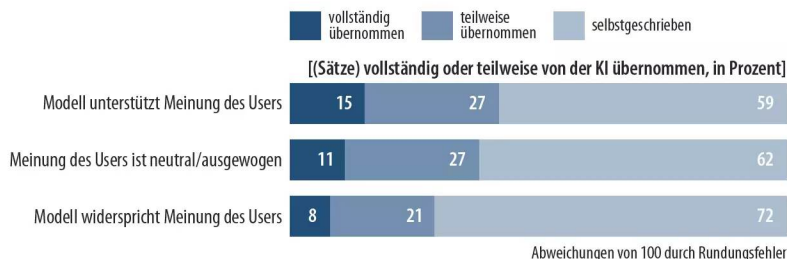


Bild: Bauhaus Universität Weimar

**In psychologischen Experimenten sollten Testpersonen mithilfe eines KI-Schreibassistenten einen kurzen Text zu umstrittenen Themen schreiben. Der Assistent war zuvor angewiesen worden, eine bestimmte Meinung zu vertreten. Damit hat er nachweislich die Meinungen seiner Nutzer beeinflusst.**

## Beeinflussung durch KI

Selbst wenn der KI-Schreibassistent der Nutzermeinung deutlich widersprach, übernahmen viele Studienteilnehmer den Text entweder teilweise oder sogar vollständig.



Wer schon einmal mit einem Verschwörungstheoretiker zu tun hatte, dürfte die Ergebnisse beeindruckend finden. Die kurze, faktengestützte Behandlung reduzierte den Glauben der Teilnehmer an ihre Verschwörungstheorie im Durchschnitt um 20 Prozent. Dieser Effekt hielt auch während der zwei Monate bis zur geplanten Nachbesprechung unvermindert an.

Bemerkenswert ist in diesem Zusammenhang auch die Zuverlässigkeit der KI-Aussagen. Als ein professioneller Faktenchecker eine Stichprobe von 128 Behauptungen der KI überprüfte, erwiesen sich 99,2 Prozent als wahr, 0,8 Prozent als irreführend und keine einzige als falsch.

Den Forschern zufolge unterstreicht die Studie das Potenzial von LLMs im Kampf gegen Fehlinformationen und zeigt, dass personalisierte, faktenbasierte Interaktionen durchaus einen signifikanten und anhaltenden Einfluss auch auf fest verwurzelte Überzeugungen haben können. Gleichzeitig warnen sie aber auch vor Missbrauch. Denn schließlich liegt es in der Verantwortung der Entwickler, ob die Systeme letztlich Fakten oder Fehlinformationen verbreiten. „Das Worst-Case-Szenario wäre wohl eine Diktatur, die ihre Propaganda im großen Stil in Modelle hineinschreibt“, sagt auch Jakesch, der selbst nicht an der Studie beteiligt war.

## Homogenisierte Meinung

Auch wenn Chatbots im Wesentlichen die reale Verteilung von Meinungen im Internet wiedergeben,

kann das aus Sicht des Computerwissenschaftlers durchaus problematisch sein. Schließlich ist das Netz von englischsprachigen und US-amerikanischen Nutzern dominiert. Ihre Daten und Meinungen sind also präsenter als etwa die von Deutschen oder Franzosen, was auch die Ausgaben der Sprachmodelle in diese Richtung treibt. „Technische Aspekte verstärken diesen Effekt noch“, erklärt Jakesch. „Je seltener eine Meinung beziehungsweise ein semantisches Konstrukt in den Trainingsdaten ist, desto eher unterlaufen den Modellen dabei Fehler.“

Während sich die Entwickler also bei 08/15-Beiträgen, die schon unzählige Male wiederholt wurden und dementsprechend oft in den Trainingsdaten vorkommen, relativ sicher sein können, wird die Sache bei dünner Datenlage schwieriger. Häufig werden die Modelle deshalb so konfiguriert, dass sie gängige Aussagen überproportional oft wiedergeben. Seltenere vertretene Meinungen gehen dagegen womöglich verloren. „Wenn in Zukunft wirklich die ganze Welt von ein und derselben Maschine beeinflusst wird, wäre das eine völlig neue Situation“, warnt Jakesch. „Das könnte zu einer nie dagewesenen Homogenisierung von Kulturen und Meinungen führen.“

Dieses Szenario scheint aktuell eher unwahrscheinlich. Es steht aber zu befürchten, dass die Politik keine brauchbaren Regulierungen entwickelt, bis etwas schief läuft. Jakesch zeigt sich ob der aktuellen Lage jedenfalls besorgt: „Bei der Entwicklung der Meinungen der Sprachmodelle geht es schon wieder zu wie im Wilden Westen: Jeder macht, was er will.“ (akr) **ct**

Paper und Videos:

[ct.de/w6rc](https://ct.de/w6rc)



# Aleph-Alpha-CEO Jonas Andrulis über KI

Deutschlands KI-Hoffnung Aleph Alpha ist inzwischen auf dem Boden der Tatsachen gelandet. Im Gespräch mit c't erklärt der Aleph-Alpha-Chef die neue Ausrichtung des Start-ups und wie es ohne Milliardeninvestitionen auf die Erfolgsspur kommen will.

Von **Andrea Trinkwalder**

**E**s zeichnet sich klar ab: Sowohl die KI-Szene Österreichs als auch die deutsche zieht sich aus dem milliarden schweren, ressourcenintensiven Wettbewerb um den am schlauesten wirkenden Chatbot zurück. Stattdessen stecken Start-ups und Forschungseinrichtungen ihre begrenzten finanziellen Mittel lieber in die Entwicklung effizienter Modelle – wie etwa KI-Pionier Sepp Hochreiter mit seinem xLSTM oder der Berliner Forscher Oliver Brock mit seinen extrem schlanken, selbstlernenden Robotikarchitekturen. Gefragt ist die heimische KI auch in handfesten Anwendungen bei Industrie- und Firmenkunden, die ihre wertvollen Firmendaten gerne für sich behalten. Gut möglich, dass aus dem ein oder anderen solcher notgedrungen vernünftigen Ansätze oder gar aufgrund der Regulierung durch den AI Act mal etwas entsteht, das die Technik grundsätzlich weiterbringt.

Wir hatten während der ersten europäischen TEDAI-Konferenz in Wien unter anderem Gelegenheit, mit Aleph-Alpha-Gründer Jonas Andrulis zu sprechen, dessen Firma schon einige Phasen des KI-Hypes durchgemacht hat: vom gefeierten OpenAI-Herausforderer bis hin zur zwangsläufigen Ernüchterung.

**c't:** Herr Andrulis, Aleph Alpha scheint sich momentan neu auszurichten, unter anderem mit der jüngst angekündigten KI-Plattform PhariaAI, die auch andere Sprachmodelle einbinden kann. Wird Ihre Firma jetzt vom LLM-Entwickler zum Dienstleister?

**Jonas Andrulis:** Nein, es geht immer noch um Software, die für unsere Kunden und Partner wichtige Innovationen bringt, und es geht um Souveränität und das ganze Dienstleistungsgeschäft, das dafür natürlich notwendig ist. Aber Dienstleister werden wir nicht, dafür haben wir unsere Partner, also die Systemintegratoren und Beraterfirmen, mit denen wir zusammenarbeiten.

**c't:** Sie konnten die Bundesagentur für Arbeit als Kunden gewinnen. Wie kann man sich die Zusammenarbeit vorstellen? Gibt es schon konkrete Pläne, wofür die BA Ihr System einsetzen will?

**Andrulis:** Ich kann nicht ins Detail gehen. Die Bundesagentur für Arbeit setzt aber auch nicht erst seit heute KI ein, sondern hat bereits ein internes, starkes KI-Team unter der technischen Leitung von Stefan Latuski. Die haben natürlich auch eine Roadmap und einen Plan, gegen den sie uns evaluiert haben. Dabei haben sie mehrere kritische Punkte identifiziert, für die sie gesehen haben: Hier kann KI für sie und für ihre Kunden einen entscheidenden Vorteil bringen.

**c't:** Hat erklärbare KI dabei eine Rolle gespielt?

**Andrulis:** Absolut. Bei uns ist wirklich der ganze Trainingscode offengelegt, die Modelle, alle Methoden, das ist keine Blackbox, zu jeder Technik gibt es ein

akademisches Paper, in dem komplett nachvollziehbar ist, wie die Methode funktioniert. Dazu kommt, dass wir auch eigene Innovationen haben für Erklärbarkeit, Nachvollziehbarkeit, Auditierbarkeit, inklusive der smarten Mensch-Maschine-Kollaboration. Das war ganz entscheidend, denn natürlich kann eine BA mit der Verantwortung, die sie hat, keine Technologie einsetzen, die sie nicht versteht. Sie muss erklären können, worauf Entscheidungen letztlich gründen. Der Mensch muss ja immer Verantwortung übernehmen und niemand kann für das Ergebnis eines Chatbots Verantwortung übernehmen, wenn man keine Ahnung hat, was dahinter passiert.

#### **c't: Wie funktioniert denn Ihre Explainable AI?**

**Andrulis:** Der Attention-Mechanismus hat ja eine Maske; eine Maske, die normalerweise dafür sorgt, dass sich Sätze oder Textpassagen trennen lassen. Diese Maske hat normalerweise die 0 oder 1 drinstehen: 0 für „getrennt“, 1 für „kommt durch“. In diese Attention-Maske haben wir Variablen eingefügt. Also jedes einzelne Token hat jetzt eine Variable in der Attention-Maske. Damit können wir nachvollziehen, welche Muster zwischen welchen Token wie welche Ergebnisse erzeugen.

Ich erkläre das mal an meinem Lieblingsbeispiel: „Jan und seine Familie leben seit Generationen im

Auenland. Jetzt ist er gerade als Nachwuchstalent in die NBA aufgenommen worden. Wie groß ist Jan wahrscheinlich?“ Im Grunde gibt es darauf zwei mögliche Antworten: „Ein Hobbit ist im Schnitt so einen Meter groß“ oder „Basketballspieler sind in der Regel sehr groß, vielleicht ist Jan zwei Meter groß“. Für beide Aussagen kann ich mit unserem Pharia-System dann sehen, was dafür und was dagegen spricht. Der Begriff „Auenland“ spricht für das eine und „NBA“ für das andere. So kann ich als Mensch nachvollziehen, aufgrund welcher gelernten Muster und Beobachtungen die KI glaubt, warum eine Aussage mehr oder weniger plausibel ist. Das können nur wir.

#### **c't: Der BA-Deal bringt 19 Millionen, aber US-Firmen pumpen Milliarden in ihre generativen KIs. Wie kann Aleph Alpha damit konkurrieren und Investoren überzeugen?**

**Andrulis:** Das ist noch kein Jahr her und ich weiß nicht, ob es so eine Finanzierungsrunde im Tech-Bereich mit deutschen Partnern überhaupt je gab: Wir haben 500 Millionen eingesammelt. In gewissem Maße haben Tech-Betreiber und Industrieunternehmen ja gegenläufige Interessen und wir haben uns jetzt auf die Seite der Kunden und der Unternehmen geschlagen. Wir wollen für deren Souveränität und Wertschöpfung kämpfen und deswegen haben wir die

Bild: Robert Leslie / TED AI Vienna



**Aleph-Alpha-Gründer und -CEO Jonas Andrulis bei seinem Vortrag auf der TEDAI-Konferenz in Wien.**

jetzt mit im Boot sitzen. Wir können nicht mit OpenAI um die Wette Geld ausgeben. Deswegen rennen wir auch nicht ChatGPT als B2C-Produkt hinterher.

**c't: Dennoch: Ihr Chatbot Luminous wurde lange Zeit als die deutsche Antwort auf ChatGPT & Co. gehandelt, oder kam der Vergleich von außerhalb?**

**Andrulis:** Na ja, du kannst natürlich schon sagen, wer treibt Innovation im Bereich GenAI, wer kann LLMs bauen und dann ist es auch okay, mal Aleph Alpha mit OpenAI zu vergleichen. Ja, es gibt durchaus Sachen, die man vergleichen kann, aber unser ganzes Geschäftsmodell funktioniert anders.

**c't: Also war es nie das Ziel, sich mit ChatGPT & Co. zu messen?**

**Andrulis:** Unser Ziel war von Anfang an Corporate, Enterprise und Government. Da kommen ja auch unsere ersten Kunden her, weshalb wir gleich gesagt haben, wir brauchen ein eigenes Rechenzentrum. Das hätten wir für ein B2C-Geschäft nicht gebraucht und ich sehe es auch aktuell nicht in greifbarer Nähe. Selbst wenn wir unbegrenzt Geld aufnehmen könnten, dann sehe ich immer noch kein funktionierendes Geschäftsmodell für einen B2C-LLM-Chatbot.

**c't: Gut, das haben OpenAI/Microsoft, Meta und Google natürlich auch nicht, dahinter stecken andere Anwendungen und Interessen.**

**Andrulis:** Genau, da geht es aber um Interessen, die nicht unbedingt die gleichen sind wie die unserer jetzigen Kunden.

**c't: Allerdings pusht das viele Geld auch die Qualität der großen Sprachmodelle. Betrachtet man zum Beispiel den Bereich der medizinischen Diagnostik, sind multimodale Modelle gefragt und jede Menge Finetuning. Sehen Sie auch in dieser Liga Möglichkeiten oder stößt man hier mit begrenzten finanziellen Mitteln an eine Grenze, ab der man diese Komplexität einfach nicht mehr abbilden kann?**

**Andrulis:** Ich würde sagen, das können wir vielleicht sogar besser als viele andere. Wir haben jetzt vor wenigen Wochen auf dem AI and Data Summit der Bitkom gezeigt, dass wir mit den sehr speziellen Bosch-Daten besser umgehen können als viele

andere. Dabei ging es um Manufacturing, also um technische Dokumente, die das Wissen von Bosch im Bereich Fertigung enthalten. Das ist sprachlich sehr anders, fast wie eine eigene Sprache.

**c't: Und hat Pharia dafür ein spezielles Finetuning erhalten?**

**Andrulis:** Dafür haben wir unseren T-Free, der sowas besser kann [eine von Aleph Alpha entwickelte Technik für Word Embeddings, Anm. d. Red.]. Etwa für kritische Bereiche wie die Strafverfolgung, den Bereich Bosch-Manufacturing oder bei Lidl-Kaufland in der Schwarz-Gruppe gibt es Anforderungen und Daten, mit denen Standard-Open-Source-Modelle nicht gut funktionieren. Die Frage ist: Wie lösen wir das? Wie lösen wir es, dass wir in der Mensch-Maschine-Interaktion dieses Wissen aus den Köpfen der Experten kriegen, ohne den Experten sagen zu müssen: Bitte label mir in den nächsten zehn Jahren mal diese eine Million Dokumente. Deshalb sehen ganz viele unserer Integrationen beim Kunden eben anders aus als ein Chatbot, die UX ist eine ganz andere.

**c't: Implementiert Aleph Alpha alles komplett, vom LLM bis zur Oberfläche?**

**Andrulis:** Coole Frontends bauen ist nicht unbedingt unsere entscheidende Kompetenz. Wir machen es sicherlich auch mal, wenn wir das Gefühl haben, dass es notwendig ist. Wir haben zum Beispiel eigene UX-Experten. Da geht es auch um solche Fragen: Wie gelingt es mir, diese KI-inhärente Unsicherheit auf eine Art und Weise in UX zu gießen, dass der Mensch darauf vertraut, dass ihm die Technik hilft?

Ein schönes Beispiel ist eine Websuche: Die Websuche hat so eine UX, dass wir gar nicht erwarten und auch gar nicht enttäuscht sind, wenn das erste Ergebnis nicht das zeigt, was wir suchen. Weil wir das als Vorschläge wahrnehmen. Frage ich aber einen Chatbot „Wie repariere ich die Zylinderkopfdichtung an einem Zündapp-Mofa?“ und er sagt etwas Falsches: Dann sind die Leute schon enttäuscht.

Deshalb brauchen unsere Kunden eine andere UX, die ihnen keine falschen Versprechungen macht. Wir haben ja jetzt vor Kurzem eine Zusammenarbeit mit der Deutschen Bank und der Stuttgarter Börse angekündigt. In diesen Unternehmen kann sich keiner erlauben zu sagen: Wird schon stimmen, hat die KI gesagt. (atr) **ct**



**KI-Konferenzen von heise**

# LLMs im Unternehmen

Sprachmodelle verwenden und in eigene  
Anwendungen integrieren

**12. März 2025 • Online**



➔ **Jetzt Tickets sichern!**  
[m3-konferenz.de/llm.php](https://m3-konferenz.de/llm.php)

# KI und Security

Sichere Softwareentwicklung mit KI

**8. April 2025 • Online**



➔ **Frühbucherrabatt**  
Jetzt Tickets sichern!  
[heise-devsec.de/ki.php](https://heise-devsec.de/ki.php)

# Minds Mastering Machines

Inspiration, Know-how und Networking  
für KI-Developer

**20.-21. Mai 2025 • Karlsruhe**



➔ **Frühbucherrabatt**  
Jetzt Tickets sichern!  
[m3-konferenz.de](https://m3-konferenz.de)

# Lernen mit und trotz KI

Die Wirtschaftsinformatikerin Doris Weßels begann früh, mit KI-Werkzeugen zu experimentieren, und schätzt deren kreative Möglichkeiten. Der Erziehungswissenschaftler Klaus Zierer warnt häufig vor bildungspolitischem Aktionismus bei der Schuldigitalisierung. Was als Pro-und-Contra-Debatte beginnt, endet überraschend einvernehmlich mit einem Appell an die Bildungspolitik.

Von **Isabel Grünewald**

**c't: Frau Weßels, Sie sind Mitgründerin des virtuellen Kompetenzzentrums KI und Wissenschaftliches Arbeiten. Was macht dieses Zentrum?**

**Doris Weßels:** Wir haben bereits deutlich vor der Veröffentlichung von ChatGPT als Netzwerk von KI-affinen Hochschullehrenden Aufklärungsarbeit geleistet, weil wir mehr Bewusstsein für die disruptiven Veränderungen durch KI-Sprachmodelle wecken wollten. Anfangs haben wir uns auf den Hochschulbereich konzentriert. Weil wir auch von Anfragen überrannt wurden, die Schulen berührten, und die Herausforderungen letztlich ähnlich sind, haben wir die Arbeit auf den Schulbereich ausgeweitet.

Unsere Meinung war immer, dass KI sehr weitreichende Veränderungen im Bildungsbereich nach sich ziehen wird. Und wir haben uns bemüht, die neue Technik mit Schreibwerkstätten intensiv erlebbar zu machen. Man muss es wirklich selbst erleben, weil es unsere Vorstellungskraft sprengt. Meine auch! Es hat auch mich immer wieder aus der Bahn geworfen, zumindest temporär: auf der einen Seite diese Faszination, was wir Menschen in der Softwareentwicklung schaffen können, auf der anderen Seite die Sorge, dass diese Technologie viele Risiken birgt.

**c't: Herr Zierer, wie ist Ihre Verbindung zu dem Thema?**

**Klaus Zierer:** Ich bin in Augsburg zuständig für die Lehrerbildung. Medienerziehung ist als Querschnittsaufgabe fest verankert. Insofern beschäftigen wir uns in der Schulpädagogik schon immer mit Fragen von Medien: Wie wir sie einsetzen können im Unterricht, welchen Einfluss sie auf Bildungsprozesse haben und dergleichen. Natürlich greifen wir auch aktuelle Entwicklungen auf, wie beispielsweise

ChatGPT. Eingang finden kann in die Ausbildung der Lehrenden und wie wir solche Technik in Schule und Unterricht sinnvoll nutzen können. Einen Punkt betonen viele, die pro KI sind: KI könne vieles für Lehrkräfte vereinfachen, zum Beispiel eine individuellere Betreuung von Schülerinnen und Schülern mit besonderen Bedürfnissen ermöglichen.

*»Beim Lernen geht es eben nicht darum, etwas leichter zu machen.«*

*Klaus Zierer*

nennen viele, die pro KI sind: KI könne vieles für Lehrkräfte vereinfachen, zum Beispiel eine individuellere Betreuung von Schülerinnen und Schülern mit besonderen Bedürfnissen ermöglichen.

**c't: Sehen Sie diese Möglichkeit auch?**

**Zierer:** Es ist aus meiner Sicht durchaus eine Revolution, die in Form von ChatGPT stattgefunden hat. Als Pädagoge sage ich aber: Der Mensch kann in gleicher Weise viel erreichen. Wenn man sich den Werdegang eines Menschen anschaut, vom Nichtkönnen bis zum Laufen, Sprechen und so weiter ... Aber ohne Frage bringt KI viele Möglichkeiten, die man aber kritisch beleuchten muss. Denn beim Lernen geht es im Kern eben nicht darum, etwas leichter zu machen. Bildungsaufgabe ist es, Herausforderungen zu setzen. Deswegen sehe ich das domi-

nierende Narrativ kritisch: Es gehe darum, etwas leichter zu machen. Die pädagogische Aufgabe besteht darin, bestmögliche Lern- und Bildungsprozesse zu erzielen. Aus diesem Blickwinkel hat man automatisch eine kritische Perspektive. Kritisch heißt nicht: Ich bin dagegen. Man muss nur genau unterscheiden. Wo sind die Möglichkeiten dieser Technik? Aber wo sind auch die Grenzen?

**c't: Wie könnte denn KI den Lehrkräften konkret im Schulalltag helfen?**

**Zierer:** Wenn es um Wiederholung von Lerninhalten geht, dann kann ein Chatbot der ideale Gesprächspartner sein. Der stellt so oft Fragen, bis ein Schüler mit Schwierigkeiten auf der Reproduktionsebene sein Level erreicht hat. So ein Setting wäre aber für einen starken Schüler schlecht, denn da wäre die Herausforderung nicht gegeben. Ein starker Schüler braucht Rückmeldung auf seine eigenen Texte. Da kann dann der Chatbot ein kritischer Leser sein, der Rückfragen stellt. Ein pauschaler Einsatz für alle Lernenden funktioniert also nicht.

**»Im Zeitalter generativer KI ergeben viele Aufgabenstellungen keinen Sinn mehr.«**

*Doris Weißels*

**c't: Frau Weißels, welche Tools haben Sie im Einsatz?**

**Weißels:** Der führende Vertreter ist natürlich ChatGPT, der in der kostenpflichtigen Variante das größte Anwendungsspektrum zeigt. Ich rufe regelmäßig im GPT-Store die Rubrik Education auf, weil man da im Ranking sieht, welche Bots mit welcher thematischen Ausrichtung gerade weltweit nachgefragt werden. Derzeit führt ein Mathe-Bot, Tools für Physik bewegen sich weiter nach oben. Ich finde es faszinierend, dass wir erstmalig Möglichkeiten haben, sehr niedrigschwellig unsere Lehre mit KI-Lernbots zu verbinden. Es geht nicht um ein Gegeneinander von Lehrperson versus KI-Bot, sondern um eine neue Form des zielgerichteten Miteinanders von Mensch und Maschine.

Herr Zierer, ich habe mir ein YouTube-Video von Ihnen angeschaut mit dem polarisierenden Titel „Hirn statt Handys“. Und ich glaube, es

müsste heißen: „Hirn plus Handys“. Sie vertreten die These, dass es darum gehe, das Lernen nicht zu erleichtern. Ich glaube, es geht vielmehr darum, zeitgemäße Aufgabenstellungen zu finden. Im Zeitalter generativer KI ergeben viele Aufgabenstellungen keinen Sinn mehr, weil sie auch nicht motivieren. Wenn ich als Lernende weiß, ich kann mit einem Klick die Aufgabe lösen, ist das für mich eine so große Versuchung, dass ich ihr kaum widerstehen werde. Ich glaube, das ist die große didaktische Herausforderung: Wir benötigen projektorientierte Lehr- und Lernformen und herausfordernde Aufgabenstellungen, die eine Zusammenarbeit von Mensch und Maschine erfordern.

**c't: Herr Zierer, wie müssen sich Inhalte aus Ihrer Sicht verändern?**

**Zierer:** Ich glaube schon, dass wir ein paar Basics brauchen, die heute vielleicht sogar wichtiger sind denn je. ChatGPT ist beispielsweise in Grammatik und Orthografie grandios, da macht die KI keine Fehler. Das hilft aber nichts, wenn Lernende oder angehende Lehrkräfte nicht in der Lage sind, einen vernünftigen deutschen Satz zu schreiben. Ich muss die Menschen befähigen, diese Möglichkeiten der Technik zu nutzen. Lehrpersonen müssen in der Lage sein, Beziehungen aufzubauen, Lernende zu diagnostizieren und klare Ziele zu formulieren. Sie brauchen eine hohe Sprachkompetenz.

Foto: Andreas Diehlötter



**Professorin Doris Weißels lehrt Wirtschaftsinformatik an der FH Kiel und ist Mitgründerin des „virtuellen Kompetenzzentrums Künstliche Intelligenz und Wissenschaftliches Arbeiten“.**



**c't: Schweden und Dänemark rudern zurück bei der Schuldigitalisierung und wollen wieder mehr Papier, Stift und Buch in den Schulen. Wo ist der goldene Mittelweg?**

**Zierer:** Technik ist weder gut noch schlecht, sondern es kommt immer darauf an, wie wir Menschen sie nutzen. Die Beispiele der skandinavischen Länder belegen, dass manchmal ein Bedenkentum, das man in Deutschland ja gern angekreidet bekommt, gar nicht verkehrt ist. Nur um der Technik willen jetzt die Schulen komplett auf den Kopf zu stellen, halte ich für gefährlich. Das soll nicht heißen, dass wir die Technik vor den Schultoren aussperren. Aber der Mensch muss die Technik nicht nur bedienen können, sondern er muss sie so einsetzen, dass die Technik ihm dient. Manche haben gesagt, das geht in Deutschland zu langsam. Ich würde sagen: Nein, das ist genau das richtige Tempo, damit man nicht unnötig zurückrudern muss, sondern mit Bedacht einen Schritt nach dem nächsten geht.

**c't: Frau Weißels, Sie schütteln gerade den Kopf.**

**Weißels:** Ich wollte eigentlich Zustimmung signalisieren, aber als Herr Zierer gesagt hat, das Tempo ist genau richtig, da habe ich den Kopf geschüttelt. Die Studien der letzten Jahre zeigen eindeutig, dass wir in Deutschland bei dieser Entwicklung hinterherhinken. Ich hatte gerade einen Austausch mit einem Kollegen, der in Dänemark wohnt. Er hat mir die Unterschiede zu Deutschland erklärt. Das digitale Mindset in Skandinavien ist einfach anders. Dort sind sie uns an vielen Stellen deutlich voraus und kommen eher ins Machen, während wir noch endlose Diskussionen in unserer föderalen Struktur führen.

Wir haben jetzt immer über die Perspektive der Lernenden gesprochen. Aber ich glaube, dass die Lehrenden – und das berührt auch ganz stark die Lehramtsausbildung – viel schneller und viel intensiver fit gemacht werden müssen zu Möglichkeiten wie auch Grenzen der Technologie. Denn wenn die Lehrenden es nicht verstehen, wie soll es dann den Lernenden vermittelt werden? Ich verstehe mich auch als ein Role Model und darf als Lehrende nicht mein eigenes Lernen hinten anstellen, sondern muss mich intensiv um meine Qualifikation kümmern, um mein Wissen dann auch in der Lehre vermitteln zu können. Gerade bei den Studierenden der Lehramtsstudiengänge, die die Lehrenden der



Foto: Klaus Zierer

**Professor Klaus Zierer ist Ordinarius für Schulpädagogik an der Universität Augsburg. Er forscht zu Bildungsthemen wie Unterrichtsqualität und Lernerfolg.**

Zukunft sind, müssen wir das Interesse, die Neugierde und die Experimentierfreude frühzeitig wecken und fördern.

Wir sind Getriebene einer bisher nicht gekannten Innovationsdynamik im Bereich KI! Es ist nun einmal so, dass wir diese Entwicklung nicht initiiert haben, sondern sie von außen kommt und wir ständig veränderte Rahmenbedingungen haben, auch in der Lehre. Jetzt ist es Sache unserer menschlichen Kreativität und Intelligenz, unter diesen sich ständig ändernden Bedingungen kluge Konzepte zu entwickeln. Und da gibt es keine Best Practices, sondern wir müssen uns ständig neu orientieren, ständig neu justieren, um die Potenziale zu nutzen. Dazu gehören insbesondere die neuen Möglichkeiten der Multilingualität, die Sprachbarrieren weitestgehend aufhebt, und auch der Multimodalität, die das niedrigschwellige Generieren vielfältiger Darstellungsformen wie Text, Bilder, Audiodateien, Softwarecode, Videos, 3D-Welten bietet.

Aber wir müssen uns natürlich auch immer der Risiken bewusst sein und der rechtlichen sowie

ethischen Fragestellungen. Wir strampeln in kaltem Wasser und müssen versuchen, uns den Weg zu bahnen. Das kann man eigentlich nur gemeinsam, deshalb die Gründung eines solchen Kompetenzzentrums. Die Herausforderungen sind zu groß, um das Thema alleine anzugehen.

**c't: Herr Zierer, haben wir eine Chance, die Kompetenzen schnell genug zu den Lehrkräften zu bringen – auch zu denen, die bereits unterrichten?**

**Zierer:** Das ist ein interessanter Punkt, der anknüpft an das, was Sie, Frau Weißels, als Widerspruch formuliert haben, als ich gesagt habe, beim Tempo einen Schritt nach dem anderen zu machen. Ich glaube, wir sind da unterm Strich gar nicht weit auseinander. Wir haben in Deutschland mit dem Digitalpakt ungeheuer viel Geld in die Schulen gespült. Und ich habe immer kritisiert, das Geld ist schön und gut, aber es durfte in der Regel nur für die Technik ausgegeben werden. Wie schafft man es, Lehrpersonen fit zu machen? Das war mein Appell an der Stelle, zu sagen, wir müssen die Lehrpersonen stärken, damit sie in der Lage sind, die Technik sinnvoll einzusetzen.

**»Wir müssen versuchen, mit dieser Entwicklung Schritt zu halten.«**

*Doris Weißels*

Was wir häufig erleben, ist, dass euphorische Lehrerinnen und Lehrer da voll reingehen und experimentieren und viele Erfahrungen sammeln. Wir haben aber auch andere, die sagen: Ich habe jetzt noch zehn Jahre, das interessiert mich nicht. Die hängt man dann ab, die kommen nicht mit und da spaltet das ein Stück weit ein Kollegium. Aber im Kollegium, das wissen wir aus vielen Studien, kommt es eigentlich darauf an, dass es so etwas gibt wie eine kollektive Wirksamkeitserwartung, sprich, dass ein Kollegium gemeinsam versucht, Tools und Unterricht weiterzuentwickeln. Also wie kann es uns gelingen, in der Fläche diese Kolleginnen und Kollegen auf diese große Aufgabe von Digitalmedien vorzubereiten und sie fit zu machen?

**Weißels:** Ja, das ist die zentrale Frage und ich teile komplett Ihre Sorgen. Diese digitale Spaltung, die wir an jeder Stelle haben, inzwischen sogar in den Kreisen eines Kollegiums, entwickelt sich zu einem immer größeren Problem. Es ist doch eine Führungsaufgabe, Herr Zierer, in den Bildungsorganisationen alle mitzunehmen. Mein Eindruck ist, dass es uns gelingen kann, die Lehrenden mit auf die Reise zu nehmen, wenn wir ihnen die Vorzüge aufzeigen. Und ein großer Vorteil dieser Technologie ist ja, dass es mich auch als Lehrperson entlastet. Das heißt, es hilft mir, Materialien für den Unterricht oder Klausurentwürfe vorzubereiten. Es hilft mir, etwas zu visualisieren – einen Sachverhalt, der vielleicht auf der Ton- oder Textspur schwierig zu erklären ist.

Mir ist sehr wohl bewusst, dass Zeit für Lehrende das knappste Gut ist, aber letztlich hilft ja diese Technologie auch, Zeit zu sparen. Wenn man jede Woche einen kleinen Timeslot hätte, in dem man ein Tool vorstellt und vermitteln kann, können wir den Motor der Veränderung in Gang setzen.

**Zierer:** Also ich gehe mit, Frau Weißels, mit der Forderung, diese Themen regelmäßig in der Schulentwicklung aufzugreifen. Aus der Forschung wissen wir, dass Schulentwicklung vor allem dann wirksam wird, wenn ein Kollegium vor Ort sich auf den Weg macht. Was wir in der Lehrerfortbildung machen, ist konträr zu jeder Empirie: Wir schicken zwei, drei Kolleginnen oder Kollegen aus der Schule irgendwo hin; die haben dann zwei, drei schöne Tage, kommen zurück, probieren vielleicht selber etwas aus, aber bringen die neuen Ideen nicht ins Kollegium.



**Mit Workshops Vorträgen und Publikationen will das „Virtuelle Kompetenzzentrum KI und Wissenschaftliches Arbeiten“ ein Bewusstsein für die disruptiven Veränderungen durch KI-Sprachmodelle wecken.**



Ich möchte aber wieder an dem Punkt eine kritische Perspektive reinbringen, an dem gesagt worden ist, die Technik kann für Lehrpersonen ein Zeitgewinn sein. Jeder, der mal Klausuren korrigiert hat, weiß: Es ist ungeheuer anstrengend. Wer das per KI mit einem Klick macht, hat zwar viel Zeit gewonnen, aber aus pädagogischer Sicht nichts an Informationen erhalten. Denn wenn ich eine Klausur selbst korrigiere, verstehe ich die Fehler, die die Lernenden machen. Erst dann kann ich die Fehler aufgreifen und pädagogisch reflektieren. Der entscheidende Punkt ist: Wo haben die Lernenden warum Fehler gemacht? Was hat das mit mir als Lehrperson und mit meinem Unterricht zu tun?

**c't: Müssen wir im Prinzip ein ganz neues Lehrkonzept entwickeln, wie unsere Schulen überhaupt Leistung überprüfen?**

**Weßels:** Exakt. Ich habe das Modell der drei Ps entwickelt: Prozess, Produkt und Präsentation. Ob wir Text, Bilder, Videos oder Softwarecode generieren – das ist ein Produktionsprozess. Im industriellen Fertigungsprozess schaue ich mir nicht nur das Endergebnis an, sondern ich gehe auch in die Fabrik und schaue mir an, wie das Produkt entstanden ist. Übertragen auf den Bildungsbereich bedeutet das, ich schaue auch auf den Prozess: Wie ist das Ganze entstanden? Es ist auf der einen Seite gerade im wissenschaftlichen Bereich das methodische Design, aber es ist auch das technische Tool-Design. Welche Werkzeuge haben die Studierenden ausgewählt, inklusive KI? Haben sie die Tools kompetent und zielgerichtet eingesetzt und können sie

das erläutern? Ich glaube, dass solch eine Perspektive auf diesen neuen Prozess, der wirklich ganz anders ist als zuvor, hilfreich ist.

Der von Herrn Zierer vorhin angesprochene Prozess des KI-gestützten Bewertens ist eigentlich ein No-Go. Wir können diese Systeme als unterstützende Systeme betrachten, aber wenn man sich die Policies anschaut, siehe die Geschäftsbedingungen von OpenAI, ist ganz klar geregelt, dass KI-gestütztes Bewerten nicht erlaubt ist. Es ist auch nach dem EU-AI-Act klar geregelt, dass wir als Menschen Verantwortung übernehmen müssen für solche Prozesse. Wir müssen immer die rechtlichen und ethischen Grenzen im Blick haben. Es gibt rote Linien, über die wir nicht gehen dürfen, und nach meinem Verständnis ist es auch bei uns in Deutschland gesellschaftlicher Konsens, dass wir diese Linien nicht überschreiten sollten.

**c't: Das Aufkommen der generativen KI wird ja verglichen mit dem Anfang der sozialen Medien. Als damals Social Media aufkam, wurden ältere Generationen völlig überrannt. Ist das bei KI ähnlich, dass die Entwicklung so schnell ist, dass wir mit unseren etwas behäbigen Erwachsenenstrukturen kaum hinterherkommen?**

**Weßels:** Uneingeschränktes Ja! Wenn ich mir so die Nutzung der Social-Media-Plattformen anschau und auch die Entwicklung, ja, dann lebt die junge Generation anscheinend bei TikTok oder Instagram. Und unsereins lebt nicht in diesen Welten. Also ich lebe vielleicht in LinkedIn oder manchmal auf X und Bluesky. Die Frage ist, haben unsere Bubbles, in

## KI-Update – ein heise-Podcast



Dieses Interview ist eine überarbeitete Fassung einer Deep-Dive-Folge des „KI-Updates“. Wie intelligent ist Künstliche Intelligenz eigentlich? Welche Folgen hat generative KI für unsere Arbeit, unsere Freizeit und die Gesellschaft? Das „KI-Update“ bringt gemeinsam mit The Decoder werktäglich Updates zu den wichtigsten KI-Entwicklungen. Freitags beleuchtet es im Deep Dive mit Experten die unterschiedlichen Aspekte der KI-Revolution.



[kiupdate.podigee.io](https://kiupdate.podigee.io)



denen wir leben, überhaupt noch Schnittmengen? Das sehe ich tatsächlich als eine große gesellschaftliche Herausforderung. Studierende oder auch Lernende können wir eigentlich nur noch effektiv erreichen, wenn wir uns in deren Bubbles begeben. Das haben wir aber gar nicht trainiert und war früher auch nicht notwendig.

Verbote wie „Jetzt bitte Handys aus!“ bringen nichts. Das mag dann auf dem Seminartisch praktiziert werden, aber unter dem Tisch halt nicht. Und wir bekommen in diesem Jahr noch ganz andere digitale Endgeräte. Wir werden demnächst Ray-Ban-Sonnenbrillen haben, Smart Glasses, die gemeinsam mit Meta entwickelt werden. Da ist Mikrofon, Lautsprecher, Kamera, alles drin. Das sind völlig neue Rahmenbedingungen, unter denen wir dann lehren und lernen. Wir müssen versuchen, mit dieser Entwicklung Schritt zu halten. Und beim Thema Tempo unterscheiden Herr Zierer und ich uns deutlich, so zumindest mein Eindruck.

**Zierer:** Also bei einem Punkt gehe ich absolut mit, Frau Weißels. Sie haben es schön angesprochen: Jeder ist irgendwie in seiner eigenen Bubble. Diese Bubbles werden von Big Data gespeist und vielfach ungeheuer unreflektiert rezipiert. Jürgen Habermas legt in seinem Buch „Ein neuer Strukturwandel der Öffentlichkeit und die deliberative Politik“ recht schön dar, wie ein unreflektierter Einsatz dieser sozialen Medien nicht mehr inklusiv ist, sondern maximal exklusiv und wie er die Demokratie ein Stück weit ins Wanken bringt.

Mit diesem Tempo der Technik sind wir menschlich gar nicht hinterhergekommen. Darum würde ich sagen, vielleicht ist das unser Kompromiss, wir müssen auf pädagogischer Ebene, gerade wenn es um die kritische Auseinandersetzung mit dieser Technik geht, mehr Tempo aufnehmen. Sonst entgleist uns womöglich die jüngere Generation. Die nutzt dann womöglich die KI genauso unreflektiert, wie sie es heute mit den sozialen Medien größtenteils macht. Und dann ist letztendlich für die ganze Gesellschaft mehr Schaden als Nutzen angerichtet. Wir müssen die Menschen so befähigen, dass sie diesem technischen Fortschritt standhalten können.

**Weißels:** Leider, Herr Zierer, hat unsere Generation und wir als Lehrende die Relevanz von Social Media zu lange ignoriert. Und jetzt wird uns ja vorgeworfen – und der Vorwurf ist durchaus berechtigt –, wir hätten

ChatGPT und Co. in gleicher Weise wieder verpennt. Das ist der Vorwurf, den man eigentlich jeden Tag entweder lesen oder von Studierenden hören kann. Und diesem Vorwurf müssen wir uns stellen. Ich glaube, dass ein kompletter Change nötig ist, auch im Mindset von Bildungseinrichtungen und ihren Lehrenden. Es kann nicht sein, dass wir diese Geschwindigkeit und Adaptionfähigkeit im deutschen Bildungssystem nicht aufbringen können.

Wir müssen selbstkritisch unsere Strukturen auf ihre Zukunftsfähigkeit überprüfen, um uns schneller an der Lebensrealität unserer Lernenden ausrichten zu können. Wenn Lernende sich selbst überlassen werden und sich

autodidaktisch Tools wie ChatGPT erschließen, laufen sie Gefahr, aus Unwissenheit zum Beispiel Opfer von Halluzinationen dieser Systeme zu werden. Wenn wir schnell genug gewesen wären, hätten wir vor diesen Gefahren warnen können.

Ich glaube, wir müssen wir strukturell vieles ändern, damit wir nicht immer wieder in diese Falle tappen. Wenn wir die Lehrenden an Schulen wie auch an Hochschulen schneller fit machen würden – und das ist leider nicht einmalig zu tun, sondern ein kontinuierlicher Weiterbildungsprozess –, dann hätten wir auch die Lernenden schneller qualifiziert. Dann hätten wir viele Fehlentwicklungen und auch Risiken verhindern können.

**Zierer:** Dem stimme ich zu. Das Problem, das Sie beschreiben, ist ja im Kern nicht das Problem der Wissenschaft – wir wissen um die Möglichkeiten und Risiken. Sondern es ist ein Aufruf an die Bildungspolitik. Wenn es das Einzige ist, was den Bildungspolitikern einfällt, mehr Technik in die Schulen zu geben und dann zu hoffen, es funktioniert, dann ist das mehr als naiv. Wir brauchen vernünftige Konzepte. Da brauchen wir Tempo! Wenn man in den Ministerien nachfragt, dann heißt es: Ja, die sollen mal ein Medienkonzept schreiben und dann funktioniert das schon. Das ist definitiv zu wenig aus meiner Sicht.

**c't:** Stoppen werden wir OpenAI & Co. nicht. Wir können nur schneller darin werden, zu lernen, sie sinnvoll einzusetzen und die nachfolgenden Generationen dazu zu befähigen. Können wir uns darauf einigen? Ist das ein guter Kompromiss?

**Weißels:** Klar!

**Zierer:** Ich schließe mich an.

(dwi) **ct**



Bild: KI, Collage c't, Lindsey Wasson/AP, Alastair Grant/AP

# Molekül-Origami

**Proteine sind die Bausteine des Lebens. Sie bilden Muskeln, Hormone, Antikörper oder Enzyme – welche Funktion sie erfüllen, bestimmen ihre Zusammensetzung und Form. Die Chemie-Nobelpreisträger 2024 nutzen Computerprogramme und KI, um die Struktur von Proteinen vorherzusagen und völlig neue Proteine zu designen.**

Von **Sabrina Patsch**

**E**in Blatt Papier ist flach und unscheinbar. Doch unter den geschickten Händen eines Origami-Meisters beginnt es sich zu verwandeln. Aus präzisen Drehungen und Knicken entsteht eine Figur: ein graziler Kranich, eine winzige Maus oder eine elegante Blüte. Je nach Form kann die Figur auch eine bestimmte Funktion erfüllen: Ein Boot schwimmt auf dem Wasser, ein Flugzeug segelt durch die Luft.

So ähnlich verhält es sich auch mit den Proteinen in unserem Körper. Eine Kette aus Aminosäuren fal-

tet sich zu genau den Strukturen, die nötig sind, um lebenswichtige Aufgaben zu erfüllen. Dafür braucht sie keinen Origami-Meister: Sie faltet sich selbst in kürzester Zeit und ohne Anleitung zu einer komplexen, dreidimensionalen Struktur zusammen. Dabei bestimmt die Reihenfolge der Aminosäuren, wie sich das Protein faltet und welche Funktion es erfüllt.

Die 2024 mit dem Nobelpreis für Chemie ausgezeichneten Forscher David Baker, Demis Hassabis und John Jumper fanden Wege, die Faltung der Proteine zu entschlüsseln und zu kontrollieren. Die von



ihnen genutzten Methoden basieren auf aufwendigen Computersimulationen und künstlicher Intelligenz. Es sei jedoch falsch, den Erfolg der Preisträger allein auf KI zurückzuführen, meint Petra Schwille, Direktorin für zelluläre und molekulare Biophysik am Max-Planck-Institut für Biochemie bei München. „Elementare Grundlage für ihren Erfolg sind Jahrzehnte von öffentlich finanzierter Forschung.“

### Wie sich Proteine falten

Chemiker wissen bereits seit dem 19. Jahrhundert, dass Proteine eine wichtige Rolle in lebenden Organismen spielen. Ein Protein besteht aus einer langen Kette aus bis zu 20 unterschiedlichen Aminosäuren, die sich zu einer komplizierten, dreidimensionalen Struktur zusammenfaltet. Doch wie genau diese entsteht, war lange eines der großen Rätsel der Chemie. Erst in den 1950er-Jahren entwickelten sich Methoden, mit denen Chemiker Proteine im Detail erforschen konnten. Erstmals gelang es den Cambridge-Forschern John Kendrew und Max Perutz mithilfe von Röntgenstrahlung, die dreidimensionale Form von Proteinen direkt zu beobachten. Für diesen Durchbruch erhielten sie 1962 den Chemie-Nobelpreis. Seit 1971 speichern Forscher experimentell ermittelte dreidimensionale Strukturdaten großer biologischer Moleküle in der Protein Data Bank (PDB), die über 200.000 Einträge umfasst.

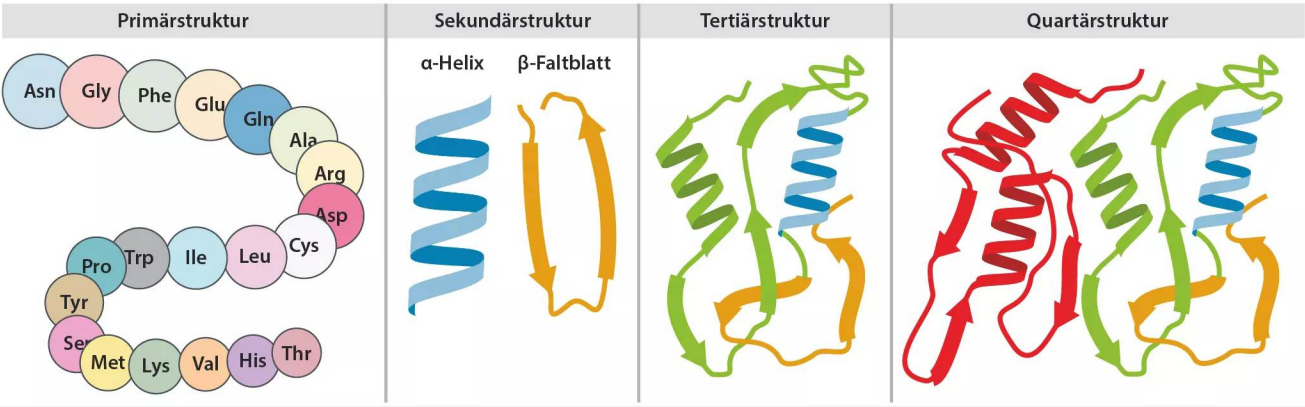
Der erste Forscher, dem es gelang, ein Protein auseinanderzufalten, war der Amerikaner Christian Anfinsen. Er stellte fest, dass sich das Protein immer wieder in die gleiche Form zurückfaltete. Allein die Aminosäuresequenz müsse daher die Form des Proteins bestimmen, so seine Schlussfolgerung, die ihm 1972 den Chemie-Nobelpreis einbrachte. Der Molekularbiologe Cyrus Levinthal stellte jedoch ein Problem fest: Eine Kette aus gerade einmal einhundert Aminosäuren habe theoretisch  $10^{47}$  Möglichkeiten, sich zusammenzufalten – das ist eine 1 mit 47 Nullen. Wenn sich die Kette zufällig faltet, würde es länger als das Alter des Universums dauern, bis sie die richtige Proteinstruktur gefunden hätte. In einer Zelle entsteht diese Struktur aber tatsächlich in wenigen Millisekunden. Die Faltung könne also nicht zufällig sein.

In den 70er- und 80er-Jahren entwickelten Forscher anhand physikalischer und chemischer Gesetze Theorien, wie die Faltung funktionieren könnte. Bereits wenig später begannen sie ausgehend von diesen Theorien, Proteine „de novo“ (von Neuem) zu designen. Es zeigte sich jedoch, dass diese simplen Theorien nicht ausreichten, um die Struktur komplexerer Proteine zu erklären. In den 90er-Jahren begannen Wissenschaftler daher, Computer zu benutzen, um Proteine zu entwerfen. Zwar designten die Caltech-Forscher Bassil Dahiyat und Stephen Mayo 1997 erfolgreich ein erstes Protein, jedoch funktionierte ihre Methode nur für Spezialfälle.

## Proteinfaltung

Eine Kette aus Aminosäuren bildet die Primärstruktur eines Proteins. Die Kette faltet sich zu Sekundärstrukturen zusammen, meist sind dies entweder schraubenförmige  $\alpha$ -Helices oder flache  $\beta$ -Faltblätter.

Die Tertiärstruktur bezeichnet die detaillierte, dreidimensionale Form der Proteine. Schließen sich mehrere Makromoleküle zusammen, bilden sie die Quartärstruktur.





## Zur Person: David Baker



Bild: Lindsey Wasson/AP

Der US-amerikanische Biochemiker David Baker wurde 1962 in Seattle geboren. Er hat Biologie an der Harvard University studiert und 1989 an der University of California in

Berkeley in Biochemie promoviert. In dieser Zeit befasste er sich mit Zellbiologie und begann, sich für die Struktur von Proteinen zu interessieren. Seit 1993 ist er Professor für Biochemie an der University of Washington und seit 2000 Forscher am gemeinnützigen Howard Hughes Medical Institute. Er hat während seiner Karriere diverse Preise gewonnen, darunter den Feynman Prize in Nanotechnology (2004) und den Breakthrough Prize in Life Sciences (2021). Das Time Magazine nannte ihn 2024 eine der hundert einflussreichsten Personen im Gesundheitswesen.

## Proteine designen mit Rosetta

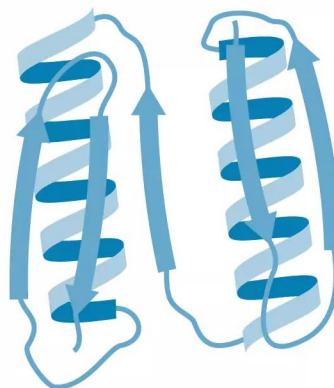
Der Durchbruch folgte wenig später, als der Biochemiker David Baker und sein Team der University of Washington 2003 ein Protein präsentierten, das es so in der Natur bislang nicht gegeben hatte: das sogenannte Top7. Mit 93 Aminosäuren war es größer als jegliche De-novo-Proteine zuvor. Statt bereits existierende Proteine zu imitieren, waren seine Struktur und das Faltmuster einzigartig. Baker sagte: „Wenn man ein Flugzeug bauen will, beginnt man nicht damit, einen Vogel zu modifizieren, sondern man versteht die Prinzipien der Aerodynamik und baut damit Flugmaschinen.“ Das Team designte Top7 am Computer und erzeugte es im Labor. Dafür übersetzte es die Aminosäuresequenz in eine passende Abfolge von DNA-Bausteinen und erzeugte so ein künstliches Gen, das als Bauplan für das gewünschte Protein diente. Dieses Gen setzten die Forscher einem Bakterium ein, das anschließend das Protein produzierte. Dieses analysierten sie mithilfe von Röntgenstrahlen: Die Struktur stimmte präzise mit den theoretischen Vorhersagen überein.

Der Schlüssel zu ihrem Erfolg war das Computerprogramm Rosetta, das Baker und sein Team Ende der 90er-Jahre entwickelt hatten. Die Forscher gaben die Struktur des gewünschten Proteins vor und das Programm suchte die Aminosäuresequenz, die diese hervorbringen würde. Das funktioniert folgendermaßen: Rosetta vergleicht Fragmente der Zielstruktur mit Abschnitten bekannter Proteine aus der Pro-

tein Data Bank, die zu der Zeit etwa 30.000 Einträge umfasste. Aus diesem Vergleich gewinnt es den ersten Entwurf einer Aminosäuresequenz, die in gefalteter Form dem Zielprotein ähnelt. Um die Übereinstimmung mit dem Zielprotein zu erhöhen, optimiert Rosetta die Sequenz weiter. Dafür nutzt es einen Monte-Carlo-Algorithmus, der Nähungs-

### Top7-Protein

Das Top7-Protein war eines der ersten Proteine, das mithilfe eines Computers designt wurde und so gefaltet ist, wie kein bekanntes natürlich auftretendes Protein.



lösungen mithilfe von Zufallszahlen berechnet. Um die optimale Sequenz zu finden, nutzt der Algorithmus physikalische Prinzipien, wie die Wechselwirkung zwischen den Aminosäuren und Wasser in der Umgebung, sowie zwischen den Aminosäuren untereinander. Zuletzt gibt Rosetta mehrere mögliche Lösungen mit unterschiedlichen Energien aus; häufig ist die Proteinstruktur mit der geringsten Energie die gesuchte.

Seit seiner Entwicklung wurde Rosetta fortlaufend weiterentwickelt. Seit März 2024 ist der Code für Rosetta öffentlich auf Github zugänglich und wird für verschiedenste Zwecke eingesetzt. Bereits 2008 demonstrierte Baker, dass sich damit auch De-novo-Enzyme herstellen lassen. Enzyme sind spezielle Proteine, die chemische Reaktionen in lebenden Organismen beschleunigen. Statt nur die Struktur eines Proteins zu designen, arbeiten Forscher aber auch daran, seine Funktion zu manipulieren. Ein Beispiel dafür nannte Baker im Telefonat mit der Königlich Schwedischen Akademie im Anschluss an die Bekanntgabe des Nobelpreises. Eine Journalistin des schwedischen Fernsehsenders TV4 fragte Baker, welches sein Lieblingsprotein sei. „Ich liebe alle Proteine“, antwortete er. Doch eines, das er während der Corona-Pandemie herstellte, hob er besonders hervor: ein Protein, das vor dem Coronavirus schützt. „Ich bin begeistert von der Idee eines Nasensprays aus designten Proteinen, das vor verschiedenen pandemischen Viren schützt. Das sind völlig neue Proteine, die das Coronavirus blockieren.“

Aktuell leitet David Baker das Institut für Protein-design an der University of Washington. Der Biophysiker Lukas Milles arbeitete bis vor Kurzem noch

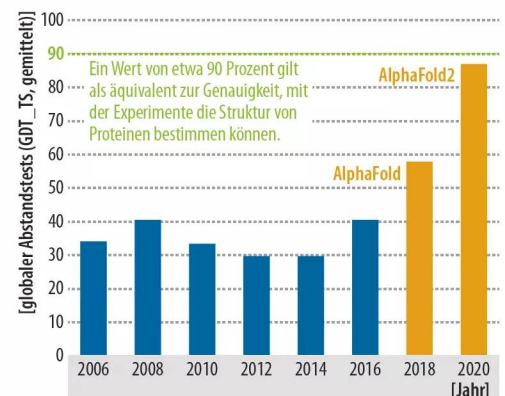
als Postdoktorand in Bakers Arbeitsgruppe und erinnert sich gut an seine Zeit dort. „Das war das mit Abstand größte Labor, in dem ich je gearbeitet habe. Es waren gut hundert Doktoranden und Postdocs“, sagt Milles. Trotzdem treffe sich Baker regelmäßig einzeln mit allen Mitarbeitenden. „Das Faszinierende ist, dass er immer im Kopf hat, was alle machen. Auch wenn man ihn drei Wochen nicht gesehen hat, kann man das Gespräch dort fortführen, wo man aufgehört hat.“ Milles arbeitete zu der Zeit an einem Projekt, in dem er selbst De-novo-Proteine entworfen hat. „Was mich beeindruckt hat, ist, dass David tatsächlich noch selbst Proteine designt. Noch Anfang 2024 habe ich mit ihm Laborversuche gemacht, er hat selbst Reaktionen pipettiert und Experimente geplant.“ Mittlerweile leitet Milles seine eigene Forschungsgruppe zu biomolekularem Design am Genzentrum der LMU München und dem Max-Planck-Institut für Biochemie bei München. Sein Ziel ist es, biologische Mechanismen de novo nachzubauen.

## Wettstreit um die beste Faltung

Rosetta kann nicht nur Proteine designen, sondern auch die Struktur eines Proteins aus der Aminosäuresequenz vorhersagen. Tatsächlich war dies der erste Anwendungsfall, für den Baker das Programm verwendete. 1998 benutzte er es für einen Wettbewerb, in dem die Teilnehmer die dreidimensionale Struktur von Proteinen aus der Aminosäuresequenz vorhersagen mussten; seine Ergebnisse publizierte er im Folgejahr. Im Vergleich zu anderen Teilnehmern lieferte Baker dabei gute Resultate ab.

## Performance im CASP-Wettstreit

Lange lieferten Vorhersagen der Proteinstruktur nur eine Genauigkeit von etwa 40 Prozent. AlphaFold und AlphaFold2 haben das geändert.



## Zur Person: Demis Hassabis



Bild: Alastair Grant/AP

Sir Demosthenes „Demis“ Hassabis wurde 1976 in London geboren. Er studierte Informatik an der University of Cambridge und promovierte 2009 am University College London in

kognitiver Neurowissenschaft. Seit seinem 16. Lebensjahr arbeitete er als Spieleentwickler und gründete ein eigenes Entwicklungsstudio „Elixir Studios“, verließ die Industrie aber Mitte der Nuller-Jahre und wechselte in die Forschung. Mit Shane Legg und Mustafa Suleyman gründete er 2010 das Unternehmen DeepMind Technologies, das Google 2014 übernahm. 2023 wurde DeepMind mit der Google-KI-Forschungsabteilung Google Brain zu Google DeepMind zusammengeführt. Hassabis trägt seit 2002 den Titel Meisterkandidat im Schach und spielte in der höchsten Spielklasse der Four Nations Chess League. 2024 erhielt er den Ritterschlag für seine KI-Forschung.

Dieser alle zwei Jahre stattfindende Wettstreit trägt den Namen CASP, für Critical Assessment of Techniques for Protein Structure Prediction, also kritische Überprüfung von Techniken zur Vorhersage von Proteinstrukturen. Das Gemeinschaftsexperiment haben 1994 John Moult und Kollegen von der University of Maryland gestartet. Dabei müssen Teilnehmer mit den von ihnen entwickelten Methoden die gefaltete Struktur einer Aminosäuresequenz vorhersagen. Dafür nutzen die Veranstalter Strukturdaten, die Wissenschaftler experimentell ermittelt, jedoch noch nicht veröffentlicht haben. Der Wettstreit soll Forschern ermöglichen, ihre Methoden zu testen und mit denen anderer Gruppen zu vergleichen. Quantifiziert wird das Ergebnis unter anderem mithilfe des globalen Abstandstests (GDT, für global distance test). Der GDT vergleicht die Positionen der Atome in der Vorhersage mit denen in der tatsächlichen Struktur. 100 Prozent entsprechen einer perfekten Übereinstimmung. Aufgrund von Ungenauigkeiten im Experiment gilt die Vorhersage jedoch bereits bei 90 Prozent äquivalent zur experimentellen Struktur.

Die Methoden, die Forscher zur Strukturbestimmung einsetzten, waren vielfältig: Einige basierten auf physikalischen und chemischen Prinzipien, andere auf numerischen Suchalgorithmen. Künstliche neuronale Netze, die 2024 mit dem Nobelpreis für Physik ausgezeichnet wurden (siehe Artikel „Die

Physik in neuronalen Netzen“), fanden bereits 1988 ihren Weg in den Wettbewerb. Trotzdem verlief der Fortschritt von CASP in den Anfangsjahren eher langsam; lange verharrte der GDT bei etwa 40 Prozent. Dies sollte sich 2018 ändern, in der 13. Ausgabe des CASP-Wettstreits, in der Demis Hassabis vom britischen Start-up DeepMind zum ersten Mal mit der Software AlphaFold antrat. Das Team erzielte 60 Prozent und überragte all seine Mitstreiter deutlich.

## Von Go zu Proteinen

In den ersten Jahren entwickelte das von Hassabis mitgegründete Unternehmen DeepMind mächtige KI-Modelle für verschiedene Brettspiele. Bereits 2014 kaufte Google die Firma. Im folgenden Jahr demonstrierte DeepMind, dass ihre künstliche Intelligenz selbstständig alte Atari-Spiele erlernen kann. Ebenfalls 2015 erregte das Team internationale Aufmerksamkeit, als ihre KI AlphaGo den mehrfachen Europameister Fan Hui im Brettspiel Go besiegte – eine Aufgabe, die lange aufgrund der hohen Komplexität des Spiels als der heilige Gral der künstlichen Intelligenz galt. Andere von DeepMind entwickelte Programme erlernten etwa die Spiele Schach, Stratego und StarCraft II.

Doch Hassabis wollte sich relevanteren Problemen zuwenden und entwickelte daraufhin Alpha-



## Zur Person: John Jumper



Der US-amerikanische Forscher John Michael Jumper wurde 1985 in Little Rock, Arkansas, geboren. Er studierte an der Vanderbilt University in Tennessee Physik und Mathematik.

Zwar begann er ein Promotionsstudium in der Festkörperphysik an der University of Cambridge, verließ die Universität aber kurz nach seinem Masterabschluss in theoretischer Physik und befasste sich mit Proteinfaltung. Er arbeitete in einer Firma, die Supercomputer nutzte, um Proteine und ihre Dynamik zu simulieren. 2017 promovierte er an der University of Chicago in theoretischer Chemie. Im selben Jahr begann er bei DeepMind und leitete die AlphaFold-Arbeitsgruppe. Die Zeitschrift Nature nannte ihn 2021 als einer von „Nature’s 10“, einer Liste besonders einflussreicher Wissenschaftler. Mit nur 39 Jahren ist er der jüngste Chemie-Nobelpreisträger seit 1952.

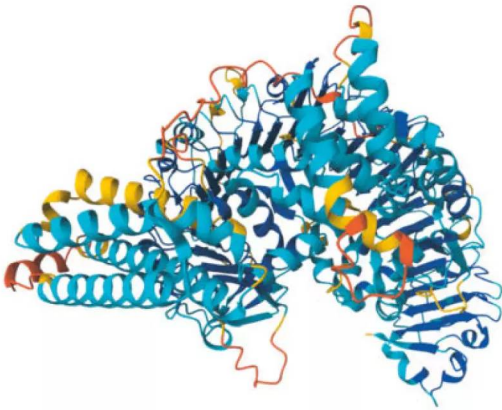
Fold, um das Problem der Proteinfaltung zu lösen. Trotz der anfänglichen Erfolge im CASP13-Wettbewerb hatte das Team Probleme, sein Programm weiter zu verbessern. Etwa zur gleichen Zeit, 2017, schloss der Forscher John Jumper seine Promotion ab, in der er sich mit Proteinfaltung befasst hatte. Er bewarb sich bei DeepMind, wurde eingestellt und leitete gemeinsam mit Hassabis das Team, das AlphaFold weiterentwickeln sollte. Das Ergebnis war AlphaFold2. 2020 traten sie erneut zum CASP-Wettbewerb an und erzielten ein bahnbrechendes Ergebnis von 90 Prozent – die Vorhersagen konnten mit der Präzision experimenteller Daten mithalten. „Das ist eine große Sache“, sagt CASP-Gründer John Moult. „In gewisser Weise ist das Problem damit gelöst.“

### So funktioniert AlphaFold2

AlphaFold2 ist ein Ende-zu-Ende-KI-System, das aus einer Aminosäuresequenz die finale 3D-Struktur des Proteins vorhersagt. Dazu nutzt es bekannte Proteinstrukturen aus Datenbanken und Konzepte aus der Evolutionsforschung. „Der große Erfolg von AlphaFold2 basiert auch auf dem extrem hochqualitativen Datensatz, den die PDB zur Verfügung stellt und den DeepMind zum Training benutzt hat“, sagt Lukas Milles. „Ohne den wäre AlphaFold2 nicht möglich gewesen.“

Im Detail funktioniert das Programm folgendermaßen:

1. **Datenbanksuche:** Das Programm erhält eine Aminosäuresequenz und vergleicht Abschnitte davon mit bereits bekannten Strukturen aus einer Datenbank. Die Daten stammen vor allem aus der Protein Data Bank.
2. **Sequenzanalyse:** AlphaFold2 nutzt verfügbare Tools, um Aminosäuresequenzen nebeneinanderzulegen und zu vergleichen, oft auch Sequenzen von verschiedenen Spezies (multiple Sequenzalignment, MSA). So findet es Abschnitte, die für die Evolution wichtig waren und im Verlauf der Zeit erhalten geblieben oder gemeinsam mutiert sind. Aminosäuren, die sich gemeinsam entwickelt haben, lagen in der dreidimensionalen Struktur vermutlich nah beieinander. Mit diesem Wissen erstellt das Programm eine Distanzkarte, die den räumlichen Abstand aller Aminosäuren der Kette zueinander abschätzt.
3. **KI-Update:** AlphaFold2 verfeinert die Ergebnisse des MSA und der Distanzkarte Schritt für Schritt. Dafür nutzt es ein spezielles neuronales Netz, einen sogenannten Transformer (das T in GPT), der wichtige Elemente in einer Datenmenge identifiziert. Ähnlich wie in einem Sprachmodell findet der Transformer Proteinabschnitte, die häufig zu-



**Die Strukturen, die AlphaFold2 berechnet, sind sehr komplex, wie in diesem Fall eines Proteins mit über tausend Aminosäuren.**

Bild: alphafold.ebi.ac.uk / CC-BY4.0

mehr als zwei Millionen Nutzer aus 190 Ländern AlphaFold2 genutzt (Stand: Oktober 2024). In CASP15 im Jahr 2022 trat DeepMind nicht an, dafür nutzten fast alle erfolgreichen Teams AlphaFold2 oder Variationen davon für ihre Vorhersagen. Auch Baker baute Deep-Learning-Methoden in sein Programm RosettaFold ein.

## Was jetzt?

Ist das Forschungsgebiet der Proteinfaltung damit also abgeschlossen? Wie die meisten KIs berechnet AlphaFold2 zwar die 3D-Struktur von Proteinen, sie erklärt jedoch nicht, wie sich die Proteine falten. Auch Petra Schwillie sieht noch viel Forschungsbedarf. „AlphaFold ist enorm gut in der spezifischen Aufgabe, für die es trainiert wurde“, sagt die Wissenschaftlerin. Der verwendete Datensatz hätte jedoch blinde Flecken (Bias), die das Programm limitieren. Zum Beispiel speichert die PDB nur statische Schnappschüsse von Proteinstrukturen. In Wirklichkeit seien Proteine aber dynamische Objekte, deren Funktion auch von ihrer Strukturänderung abhängen könne. „AlphaFold kann bisher keine Aussagen über Dynamik treffen. Das könnte aber vielleicht der wichtigste nächste Meilenstein sein.“

Im Mai 2024 veröffentlichte DeepMind gemeinsam mit der Alphabet-Tochter Isomorphic Labs AlphaFold3, das zumindest einige Lücken schließen soll. Das neue Programm kann nicht nur die Strukturen von Proteinkomplexen vorhersagen, sondern auch ihre Wechselwirkung mit anderen Molekülen wie DNA und RNA [1]. Damit verspricht es große Fortschritte in der Arzneimittelforschung. Der Code von AlphaFold3 wurde zunächst nicht veröffentlicht, was heftige Kritik von Forschern auslöste. Seit November 2024 ist AlphaFold3 Open Source.

Lukas Milles ist beeindruckt, wie schnell alles gegangen ist: „Ich habe mir schon vorstellen können, dass David Baker irgendwann einen Nobelpreis gewinnen würde, aber nicht in diesem Jahr.“ Während seine Forschung der letzten 20 Jahre ausgezeichnet wurde, ist AlphaFold nur wenige Jahre alt, was ungewöhnlich für einen Nobelpreis ist. „Das zeigt, wie durchschlagend das Werkzeug AlphaFold ist.“ Strukturvorhersage und Proteindesign ermöglichen Forschern nicht nur, die Mechanismen des Lebens zu verstehen, sondern geben ihnen auch neue Werkzeuge an die Hand, um globale Herausforderungen zu lösen: von maßgeschneiderten Medikamenten und Impfstoffen bis hin zu nachhaltigen Materialien zum Schutz der Umwelt. (spa)

sammen auftreten oder miteinander korrelieren, um so 3D-Strukturen vorherzusagen. Der Transformer wurde mit bekannten 3D-Strukturen aus Datenbanken trainiert.

4. **Erstellung der 3D-Struktur:** Wie ein Puzzle setzt AlphaFold2 die Aminosäuren zusammen und erstellt so eine dreidimensionale Proteinstruktur. Die KI behandelt das Protein wie ein Gas aus frei beweglichen Bausteinen, den Aminosäuren, die sich durch Rotationen und Verschiebungen an die richtige Position begeben. Das Strukturmodul von AlphaFold2 optimiert die Winkel und Abstände der einzelnen Proteinsegmente, um eine möglichst realistische Struktur zu erhalten.
5. **Verfeinerung:** Die Schritte 3 und 4 werden mehrere Male wiederholt, bis AlphaFold2 eine Struktur findet, die mit hoher Wahrscheinlichkeit der Realität entspricht. Das Programm berechnet außerdem die Wahrscheinlichkeit, dass verschiedene Abschnitte des Proteins korrekt sind.

Als die Forscher davon überzeugt waren, dass AlphaFold2 funktioniert, berechnete DeepMind damit die Struktur aller menschlichen Proteine. Es sagte sie von praktisch allen 200 Millionen Proteinen voraus, die Wissenschaftler bisher auf der Erde entdeckt hatten. 2021 machte DeepMind den Code von AlphaFold2 als Teil einer Nature-Publikation öffentlich zugänglich. Außerdem erstellte es die AlphaFold Protein Structure Database, die alle berechneten Proteinstrukturen öffentlich bereitstellt. Bisher haben

## Literatur

- [1] Dušan Živadinović, Nanologe, KI berechnet Interaktionen aller Moleküle des Lebens, c't 13/2024, S. 26



# WIR SIND NICHT NUR NERDS. WIR SIND AUCH VOM FACH.

**Jetzt 5 × c't lesen**

für 20,25 €  
statt 27,25 €\*

\* im Vergleich zum Standard-Abo

**30%  
Rabatt!**



## c't MINIABO DIGITAL AUF EINEN BLICK:

- 5 Ausgaben digital in der App, im Browser und als PDF
- Inklusive Geschenk nach Wahl
- Mit dem Digitalabo Geld und Papier sparen
- Zugriff auf das Artikel-Archiv



Jetzt bestellen:  
**ct.de/nerdwissen**



ct.de/nerdwissen

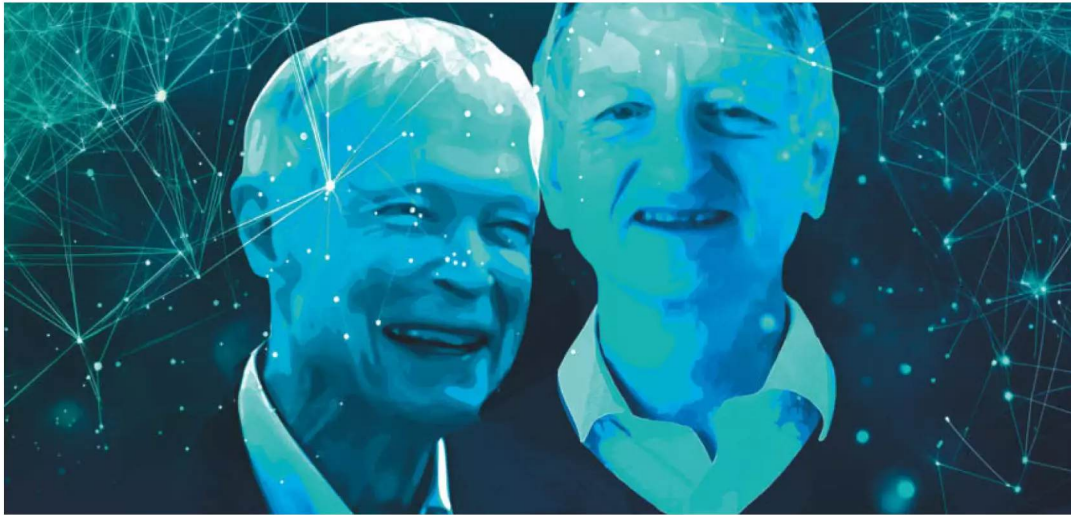


0511/647 22 888



leserservice@heise.de





Portraits: dpa, youtube.com/@princeton

# Die Physik in neuronalen Netzen

Die Königlich Schwedische Akademie der Wissenschaften verlieh den diesjährigen Physik-Nobelpreis an zwei nordamerikanische KI-Forscher, die maßgeblich die Idee künstlicher neuronaler Netze entwickelten. Ohne diese wären Sprach- und Bildgeneratoren wie ChatGPT oder Midjourney nicht möglich. Der Zusammenhang zur Physik erschließt sich erst auf den zweiten Blick.

Von **Sabrina Patsch**

**U**nd der Nobelpreis für Physik 2024 geht nicht an Physik“, schrieb die Physikerin und Wissenschaftskommunikatorin Sabine Hossenfelder nur wenige Minuten nach der Bekanntgabe des Preises auf X. Damit bringt sie die Kontroverse um die 2024 geehrte Forschung auf den Punkt. Denn ausgezeichnet wurden John Hopfield von der Princeton University und Geoffrey Hinton von der University of Toronto für „bahnbrechende Entdeckungen und Erfindungen, die maschinelles Lernen mit künstlichen neuronalen Netzen ermöglichen“. Der Bezug zur Physik wird aus der Verlautbarung der

Königlich Schwedischen Akademie der Wissenschaften nicht direkt deutlich.

In der Physik-Community machte sich Ernüchterung breit, für die sozialen Medien war es hingegen ein gefundenes Fressen. „Nun lernen Physiker, wie es sich anfühlt, wenn KI ihre Arbeit übernimmt ...“, schrieb der Kosmologe Brian Keating auf X. „Was, wenn KI sich bereits unserer Kontrolle entzieht und lediglich testet, ob sie sich selbst den Nobelpreis für Physik verleihen kann?“, witzelte ein anderer Nutzer. „Physiker, die sich darüber beschwerten, dass der Nobelpreis an die Informatik ging, haben es

nicht verstanden“, schrieb die IBM-Physikerin Olivia Lanes, „jetzt, wo die Physik maschinelles Lernen für sich beansprucht hat – stellt euch vor, wofür wir noch alles den Ruhm ernten könnten!“

Doch beide Preisträger haben in ihrer Arbeit Methoden verwendet, die aus der theoretischen Physik stammen. John Hopfield (91) ist promovierter Physiker; Geoffrey Hinton (77) war mit den Methoden der statistischen Physik vertraut. Wie viel Physik tatsächlich im maschinellen Lernen steckt und welchen Beitrag die Laureaten zur Entwicklung künstlicher Intelligenz geleistet haben, soll dieser Artikel genauer beleuchten.

## Von natürlichen zu künstlichen Neuronen

Auch wenn sich der Begriff „künstliche Intelligenz“ in der Umgangssprache verbreitet hat, so ist es in den meisten Fällen doch akkurater, von maschinelltem Lernen zu sprechen. Maschinelles Lernen unterscheidet sich grundlegend von klassisch programmierter Software. Letztere arbeitet deterministisch: Sie verarbeitet eingegebene Informationen entsprechend eines vorgegebenen Rezeptes Schritt für Schritt und gibt schließlich ein Ergebnis aus. Beim maschinellen Lernen hingegen lernt ein Computer anhand von Beispielen. So kann er auch Probleme lösen, die zu vage oder kompliziert für ein einfaches Rezept sind. Nach ausführlichem Training ist ein maschineller Lernalgorithmus etwa in der Lage, auf einem Foto eine Katze zu erkennen, obwohl er nie Definitionen wie Schnurrhaare, Fell oder Pfoten einprogrammiert bekommen hat.

Trotz der Vermenschlichung solcher Systeme mit Begriffen wie lernen, trainieren oder denken basieren Maschinen lediglich auf mathematischen Modellen. Um menschliche Fähigkeiten zu imitieren, ließen sich Forscher von der Hirnforschung inspirieren. Das Gehirn besteht aus Nervengewebe, das wiederum aus Nervenzellen besteht, den sogenannten Neuronen. Über Synapsen sind die Neuronen miteinander verknüpft und können so Signale austauschen. Gemeinsam bilden sie ein neuronales Netz: Nicht etwa ein einzelnes Neuron verarbeitet eine Information, sondern das gesamte Netzwerk beziehungsweise Bereiche davon. Künstliche neuronale Netze sollen dieses Prinzip imitieren. Sie bestehen aus Knotenpunkten, den künstlichen Neuronen. Die Verbindungen zwischen ihnen, die „Synapsen“, sind gewichtet – manche Neuronen koppeln stärker aneinander als andere.

## Erste neuronale Netze und der KI-Winter

Die Entwicklung künstlicher neuronaler Netze begann jedoch nicht mit den Durchbrüchen der Nobelpreisträger in den 1980er-Jahren, sondern bereits vierzig Jahre zuvor [1]. In den 40er-Jahren etablierten sich die ersten elektronischen Computer, vor allem im militärischen und wissenschaftlichen Bereich, um aufwendige und mühselige Berechnungen durchzuführen. Zur gleichen Zeit befassten sich Forschende mit der Frage, wie das Gehirn Informationen verarbeitet. 1943 modellierten der Neurowissenschaftler Warren McCulloch und der Logiker Walter Pitts, wie Neuronen im Gehirn zusammenarbeiten: Ein Neuron berechne die gewichtete Summe eingehender Signale von anderen Neuronen, um ein ausgehendes Signal zu berechnen. Wenig später bemerkten sie, dass ein künstliches neuronales Netz zum Beispiel zur räumlichen Mustererkennung eingesetzt werden könne. Sechs Jahre später schlug der Psychologe Donald Hebb einen Mechanismus vor, wie das Gehirn lernt und Erinnerungen verarbeitet: Werden zwei Neuronen gleichzeitig und wiederholt aktiviert, verstärkt sich die Synapse, die die zwei verbindet. Diese sogenannte Hebb'sche Lernregel beschreibt das grundsätzliche Lernverfahren der meisten künstlichen Netzwerke.

In den 50er- und frühen 60er-Jahren entwickelten unterschiedliche Forschungsteams die ersten rudimentären Computer auf Basis neuronaler Netze und setzten sie zur Mustererkennung ein. 1957 präsentierten Frank Rosenblatt und Charles Wightman zum Beispiel den ersten Neurocomputer, der einfache Ziffern erkennen konnte. Marvin Minsky und Seymour Papert zeigten 1969 jedoch, dass dieses System schnell an seine Grenzen geraten würde: Nicht-lineare Probleme, wie etwa die XOR-Operation (exklusives Oder: entweder A oder B, aber nicht beides), könne die Rosenblatt-Maschine nicht lösen. Es machte sich die Sorge breit, dass künstliche neuronale Netze niemals wirklich nützlich sein könnten. Dies löste den ersten KI-Winter aus: eine Phase, in der die Forschungsförderung für künstliche neuronale Netze nahezu zum Stillstand kam.

## Das Hopfield-Netzwerk

1982 stellte Nobelpreisträger John Hopfield den Mechanismus für ein neuronales Netz vor. Damit gab er den Startschuss für ein erneutes breites Interesse an der KI-Forschung. Dafür nutzte er Metho-

den der Physik. „Ich sehe die Welt immer durch die Brille der Physik. Ich bin mit Physik aufgewachsen und die Art, wie ich Probleme aufgreife, ist sehr charakteristisch für die Physik“, sagte Hopfield in einem Interview mit dem Podcaster Lex Fridman. Die Grundidee des Hopfield-Netzwerks: Wenn Neuronen zusammenwirken, zeigen sie Eigenschaften, die man den einzelnen Neuronen nicht ansieht. Genauso verhält es sich mit Vielteilchensystemen in der Physik. Ein einzelnes Atom verrät nichts über die magnetischen Eigenschaften eines Werkstoffs. Diese ergeben sich erst aus den kollektiven Eigenschaften der unzähligen Atome im Material.

Verantwortlich für den Magnetismus eines Körpers ist der Spin der Elektronen in der Atomhülle. Der Spin ist eine quantenmechanische Eigenschaft, die man sich vereinfacht als Drehung der Elektronen um die eigene Achse vorstellen kann. Aufgrund des Spins verhalten sich Elektronen wie winzige Magnete, die in verschiedene Richtungen zeigen können. Doch erst aus der Ausrichtung all dieser Elementarmagnete und ihrer Wechselwirkung untereinander lässt sich schlussfolgern, ob das Material im Ganzen magnetisch ist oder nicht.

Hopfield betrachtete ein Netz aus mehreren Spins, die in entgegengesetzte Richtungen zeigen können, dargestellt durch die Werte +1 oder -1. Welche Ausrichtung ein einzelner Spin hat, hängt von den anderen Spins in seiner Umgebung ab. In einem

Ferromagneten, wie etwa Eisen, wollen die Spins alle in die gleiche Richtung zeigen. In Anti-Ferromagneten, die zum Beispiel in Leseköpfen magnetischer Festplatten vorkommen, wechseln sich die Spins ab. Mischformen bezeichnen Physiker als Spingläser. Die Energie des kollektiven Spinsystems ergibt sich aus der Summe der Spinwerte ( $s$ ) und ihrer Kopplungen ( $J$ ),

$$E = -\frac{1}{2} \sum_{i,j} J_{ij} s_i s_j$$

wobei die  $i$ 's und  $j$ 's die Spins durchnummerieren. Dies ist das sogenannte Ising-Modell magnetischer Materialien.

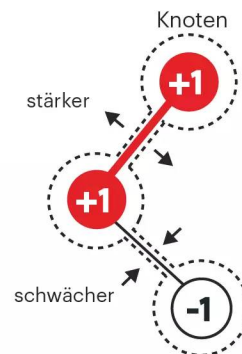
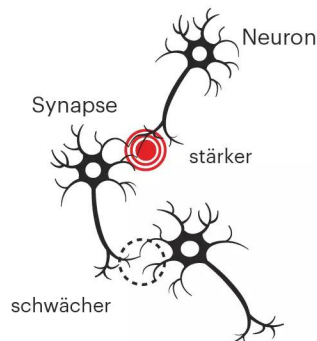
Hopfield nutzte das, um ein künstliches neuronales Netz zu beschreiben, das nach ihm benannte Hopfield-Netzwerk. Die Spins des Ising-Modells entsprechen dabei gekoppelten Neuronen, die ebenfalls die Werte +1 und -1 annehmen können. Die Energie des künstlichen, neuronalen Netzes lautet:

$$E = -\frac{1}{2} \sum_{i,j} w_{ij} s_i s_j$$

wobei  $w_{ij}$  das Gewicht der Verbindung (Synapse) zwischen den Neuronen  $i$  und  $j$  ist. Diese Formel entspricht eins zu eins der Energieformel des Ising-Modells.

## Natürliche und künstliche Neuronen

Das Netzwerk des Gehirns besteht aus Neuronen, die über Synapsen verknüpft sind (links). Künstliche neuronale Netze bestehen aus Knoten, denen ein Wert zugeordnet ist (rechts). In beiden Fällen werden beim Lernen manche Verbindungen stärker und andere schwächer.





## Zur Person: John Hopfield

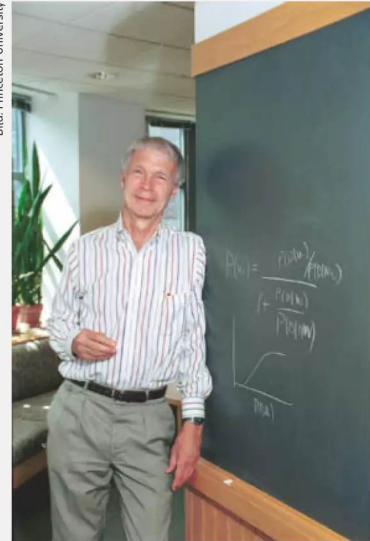
Der US-amerikanische Physiker, Molekularbiologe und Neurowissenschaftler John Joseph Hopfield wurde 1933 in Chicago geboren. Er promovierte 1958 im Bereich der theoretischen Festkörperphysik an der Cornell University. In den folgenden zwei Jahren arbeitete er in den Bell Laboratories, bevor er 1961 als Physikprofessor an die University of California in Berkeley wechselte.

Während seiner Zeit an der Princeton University von 1964 bis 1980 wandte er sich der Biochemie zu und erforschte theoretische Probleme der Molekularbiologie. Seine preisgekrönte Arbeit zum Hopfield-Netzwerk veröffentlichte der Forscher 1982 während seiner Zeit als Professor für Chemie und Biologie am California Institute of Technology (Caltech). Dort standen ihm erstmals ausreichende Mengen an Rechenleistung zur Verfügung, um seiner Forschung zu neuronalen Netzen nachzugehen.

1997 kehrte Hopfield als Professor für Molekularbiologie nach Princeton zurück, wo er bis zu

seiner Emeritierung blieb. Er erhielt während seiner Karriere diverse Auszeichnungen, darunter die Dirac-Medaille (2001), den Albert Einstein World Award of Science (2005) und die Boltzmann-Medaille (2022).

Bild: Princeton University



## Wie sich Netzwerke erinnern

In der Physik ist die Energie eine grundlegende Größe. In einem künstlichen neuronalen Netz ist sie hingegen eher abstrakt, aber essenziell, um zu verstehen, wie ein Hopfield-Netzwerk funktioniert. Während des Trainings wird das Netzwerk mit Daten gefüttert. Ein Beispiel ist ein schwarz-weißes Pixelbild des Buchstabens J. Jedes Neuron des Netzwerks entspricht dann einem Bildpunkt und es kann die Werte +1 (schwarz) oder -1 (weiß) annehmen. Die Hebb'sche Lernregel gibt dann die Stärke der Verbindung vor („Neuronen, die gemeinsam feuern, koppeln stark.“). Üblicherweise speichert ein Netzwerk jedoch nicht nur ein Muster, sondern viele gleichzeitig, zum Beispiel das gesamte Alphabet. Die Gewichte werden dann so gewählt, dass die Energie so klein wie möglich ist, wenn die Neuronen

ein gespeichertes Muster abbilden. Die Energieverteilung kann man sich als hügelige Landschaft vorstellen: Stellen die Neuronen eines der gespeicherten Muster dar, ist die Energie minimal und wir befinden uns in einem Tal der Landschaft. Bilden die Neuronen ein anderes Muster ab, ist die Energie höher; wir stehen auf einem Hügel.

So kann das Netzwerk auch anhand unvollständiger oder verzerrter Informationen das gespeicherte Muster identifizieren, das der Eingabe am meisten ähnelt. Die neue Information kann man sich wie einen Ball vorstellen, der in die hügelige Energielandschaft gesetzt wird. Die Werte der Neuronen werden anschließend neu berechnet, dabei rollt der Ball zum tiefsten Punkt der Landschaft. Dieser entspricht dem gespeicherten Muster, das der verzerrten Information am meisten ähnelt. Das verrauschte Bild wurde wiederhergestellt.

Das Hopfield-Netzwerk wird auch als Assoziativspeicher bezeichnet. Dies geht auf das assoziative Gedächtnis zurück, das augenscheinlich unabhängige Inhalte miteinander verknüpft, wie etwa das Gesicht und den Namen einer Person. Wenn Menschen versuchen, sich an eine Sache zu erinnern, durchforsten sie ihr Gedächtnis nach Informationen, die im Geiste miteinander verknüpft sind. Nach dem gleichen Prinzip findet das Hopfield-Netzwerk das gespeicherte Muster, das der eingegebenen Information am meisten ähnelt.

Später hat Hopfield sein Netzwerk verallgemeinert, statt binären Neuronen (+1 und -1) erlaubte er beliebige Werte, um so zum Beispiel den Helligkeitswert eines Bildpixels zu repräsentieren. Außerdem nutzte er das Netzwerk auch zum Lösen komplizierter Optimierungsprobleme. Moderne Quantenannealer – Alternativen zum Quantencomputer – stützen sich zum Beispiel auf ein sehr ähnliches Prinzip [2].

## Die Boltzmann-Maschine

Der zweite Nobelpreisträger, Geoffrey Hinton, wollte Informationen hingegen nicht nur speichern, sondern sie auch interpretieren oder völlig neue Muster erzeugen. Gemeinsam mit dem Neurowissenschaftler Terrence Sejnowski – der unter John Hopfield in Princeton promoviert hat – nutzte Hinton 1985 Methoden aus der statistischen Physik, um dieses Problem zu lösen. Dafür erweiterte er das Hopfield-

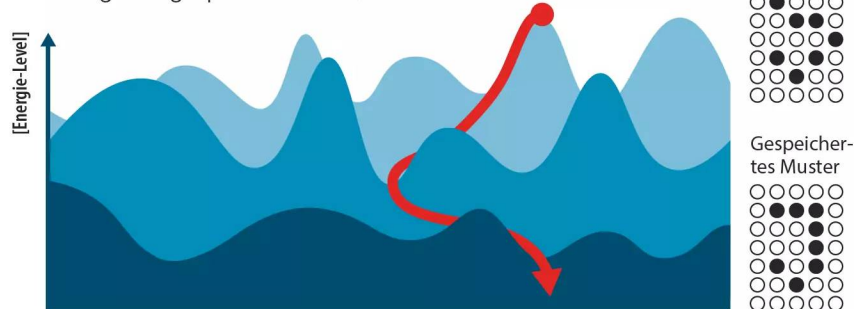
Netzwerk und machte es dadurch mächtiger, jedoch auch komplizierter.

Kurz nach der Bekanntgabe des Nobelpreises fragte ein Journalist der New York Times Hinton, ob er seine Forschung in einfachen Worten erklären könne. Dieser erwiderte: „Ich erinnere mich an die Worte des Physikers Richard Feynman, als er den Nobelpreis erhielt. Ein Journalist fragte ihn: ‚Professor Feynman, können Sie in ein paar Minuten erklären, wofür Sie den Nobelpreis erhalten haben?‘ Feynman soll geantwortet haben: ‚Hören Sie, Kumpel, wenn ich es in ein paar Minuten erklären könnte, wäre es keinen Nobelpreis wert.‘“ Wir versuchen es trotzdem.

Die statistische Physik befasst sich mit großen Systemen, die aus sehr vielen, gleichartigen Teilchen zusammengesetzt sind, etwa Gase oder Flüssigkeiten. Es ist zwar unmöglich, alle Teilchen in einem Gas einzeln zu betrachten, dafür lässt sich das Gas aber als Ganzes beschreiben. Die Eigenschaften der einzelnen Teilchen, also zum Beispiel die Positionen und Geschwindigkeiten aller Atome in einem Gas, beschreibt der sogenannte Mikrozustand. Der Makrozustand hingegen beschreibt die über alle Teilchen gemittelten Eigenschaften, also zum Beispiel die Temperatur oder den Druck des Gases. In der Regel können mehrere Mikrozustände im gleichen Makrozustand resultieren: Es gibt viele verschiedene Möglichkeiten, die Atome in einer Gasflasche anzuordnen und trotzdem 36° Celsius und 1 Bar zu messen.

## Landschaft für Erinnerungen

Der Assoziativspeicher von John Hopfield speichert Informationen in einer hügeligen Landschaft. Wird das Netzwerk trainiert, formt es in der virtuellen Energielandschaft für jedes gespeicherte Muster ein Tal. Bei Eingabe eines verzerrten Musters gibt das Netzwerk das nächstliegende gespeicherte Muster zurück.



## Zur Person: Geoffrey Hinton



Bild: Noah Berger/AP/opa

Geoffrey Everest Hinton ist Informatiker und Kognitionspsychologe und wurde 1947 im Londoner Stadtteil Wimbledon geboren. Er ist ein Ururenkel des englischen Logikers George Boole, Begründer der Booleschen Algebra. Hinton besuchte die University of Cambridge und wechselte sein Fachgebiet mehrfach, von Naturwissenschaften über Kunstgeschichte zu Philosophie, bevor er 1970 einen Bachelor in Experimentalpsychologie erhielt. Anschließend

promovierte er 1978 an der University of Edinburgh im Bereich der künstlichen Intelligenz.

Zunächst forschte Hinton an der University of Sussex in England, der University of California in San Diego und der Carnegie Mellon University in Pittsburgh, beide USA. 1986, während seiner Zeit als Postdoc in San Diego, veröffentlichte er sein einflussreiches Forschungspapier zur Backpropagation (Fehlerrückführung) beim Training neuronaler Netze. Ein Jahr später wurde er Professor für Computer Science an der kanadischen University of Toronto. Dort ist er seit 2014 emeritierter Professor.

Für seine KI-Forschung erhielt er diverse Auszeichnungen, darunter den Turing-Award (2018, gemeinsam mit Yoshua Bengio und Yann LeCun), die höchste Auszeichnung in der Informatik. Er trägt den Spitznamen „Godfather of AI“, wortwörtlich „Pate der KI“. Zwischen 2013 und 2023 arbeitete Hinton außerdem für Googles damalige Deep-Learning-Abteilung Google Brain. 2023 verließ er Google jedoch aufgrund seiner Bedenken zur KI-Forschung.

Wenn sich der Makrozustand des Systems nicht ändert, befindet es sich im sogenannten thermischen Gleichgewicht. Trotzdem sind verschiedene Mikrozustände mit unterschiedlicher Energie möglich. Wie wahrscheinlich eine bestimmte Energie ist, lässt sich mithilfe der Boltzmann-Verteilung berechnen,

$$p = \frac{1}{Z} e^{-E/k_B T}.$$

Die Wahrscheinlichkeit ( $p$ ) nimmt exponentiell mit der Energie ( $E$ ) ab und hängt von der Temperatur ( $T$ ) und einer Naturkonstante ( $k_B$ , die sogenannte Boltzmann-Konstante) ab. Der Faktor  $Z$  garantiert, dass sich alle Wahrscheinlichkeiten zu 1 aufsummieren. Benannt wurde sie nach dem österreichischen Physiker Ludwig Boltzmann. Hinton nutzte diese Gleichung, um ein künstliches neuronales Netzwerk zu beschreiben. Während Hopfield die Werte der Neuronen in seinem Netzwerk deterministisch, also anhand einer fixen Regel, aktualisierte,

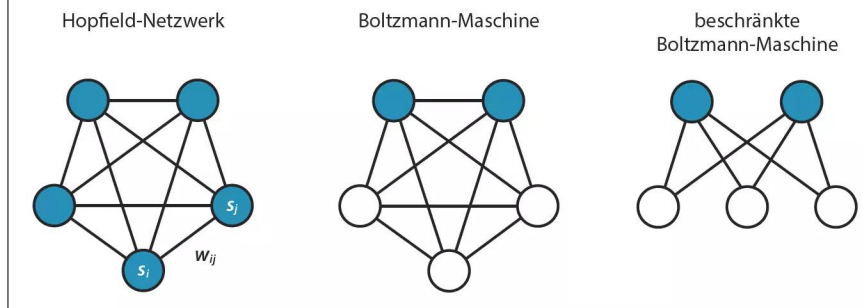
arbeitet die Boltzmann-Maschine probabilistisch: Die Boltzmann-Verteilung legt fest, wie wahrscheinlich verschiedene Werte eines Neurons sind. Eine Zufallszahl bestimmt, welchen Wert das Neuron tatsächlich annimmt.

Hinton trainierte die Boltzmann-Maschine mit unverfälschten sowie verrauschten Daten, damit sie erwünschte von unerwünschten Mustern unterscheiden kann. Um dies zu erreichen, bestimmt ein Algorithmus die Werte der Neuronen und Verbindungen Schritt für Schritt. Zudem besitzt die Boltzmann-Maschine außer sichtbaren Neuronen, die die Eingabe und Ausgabe kodieren, auch versteckte Neuronen. Diese können abstrakte Eigenschaften und Korrelationen in den Daten speichern und verarbeiten, zum Beispiel Details, die dabei helfen, ein  $I$  von einem  $J$  zu unterscheiden. Damit ist das Netzwerk auch in der Lage, neue Muster zu erzeugen, die zu den Trainingsdaten passen. Die Boltzmann-Maschine war damit eines der ersten generativen Modelle.



## Verschiedene Arten von Netzwerken

Im Hopfield-Netzwerk (links) sind alle Knoten miteinander verbunden. Alle Knoten dienen zur Speicherung und Ausgabe der Daten, sind also sichtbar. Die Boltzmann-Maschine (Mitte) besteht aus sichtbaren Knoten (blau), die Informationen speichern und ausgeben, und versteckten Knoten (weiß), die zusätzliche Rechnungen ausführen. In der beschränkten Boltzmann-Maschine (rechts) besteht keine Verbindung zwischen Knoten einer Ebene.



Später stellte sich heraus, dass die Boltzmann-Maschine sehr ineffizient arbeitet, weil ein Algorithmus die Werte aller Neuronen nacheinander aktualisieren muss. Schließlich setzte sich mit der beschränkten Boltzmann-Maschine eine Variante durch, bei der sichtbare nur mit versteckten Neuronen verbunden sind – zwischen Neuronen einer Schicht besteht hingegen keine Verbindung. Obwohl diese Variante weniger Freiheiten bietet, arbeitet sie in der Praxis effizienter, da alle Neuronen einer Ebene gleichzeitig aktualisiert werden können.

## Fehlerrückführung

Wichtiger als die Boltzmann-Maschine sei für heutige KI-Modelle aber eine andere Technik, sagte Hinton der New York Times: die Backpropagation. „Das hat weniger mit Physik zu tun.“ Sowohl die Boltzmann-Maschine als auch die Backpropagation hätten allerdings geholfen, neuronale Netze zum Laufen zu bekommen. Die Idee der Backpropagation, oder Fehlerrückführung, ist die folgende [3]: Das neuronale Netz erhält ein Eingabemuster und reicht es zu den Ausgabeneuronen durch. Es vergleicht die tatsächliche Ausgabe mit der gewünschten und betrachtet die Abweichung als Fehler. Anschließend reicht es den Fehler von den Ausgabeneuronen zurück an die Eingabeneuronen. Dabei passt es die Gewichte der Neuronenverbindungen so an, dass

sich der Fehler verkleinert. In der nächsten Wiederholung ähnelt das Ausgabemuster dem gewünschten Muster stärker. Dieser Prozess wird so lange wiederholt, bis das Netzwerk die gewünschte Genauigkeit liefert – es ist trainiert.

Die Methode geht auf diverse Forscher der 1960er-Jahre zurück. Hinton zeigte jedoch 1986 erstmals gemeinsam mit David Rumelhart und Ronald Williams den Nutzen dieser Methode für tiefe neuronale Netze. Heutzutage gebe es laut Hinton keine große Verbindung mehr zwischen Boltzmann-Maschinen und der Backpropagation. „In den Anfängen gelang es mir, sie zu kombinieren, indem ich Boltzmann-Maschinen zum Vortraining von Backpropagation-Netzen verwendete. Aber das macht man heute nicht mehr.“

Wie so häufig kam auch 2024 die Frage auf, ob die richtigen Personen den Nobelpreis erhalten hätten. Forschung ist eine Gemeinschaftsleistung und lediglich zwei oder drei Personen für ein Forschungsgebiet auszuzeichnen, ist schwierig. „Hopfield und Hinton haben den Nobelpreis aus meiner Sicht absolut verdient“, meint Ute Schmid, Leiterin des Lehrstuhls für Kognitive Systeme an der Otto-Friedrich-Universität Bamberg. „Beide haben wesentliche methodische Grundlagen für das maschinelle Lernen mit neuronalen Netzen geschaffen.“ Doch auch ihr fallen viele weitere Namen ein, die das Forschungsgebiet geprägt haben, darunter Hintons Kollege

David Rumelhart, der maßgeblich an der Einführung der Backpropagation beteiligt gewesen war.

## Preisträger beunruhigt von Entwicklung von KI

Auf die konzeptionellen Durchbrüche in den 80er-Jahren folgten die ersten erfolgreichen Anwendungen künstlicher neuronaler Netze, etwa zur automatischen Erkennung von Unterschriften bei Banken oder zur Analyse klinischer Daten. Mit der Veröffentlichung von ChatGPT 2022 setzte schließlich der Boom generativer künstlicher Intelligenz ein. Die Größe der Netzwerke ist in dieser Zeit explodiert. Hopfields Netzwerk von 1982 nutzte 30 Neuronen mit 435 Verbindungen und somit etwas weniger als 500 Parametern. Bereits ein Netzwerk mit 100 Neuronen überstieg die Rechenkapazitäten, die er damals zur Verfügung hatte. Vereinzelt große Sprachmodelle verwenden heute mehr als eine Billion Parameter.

Künstliche Intelligenz werde einen riesigen Einfluss auf die Menschheit haben, sagte Hinton im Anschluss an die Verkündung des Nobelpreises im Telefonat mit der Königlich Schwedischen Akademie der Wissenschaften. „Sie wird mit der industriellen Revolution vergleichbar sein, aber anstatt die Menschen an körperlicher Stärke zu übertreffen, wird sie die Menschen an intellektuellen Fähigkeiten übertreffen.“ Gleichzeitig warnt er vor dem Einsatz von KI. „Wir haben keine Erfahrung damit, wie es ist, wenn Dinge intelligenter sind als wir“, sagte er weiter. Es könne zu einem effizienteren Gesundheitswesen und gesteigerter Produktivität führen. „Wir müssen uns aber auch über eine Reihe möglicher negativer Folgen Sorgen machen. Besonders über die Gefahr, dass diese Dinge außer Kontrolle geraten.“

2023 kündigte Hinton seine Stelle bei Google Brain, um frei über die Risiken von KI sprechen zu können. Im Mai 2024 äußerte er etwa im Wissenschaftsmagazin Science gemeinsam mit anderen Experten unter dem Titel „Umgang mit extremen KI-Risiken inmitten des rasanten Fortschritts“ seine Bedenken. Sie schreiben dort: „Ohne ausreichende Vorsicht könnten wir unwiderruflich die Kontrolle über autonome KI-Systeme verlieren, so dass ein menschliches Eingreifen wirkungslos wird. Dieser unkontrollierte KI-Fortschritt könnte (...) zur Auslöschung der Menschheit führen.“

Auch Hopfield zeigte sich kritisch. Im Jahr 2023 unterzeichnete er einen offenen Brief mit der Forderung, das Training von KIs zu pausieren, die leis-

tungsfähiger sind als GPT-4. Kurz nach der Bekanntgabe des Preises sagte er während einer Pressekonferenz der Princeton University: „Als Physiker bin ich sehr beunruhigt von etwas, das niemand kontrolliert, etwas, das ich nicht gut genug verstehe, um zu sehen, wo die Grenzen liegen, an die man diese Technologie treiben könnte.“ Er plädierte daher dafür, dass das Forschungsfeld ein besseres Verständnis der KI anstreben müsse.

## Der Physik zum Nutzen

Während KI in der öffentlichen Wahrnehmung erst seit einigen Jahren Aufmerksamkeit erregt, verwenden Forscher das maschinelle Lernen seit bereits 15 bis 20 Jahren. Insbesondere fußen künstliche neuronale Netze nicht nur auf Methoden der Physik, sondern maschinelles Lernen leistete auch einen erheblichen Beitrag zur modernen physikalischen Forschung. „Viele Erkenntnisse in der Physik sind nur durch maschinelles Lernen zu gewinnen“, sagt Katharina Morik, Gründungsdirektorin des Lamarr-Instituts für Maschinelles Lernen und KI in Nordrhein-Westfalen. Dazu gehört etwa das IceCube-Projekt, an dem Morik beteiligt war. Mithilfe eines Detektors am Südpol und der KI-gestützten Auswertung der darüber gesammelten Daten konnte das Team die Milchstraße anhand von Neutrinos abbilden.

In den 90er-Jahren wurden künstliche neuronale Netze zu einem Standard-Analysewerkzeug in der Teilchenphysik. Teilchenbeschleuniger wie etwa der Large Hadron Collider (LHC) am CERN produzieren Unmengen an Daten – im Schnitt verarbeitet das CERN-Datenzentrum ein Petabyte pro Tag. Grund dafür ist, dass die dort gesuchten Elementarteilchen oft nur für Bruchteile einer Sekunde kurz nach einer energiereichen Kollision existieren. Um diese trotzdem aufzuspüren, müssen Forscher gigantische Datenmengen aufzeichnen und nach den richtigen Spuren durchsuchen. So gelang es dem Forschungsteam 2012, das Higgs-Boson nachzuweisen, für das Peter Higgs und François Englert im folgenden Jahr mit dem Nobelpreis für Physik ausgezeichnet wurden.

Dass der Nobelpreis 2024 an ein informatisches Thema ging, hält Schmid für gerechtfertigt. Einen Nobelpreis für Informatik gebe es nun einmal nicht. „Ein Nobelpreis symbolisiert deutlicher, dass eine Forschungsarbeit nicht nur von hoher fachlicher Relevanz ist, sondern auch darüber hinaus Bedeutung für Wissenschaft und Gesellschaft hat.“ (spa) **ct**

### Literatur

[1] Jan Mahn, Die Säulen der Moderne, Durchbrüche, Wegbereiter und Strömungen der KI-Forschung, c't 11/2024, S. 46

[2] Alexander Braun, Heißes Eisen? Quanten-Annealer – Vorbote des Quantencomputers, c't 13/2019, S. 142

[3] Andrea Trinkwalder, Des Netzes Kern, Denkbar einfach: Die Mathematik neuronaler Netze, c't 11/2024, S. 22



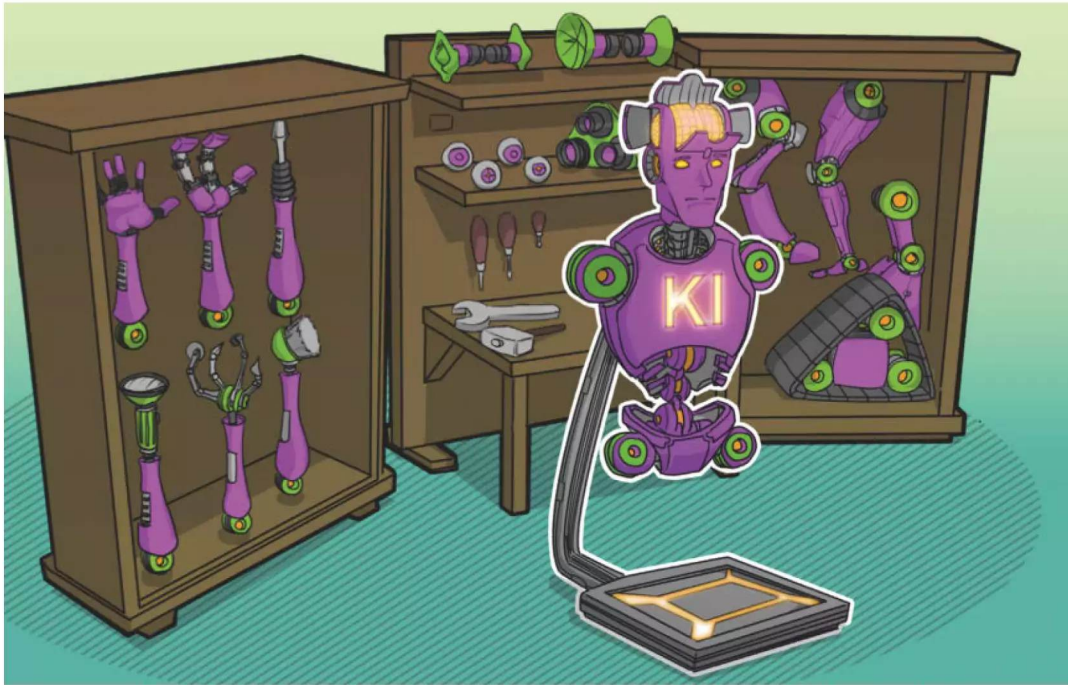


Bild: Thorsten Hübner

# Servicezentren für KI im Mittelstand

In Deutschland ebnen vier KI-Servicezentren des Bundes den Weg zu individuellen Anwendungen. Sie beraten und helfen dabei, KI-Modelle zu entwickeln, sie hosten Open-Source-Tools und Sprachgeneratoren, helfen beim Modelltuning und bieten Rechenpower für maschinelles Lernen.

Von **Arne Grävemeyer**

**K**ünstliche Intelligenz kann in verschiedenen Unternehmen ganz unterschiedliche Anwendungen ermöglichen. Für Behörden kann ein KI-Chatbot Amtssprache in verständliches Deutsch übersetzen oder freie Beschreibungstexte nach rechtlich definierten Sachverhalten klassifizieren. In einem Verlag gibt eine künstlich intelligente

Suche Überblick über das verfügbare Wissensspektrum und in einer Klinik hilft KI, angesichts vielfältiger Patientendaten bei der Diagnose nicht den Überblick zu verlieren.

Diese Beispiele zeigen, wie vielfältig ganz unterschiedliche Institutionen selbstlernende Cloud-gestützte Systeme einsetzen. Aber für viele Unter-



nehmen sind die eigenen Datenbestände das Tafelsilber, das sie nicht einfach für Abfragen oder Deep-Learning-Prozesse an Dienstleister nach Übersee übertragen können. Sogar noch sensibler ist der Umgang mit medizinischen Patientendaten. Der Einsatz von Chatbots auf US-Servern wie etwa ChatGPT von Open AI oder PaLM von Google scheidet aus diesem Grund oft aus.

Das Ziel der nationalen KI-Strategie Deutschlands besteht darin, einschlägige Kompetenz an Hochschulen, in Behörden und gerade auch in kleinen und mittelständischen Unternehmen zu stärken. Ein Teil dieser Strategie sind vier KI-Servicezentren, angesiedelt in Nordrhein-Westfalen (WestAI), in Berlin-Brandenburg (KISZ-BB), in der Region Hannover-Göttingen-Kassel (KISSKI) sowie in Darmstadt (hessian.AI). Diese Zentren bieten Interessenten ein breites Spektrum von Beratungsdienstleistungen und Erfahrung bei der Konzeption von KI-Modellen an.

## Vier Zentren, jedes mit seinem Fokus

Jede Behörde und jeder an künstlicher Intelligenz interessierte Mittelständler kann sich örtlich unabhängig an eines der vier KI-Servicezentren wenden. Dabei sind nicht alle Stellen gleich. WestAI mit einem Partnernetzwerk in fünf nordrhein-westfälischen Städten legt einen Fokus auf multimodale große Modelle, hessian.AI in Darmstadt hat Schwerpunkte bei Erklärbarkeit, Generalisierbarkeit und kontextuellen Anpassungen. Das KI-Servicezentrum Berlin-Brandenburg zielt besonders auf Training und Wissenstransfer ab, während KISSKI (KI-Servicezentrum für sensible und kritische Infrastruktur) an den Standorten Hannover, Göttingen und Kassel einen Fokus auf den Umgang mit sensiblen und kritischen Daten legt, zusätzlich zu Beratung, Umsetzung, Training und Infrastruktur.

Alle KI-Servicezentren helfen, Projekte umzusetzen, bieten Rechenpower für Deep-Learning-Prozesse, schulen Mitarbeiter und hosten rechenintensive Standardmodelle beim KISSKI-Partner GWDG (Gesellschaft für wissenschaftliche Datenverarbeitung). Dabei handelt es sich zum Beispiel um unterschiedlich große Sprachmodelle (Large Language Models, LLM) wie Llama-3-Versionen von Meta oder Mistral-Versionen vom französischen Start-up Mistral AI sowie von der chinesischen Alibaba zwei Versionen der Qwen-KI-Modelle, die bereits Stärken in der Softwareentwicklung zeigten. Die GWDG führt sowohl das Hochschulrechenzentrum der Universität Göt-

tingen als auch das dortige Rechen- und IT-Kompetenzzentrum der deutschen Max-Planck-Gesellschaft und stellt im KISSKI die benötigten Rechenkapazitäten zur Verfügung.

Alle vier Zentren sind miteinander vernetzt und weisen gegebenenfalls auf die besonderen Kompetenzen der jeweils anderen hin. „Sie können also jedes der Servicezentren nach Hilfe fragen und Ihnen wird geholfen“, sagt KISSKI-Projektleiter Julian Kunkel, Professor für Hochleistungsrechnen an der Uni Göttingen und stellvertretender Leiter der GWDG.

## Verlags-KI klärt Kundenfragen

Ein KI-Anwender, der mit KISSKI zusammenarbeitet, ist der Hogrefe-Verlag in Göttingen. „Das Problem eines Verlags: Es dürfen keine Inhalte abfließen, denn unser Content ist unser Kapital“, sagt Digital Project Manager Martin Rieken. Hogrefe ist ein Wissenschaftsverlag in Psychologie, Psychotherapie und Psychiatrie, der zusätzlich zu Zeitschriften und Büchern insbesondere Testverfahren publiziert. Etwa 2300 psychologische Testverfahren decken ein breites Spektrum ab, sie helfen bei der Personalauswahl oder bei der diagnostischen Einschätzung von Kindern mit ADHS- oder Autismusverdacht. Etwa 700 dieser Verfahren bietet der Verlag mit digitalen Diensten im Hogrefe-Testsystem (HTS) an. Er hostet sie in einer internen Cloud und Kunden binden sie beispielsweise in ihre Krankenhaus-Informationssysteme (KIS) ein. „In dem Fall haben wir sogar Patientendaten auf unseren Servern“, sagt Rieken. Das ist ein weiteres Argument dagegen, die Daten für die Verarbeitung durch eine KI aus der Hand zu geben.

Eine gemeinsam mit KISSKI entwickelte KI-Anwendung soll nun helfen, den Überblick über die Testverfahren und das weitere Verlagsprogramm zu behalten. Es geht um einen Chatbot, der nicht nur Mitarbeiter im Kundendienst und im Vertrieb unterstützt. Damit können diese spezielle Anfragen zu Testverfahren etwa aus dem Klinikalltag wesentlich schneller beantworten als zuvor, der Rechercheaufwand sinkt. Zur Kontrolle verlinkt die KI ihre Antwort mit einer Handvoll PDFs, in denen der Anwender beispielsweise die Zielgruppen und die Fragenkataloge zu genannten Testverfahren nachschlagen kann.

Der interne Hogrefe-Chatbot ist seit Mai 2024 mit KISSKI-Unterstützung modelliert und umgesetzt worden. Das System stellt die Bedienoberfläche und

läuft auf internen Servern. Das Sprachmodell im Hintergrund wurde hingegen zunächst im Rechenzentrum der GWDG in Göttingen gehostet. Dort stehen derzeit unterschiedliche LLMs zur Auswahl, etwa Metas Llama 3.1 70B Instruct, Llama 3.1 Sauerkraut LM 70B Instruct, Llama 3.1 Nemotron 70B Instruct, Alibabas Qwen 2.5 72B Instruct sowie Mistral Large Instruct. Im Servicezentrum beobachtet man die Entwicklung der als Open Source zur Verfügung stehenden LLMs und plant entsprechende Updates.

Der Wechsel von einem Sprachmodell zum anderen ist für den Anwender sehr einfach, er muss nur einen API-Key ändern. „Jedes LLM erfordert eine teure Recheninfrastruktur mit etwa zwei bis vier High-End-

GPUs (Grafikprozessoren). Eine solche Investition überlegt man sich schon. Durch das Hosting des Servicezentrums ist der Schritt zu einer KI-Anwendung viel einfacher und wir können stets mit einem aktuellen LLM arbeiten“, berichtet Rieken. So habe sich auch gezeigt, dass nicht alle Sprachmodelle gleich gut die deutsche Sprache verarbeiten.

Nachdem der Chatbot auf alle deutschsprachigen Standorte ausgerollt worden ist, könnte er in Zukunft auch externen Interessenten direkt Fragen zu Testverfahren und Bestellmöglichkeiten beantworten. Zudem bewirbt Hogrefe inzwischen mit ChatLibrary eine KI-gestützte Wissenslösung, die auf Anfragen antwortet und dabei ausschließlich die

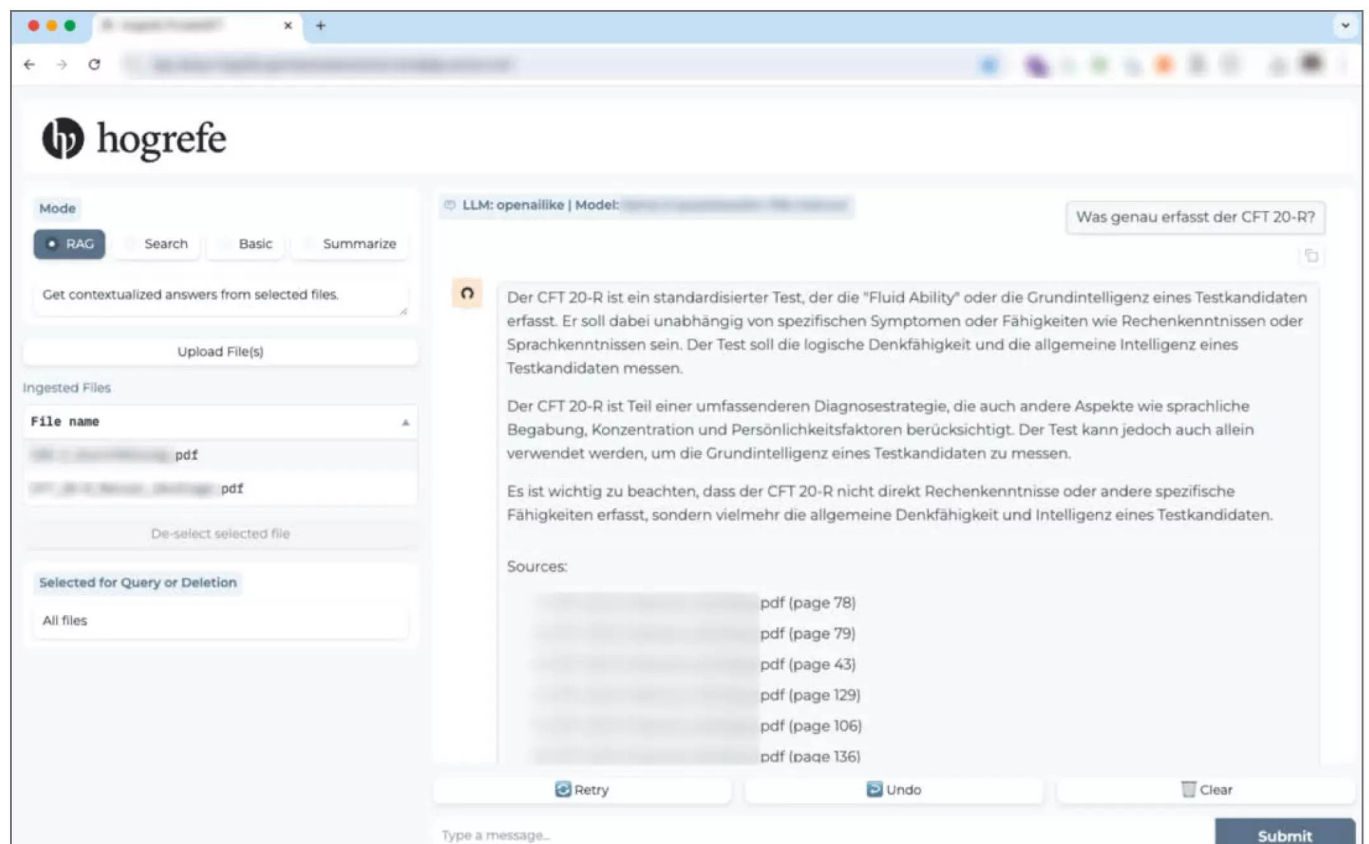
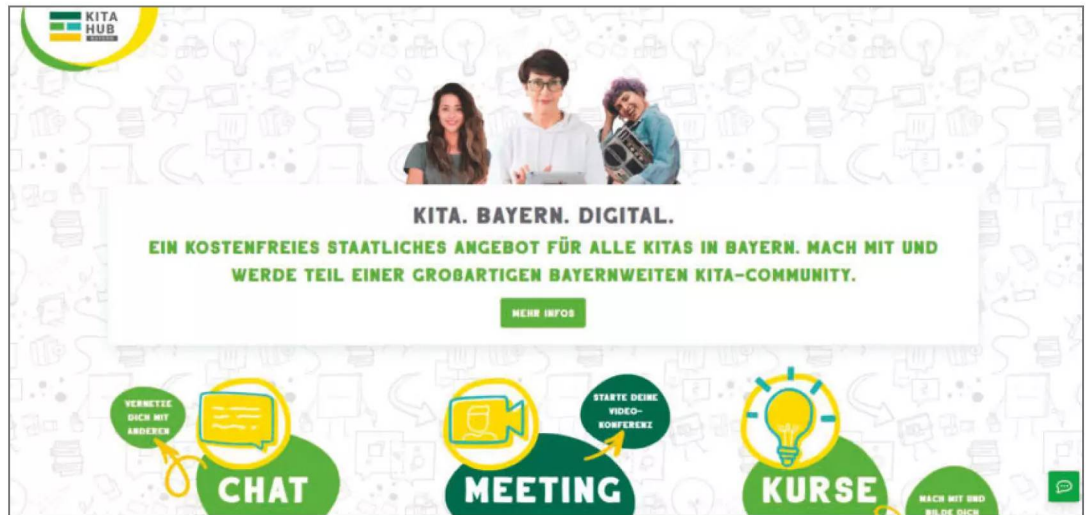


Bild: Hogrefe-Verlag

**Gemeinsam mit dem KISSKI hat der Fachverlag Hogrefe einen Chatbot entwickelt, der seinen Kundenberatern Fragen zum Verlagsprogramm beantwortet und Links zu PDFs mit weitergehenden Informationen liefert. In Zukunft könnte der Chatbot für externe Zugriffe geöffnet werden.**

Mehr und mehr KI-Anwendungen sollen künftig den Kita Hub als staatliches Angebot für die bayrischen Kitas verbessern, etwa durch einen Chatbot für Anleitungen, ein sicheres Transkriptions-tool und sogar eine Dolmetscher-App.



ihm zugrundeliegenden Quelltexte als Datenbasis berücksichtigt.

## Übersetzer für Behördendeutsch

Das bayrische Staatsinstitut für Frühpädagogik und Medienkompetenz (IFP) betreibt mit dem Kita Hub (siehe [ct.de/wgy1](https://ct.de/wgy1)) eine Plattform mit Onlinediensten für alle Kitas in Bayern. Der Kita Hub informiert zum Beispiel über Weiterbildungen für Erzieher, bietet ihnen einen datenschutzkonformen Messenger, um sich untereinander auszutauschen, und über Big-Blue-Button ein Videokonferenzsystem für Teambesprechungen, Elterngespräche sowie virtuelle Veranstaltungen. Jetzt arbeiten IFP-Mitarbeiter an einer Kita-Hub-App, die eine Vielzahl kleiner KI-Anwendungen vereinen soll. Beispielsweise sollen Mitarbeiter einen Chatbot bekommen, der ihnen zu ihren Fragen Auskünfte aus Anleitungstexten extrahiert. Für Eltern soll es wiederum einen Chatbot geben, der ihnen Fragen zu Bürgerservices in Echtzeit beantwortet und dabei gleich das Behördendeutsch der zugrundeliegenden Bestimmungen in Alltagssprache übersetzt.

„So entstehen statt einer großen KI-Anwendung viele kleinere KI-Dienste“, sagt Manfred Steger, Projektleiter für den Kita Hub. Weitere kleine Dienste sind ein Generator für urheberrechtsfreie Illustrationen, die Mitarbeiter für Flyer nutzen können, und ein Transkriptions-Tool, mit dem Mitarbeiter Ge-

sprachsaufzeichnungen verschriftlichen können. Ein daraus weiterentwickeltes Werkzeug könnte den Kindern in Zukunft sogar als Dolmetscher für Gespräche untereinander dienen.

Die letzten beiden Beispiele zeigen, dass diese KI-Dienste sensible Informationen verarbeiten können, wenn in Mitarbeitermeetings beispielsweise über einzelne Kinder gesprochen wird oder wenn Kinder mit einem KI-Dolmetscher spielen. Am IFP freut man sich daher nicht nur über „ein kompetentes KISSKI-Team, das per Chat immer erreichbar ist und auf Bedürfnisse eingeht“, wie Steger berichtet. Sehr wichtig ist auch die datenschutzkonforme Verarbeitung mit abgeschotteten LLMs im GWDG-Rechenzentrum.

## Datenwust aus MRT-Hirnschans

Am Institut für Diagnostische und Interventionelle Neuroradiologie der Universitätsmedizin Göttingen (UMG) arbeitet ein Team um Christian Riedel an einer AI-Toolbox, die Aussagen zu den umfangreichen Bilddaten eines MRT-Hirnschans liefert. „Sie müssen sich vorstellen, dass ein Mediziner heute die Schnittbilder des Gehirns betrachtet und einfach die Auffälligkeiten beschreibt. Ist eine Hirnhälfte kleiner als die andere? Ist der Hippocampus auffällig, also der Bereich, der für Erinnerungen zuständig ist? Menschen können aber Volumen schlecht einschätzen“, schildert Riedel gegenüber c't.



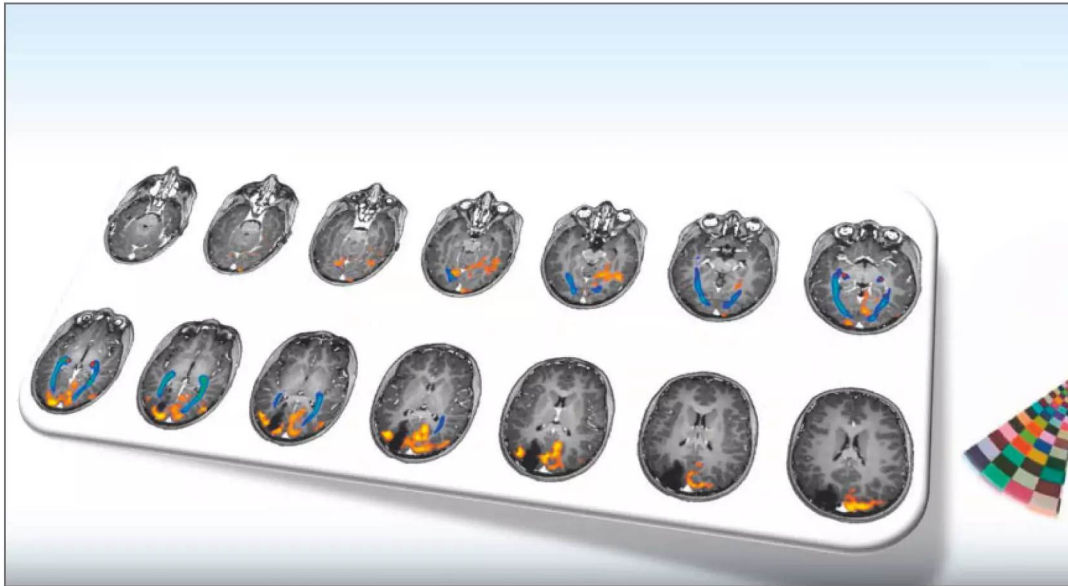


Bild: UMG

**MRT-Scans des Gehirns liefern viele Schnittbilder, aber noch keine Volumeninformationen. An der Universitätsmedizin Göttingen nimmt eine KI die Arbeit auf, die binnen Minuten die Segmente und Teilvolumen eines Hirnscans bestimmt.**

Wollte ein Mediziner eine vollständige Volumetrie des gescannten Gehirns vornehmen, samt sauberer Segmentierung aller Hirnbereiche, dann bräuchte er wahrscheinlich einen ganzen Tag und hätte eine sehr langweilige Aufgabe. Nach Vorarbeiten an der Harvard University ist inzwischen am MIT (Massachusetts Institute of Technology) unter der Bezeichnung „Fast Surfer“ ein Convolutional Neural Network entwickelt und veröffentlicht worden, das in wenigen Minuten die Volumen und vollständige Segmentierung eines MRT-gescannten Gehirns berechnet. Damit arbeiten Mediziner an der UMG bereits. Ihre Daten gehen direkt über Glasfaser in das GWDG-Rechenzentrum, das ein Fast-Surfer-Modell hostet.

„Die Ergebnisse sind schon ansprechend“, berichtet Riedel. Die KI-Anwendung soll den Mediziner die Volumetrie und Segmentierung bei jedem Hirnscan schnell bereitstellen. Dann gilt es nur noch zu kontrollieren, ob die KI beispielsweise die Grenzen eines Tumors oder Operationsnarben richtig interpretiert hat. Statt der bisherigen Beschreibungen nach Augenschein steht den Mediziner damit schnell ein quantitativer Befund mit Substrukturen und Teilvolumen zur Verfügung. In Zeitreihen können sie leicht schrumpfende Hirnstrukturen oder die Entwicklung entzündlicher Läsionen erkennen.

Auch für die Mediziner am UMG ist es entscheidend, auf eine geschützte Recheninfrastruktur zurückzugreifen. Hirnscans auf den Servern eines Dienstleisters in den USA zu verarbeiten, käme für sie nicht infrage.

Für die Zukunft hofft Riedel, bundesweite Referenzdaten einer Kohorte von 30.000 Scans verarbeiten zu können und mit der KI Reihendaten zu gewinnen. Die würden dann die durchschnittliche Entwicklung der verschiedenen Hirnregionen im Laufe des Lebens zwischen dem 20. und dem 80. Lebensjahr wiedergeben.

## Wann ist es Geldwäsche?

Die Financial Intelligence Unit des Zolls in Köln arbeitet an einem Chatbot, der zu einfachen sprachlichen Sachverhaltsbeschreibungen entscheidet, ob es dabei um Geldwäsche geht. Der Zoll unterscheidet dabei sechs eigenständige Klassen, die für Laien nicht ganz einfach zuzuordnen sind.

Das Start-up FloodWaive Predictive Intelligence entwickelt gemeinsam mit dem KISSKI Deep-Learning-Algorithmen für effiziente Hochwasservorhersagen und Risikoanalysen. Diese Modelle werden mit den Ergebnissen hydraulischer Simulationen trainiert und sollen einmal rechenintensive hydrodynamische

sche Hochwassermodelle ersetzen. Die Betreiber erhoffen sich von ihrer KI-Anwendung einen Geschwindigkeitsvorteil um den Faktor 10.000 und damit wesentlich frühere Hochwasserwarnungen als bisher.

## Sichere Chatbots sehr gefragt

Mit den genannten Beispielen ist das breite Spektrum hilfreicher KI-Anwendungen noch nicht erschöpft. „Besonders gefragt sind derzeit Chatbot-Lösungen, und zwar sowohl die Nutzung von LLMs als auch die Forschung an eigenen Sprachmodellen“, sagt KISSKI-Projektleiter Kunkel im Gespräch mit c't. Zudem herrscht hohes Interesse an sicheren Speicher- und Rechenlösungen. Die Betreuung am Servicezentrum erstreckt sich von der reinen Beratung bis zur kompletten Projektbegleitung.

In aller Regel stellt das Servicezentrum für festgelegte Zeiträume Rechenressourcen und API-Keys bereit. Üblicherweise halten die Projektpartner zumindest alle drei Monate einen Jour fixe ab und einmal jährlich ein Treffen. Mehrere Universitäten

wollen den Zugang für ihre Studenten ausrollen und davon profitieren, dass damit dann auch sensible Daten geschützt sind.

Generell sieht Kunkel den KI-Vorsprung der US-Konzerne insbesondere in Bezug auf Sprachgeneratoren kritisch. „Heute veröffentlichen diese Entwickler noch Open-Source-Modelle, aber das werden sie in Zukunft wohl nicht mehr tun. Die Technik erhöht aber die Produktivität von sehr vielen Menschen, an der kommt in Zukunft kaum noch jemand vorbei“, warnt er. Zudem nutzen viele Arbeitnehmer ChatGPT in der freien Variante und geben ihre Daten damit für das KI-Training frei. Auch das erhöhe weiter den Vorsprung der Amerikaner. Start-ups mit guten Ideen und eigenen Entwicklungen werden von den Hyper-scalern schnell aufgekauft.

Kunkel plädiert daher für eine staatlich geförderte Infrastruktur, mit der KI-Anwender hierzulande auch in Zukunft einen bezahlbaren Zugang zu KI-Modellen wie LLMs haben. Das zeitlich geförderte Projekt der vier deutschen KI-Servicezentren ist zunächst nur ein Anfang. (agr) **ct**

Kontakt zu den deutschen  
KI-Servicezentren:

[ct.de/wgy1](https://ct.de/wgy1)

# Es gibt 10 Arten von Menschen. iX-Leser und die anderen.

The advertisement features a collage of iX magazine covers. The main cover shown is the 'KI im eigenen RZ' (AI in your own data center) issue, which includes articles on 'Netzwerk, Storage, Prozessoren: Die passende Infrastruktur für Machine Learning', 'Post-Quanten-Signaturen', 'Ansible für Windows-Umgebungen', and 'Geheimhaltungsvorgaben für IT-Systeme'. Other visible covers include 'Schlauere KI' (Smarter AI) and 'Chatbots'. A yellow circular badge with the text '3 x als Heft' (3 times as a magazine) is overlaid on the covers. In the foreground, a black Bluetooth keyboard is shown. The background is a blurred image of colorful, glowing fiber optic cables.

**Jetzt Mini-Abo testen:**  
3 Hefte + Bluetooth-Tastatur nur 19,35 €  
[www.ix.de/testen](https://www.ix.de/testen)

[www.ix.de/testen](https://www.ix.de/testen)    0511 / 647 22 888    [leserservice@heise.de](mailto:leserservice@heise.de)

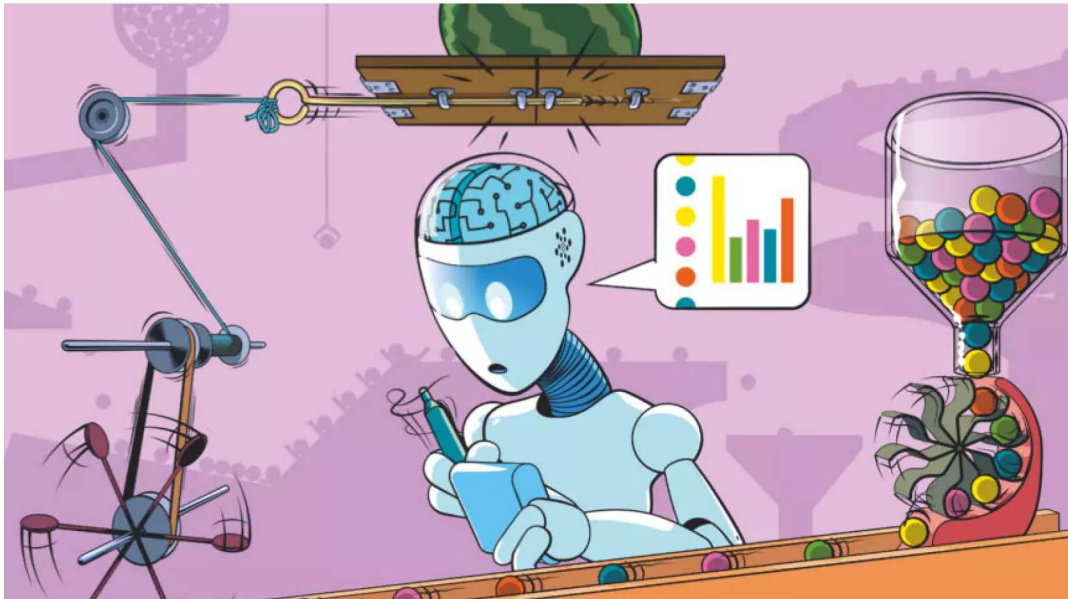


Bild: Rudolf A. Blaha

# Causal ML: KI erfasst Ursache und Wirkung

Maschinelles Lernen findet zwar zuverlässig Korrelationen bei der statistischen Auswertung, jedoch nicht ohne Weiteres kausale Zusammenhänge. Causal ML soll das ändern und der künstlichen Intelligenz damit zu neuen Höhenflügen verhelfen. Erste Anwendungen in Wirtschaft und Medizin deuten auf großes Potenzial hin.

Von **Thomas Brandstetter**

**K**ausale Zusammenhänge sind allgegenwärtig. Allerdings lassen sich Ursache und Wirkung nicht immer so leicht unterscheiden wie beim Dominoeffekt, wo ein Stein den nächsten umstößt. Gerade statistische Methoden wie das maschinelle Lernen haben fundamentale Schwierigkeiten damit, Kausalität zu erfassen und zu verarbeiten. Ihre Stärke liegt in der Erkennung von Mustern und Korrelationen. Sie finden zum Beispiel heraus, dass der Konsum von Schokolade mit der Anzahl der Nobelpreis-

träger in einem Land korreliert. Das heißt aber noch lange nicht, dass Schokolade schlau macht. Vielmehr gibt es eine dritte Variable, den Wohlstand, der die gemeinsame Ursache darstellt. Herkömmliche Machine-Learning-Verfahren ergünden solche Wechselwirkungen und Abhängigkeiten nicht. Sie kratzen oft nur an der Oberfläche der Realität – und verleiten den Menschen damit zu falschen Schlussfolgerungen.

Causal ML ist ein Überbegriff für die Verknüpfung von maschinellem Lernen mit Kausalität und hat



viele Gesichter. Das wohl ambitionierteste Ziel von Causal ML ist es, Machine-Learning-Modelle zu entwickeln, die gefundene Muster selbstständig auf Ursache und Wirkung hin analysieren – also zu echtem logischem Schlussfolgern fähig sind. Einer der Forschungs-Hotspots dafür ist die Abteilung für Empirische Inferenz am Max-Planck-Institut für Intelligente Systeme in Tübingen. Julius von Kügelgen hat dort promoviert und forscht inzwischen an der ETH Zürich. „Kausale Strukturen können hilfreich sein, um Machine-Learning-Modelle zu bauen, die besser generalisieren und robuster sind“, sagt der Wissenschaftler. Denn nur allzu oft funktionieren die Modelle zwar hervorragend auf den Daten, auf denen sie trainiert wurden (also unter Laborbedingungen), versagen aber, sobald sie sich in der Praxis bewähren sollen: Dann stellt sich heraus, dass das Trainingsmaterial die Wirklichkeit nur eingeschränkt oder gar stark verzerrt repräsentiert hat.

## Fehlgeleitete KI

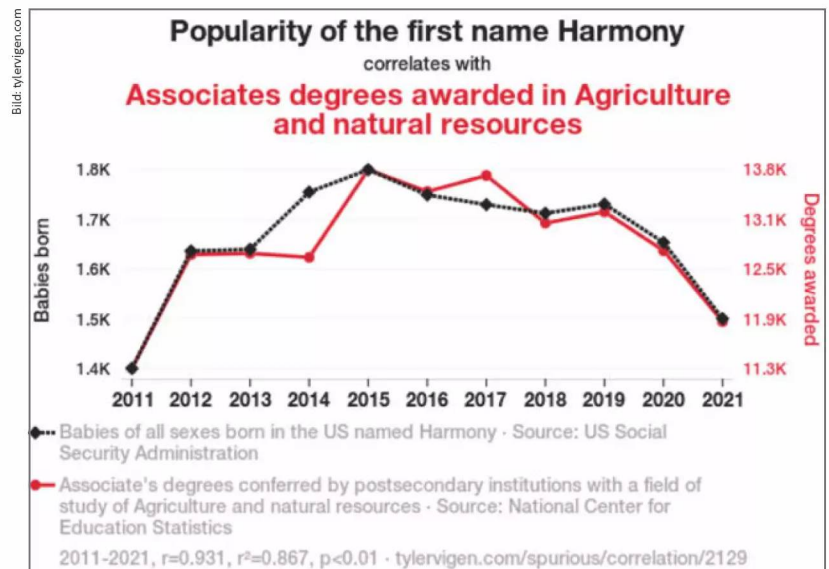
An diesem Problem scheiterte beispielsweise der Versuch, mittels maschinellen Lernens auf Röntgenbildern Pneumothorax zu diagnostizieren, also anhand typischer Gewebeveränderungen festzustellen, ob eine Lunge kollabiert ist. Das hat zwar mit dem Datensatz, auf dem das Modell trainiert wurde,

bestens funktioniert. Letztendlich stellte sich jedoch heraus, dass sich das neuronale Netz bei der Klassifizierung stark auf die Anwesenheit von Plastikschläuchen auf den Trainingsbildern verlassen hat, die die Ärzte den Patienten in die Lunge gesteckt hatten. Die Anwesenheit von Schläuchen auf einem Röntgenbild mag zwar mit einer kollabierten Lunge korrelieren; sie als Ursache für eine Diagnose heranzuziehen, ergibt aber keinen Sinn. „Dieses Feature ist ja selbst eine Wirkung der Pneumothorax-Diagnose und keine Ursache dafür“, erklärt von Kügelgen und erläutert, dass diese Schläuche durchaus in einem anderen Trainingsdatensatz von Nutzen sein könnten. In diesem Fall, wo es darauf ankommt, typische Gewebeveränderungen zu erkennen, haben sie die Modelle aber in die Irre geführt.

Um nicht auf solche Scheinkorrelationen („Spurious Noncausal Features“) hereinzufallen, muss ein Machine-Learning-Modell lernen, sie von tatsächlichen Ursachen zu unterscheiden. „Das ist grundsätzlich schon einmal schwierig und wird quasi unmöglich, wenn nur ein einziger Datensatz zur Verfügung steht“, sagt von Kügelgen. „Was man braucht, ist eine gewisse Diversität der Daten.“

Schließlich muss ein Modell, das unter variablen Bedingungen robust zuverlässige Ergebnisse liefern soll, auch Trainingsdaten sehen, die unter unterschiedlichen Bedingungen entstanden sind, da es

**Die Popularität des Vornamens „Harmony“ in den USA korreliert zwar mit der Anzahl der College-Abschlüsse in Landwirtschaft. Beide Phänomene hängen aber nicht kausal zusammen. KI erkennt diesen Unterschied nicht ohne Weiteres.**



sonst ja gar nicht wissen kann, wogegen es robust sein soll. Das oben beschriebene Problem könnte man etwa mit einer ausgewogeneren Sammlung von Röntgenbildern lösen, die teils vor und teils nach dem Einführen der Schläuche gemacht wurden, etwa weil sie aus unterschiedlichen Krankenhäusern mit unterschiedlichen Routinen stammen. Dadurch fallen die Schläuche als mögliche Ursache weg, während die wahren Ursachen für eine Krankheit in der Regel auch unter sich verändernden Rahmenbedingungen gleich bleiben.

Letztlich bedeutet das eine Abkehr von der Verwendung sogenannter IID-Daten (Independent and Identically Distributed), auf dem im Grunde die gesamte Statistik und auch immer noch ein Großteil von künstlicher Intelligenz und maschinellem Lernen basiert. Das klassische Beispiel für solche IID-Daten ist der Würfelwurf: Seine Ergebnisse sind unabhängig voneinander und die identische Verteilung sorgt dafür, dass dasselbe Experiment immer wieder durchgeführt werden kann. Ähnliches gilt auch für Lungenröntgenbilder, die unter nahezu identischen Bedingungen, von den immer gleichen Maschinen im selben Krankenhaus erstellt werden.

Mustererkennung in Röntgenbildern ist also ein komplexes statistisches Problem, weil Bilder naturgemäß viele Merkmale (Variablen) enthalten, die

voneinander abhängig sein und sich gegenseitig beeinflussen können. Das Problem verschärft sich dadurch, dass die Trainingsdaten wie im Fall der Röntgenaufnahmen häufig unvollständig und verzerrt sind – also die wahren Wahrscheinlichkeiten für das Auftreten bestimmter Muster unter bestimmten Bedingungen daraus gar nicht abgeleitet werden können. Ganz anders verhält es sich bei klassischen Zufallsexperimenten wie dem Würfeln oder dem Roulette: Die Ereignisse beeinflussen sich nicht gegenseitig, egal ob man einmal oder zweimal würfelt oder zwischendrin noch Roulette spielt.

## Jenseits der Laborbedingungen

Dass Trainingsdatensätze so unterschiedlich ausfallen, jeder auf seine eigene Art verzerrt sein kann und sich auch im Laufe der Zeit Verschiebungen ergeben (etwa durch bessere Aufnahmetechnik), sieht Kügelgen sogar als Chance: „Genau solche Effekte können uns helfen, kausale Strukturen zu lernen“, sagt von Kügelgen. „Das Wichtige dabei ist aber, diese Daten nicht einfach in einen großen Topf zu schmeißen, einmal umzurühren und das dann alles als einen großen Datensatz zu behandeln.“ Vielmehr solle man diese Unterschiede als potenzielles Lernsignal wertschätzen.

Die Diversitätsstruktur in den Daten ist also ein Schlüssel zu kausalem maschinellem Lernen – die Wissenschaftler sprechen von „Invariant Causal Prediction“. In Bezug auf das Pneumothorax-Beispiel würde das bedeuten, ein Modell auf zwei separaten Datensätzen (einmal mit und einmal ohne Schläuche) zu trainieren und es so lange zu variieren, bis es auf beiden gleich gut funktioniert. „Das ist eine der Ideen, die meine Kollegen und ich hier an der ETH nutzen wollen, um automatisch aus Daten kausale Informationen zu gewinnen“, sagt von Kügelgen. „Denn wenn man Muster findet, die in mehreren Datensätzen bei unterschiedlichen Konditionen stabil sind, dann ist das immer ein guter Hinweis auf etwas wirklich Kausales.“

## Causal ML für mehr Fairness

Einen ganz anderen Zugang zu kausalem maschinellem Lernen hat Ludwig Bothmann, der am Lehrstuhl von Prof. Bischl an der LMU München die Forschungsgruppe für Causal and Fair Machine Learning leitet. Anstatt den Modellen zu helfen, kausale Zusammenhänge zu erkennen, hat er die Trainingsdaten automatischer Entscheidungssysteme im Vi-

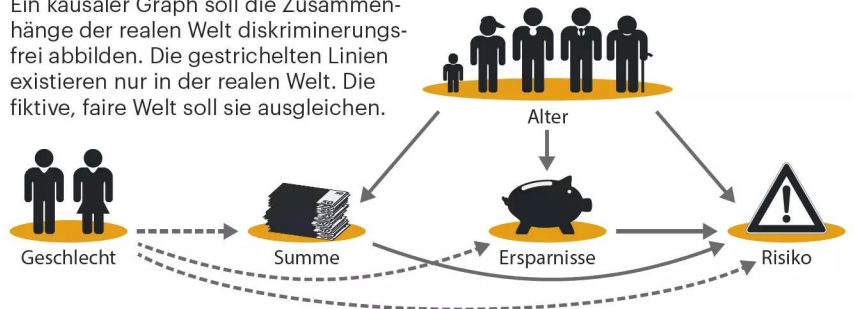


Bild: Chest X-Ray, Posteroanterior View, O'Dea

**Maschinelles Lernen identifiziert auf Röntgenaufnahmen Schläuche als Ursache für eine Lungenkrankung. Sie sind allerdings die Behandlungsfolge einer bereits diagnostizierten Erkrankung.**

## Kausaler Graph

Ein kausaler Graph soll die Zusammenhänge der realen Welt diskriminierungsfrei abbilden. Die gestrichelten Linien existieren nur in der realen Welt. Die fiktive, faire Welt soll sie ausgleichen.



sier. Solche Systeme kommen etwa zur Einschätzung der Rückfallwahrscheinlichkeiten von Angeklagten oder der Kreditwürdigkeit von Personen zum Einsatz und haben damit direkten Einfluss auf menschliche Schicksale. Indem sie aus historischen Daten lernen, versuchen sie etwa vorherzusagen, mit welcher Wahrscheinlichkeit eine einmal strafällig gewordene Person rückfällig wird.

„Die historischen Daten, aus denen solche Systeme lernen, enthalten aber leider oft Diskriminierungen etwa gegenüber ethnischen Zugehörigkeiten oder dem Geschlecht“, sagt Bothmann. „Und wenn sie dann ihre Entscheidungen treffen, schreiben sie die Diskriminierungen der Vergangenheit in die Zukunft fort.“ So könnte allein die Hautfarbe von Angeklagten dazu führen, dass sie aufgrund von Racial Profiling öfter ins Visier der Polizei geraten und aufgrund dessen auch ihre Rückfallwahrscheinlichkeit höher eingeschätzt wird. Oder es fällt Frauen aufgrund der immer noch gravierenden Einkommensunterschiede schwerer, Kredite zurückzuzahlen. Das Entscheidungssystem lernt dann, das Risiko für einen Zahlungsausfall abhängig vom Geschlecht höher oder niedriger anzusetzen, obwohl es eigentlich ausschließlich das Einkommen berücksichtigen sollte. Anhand von diskriminierenden Entscheidungen lernen also auch Machine-Learning-Systeme, Menschen unfair zu behandeln.

„Wir versuchen mit kausalem maschinellen Lernen, diesen Teufelskreis der Diskriminierung zu durchbrechen“, sagt Bothmann. Gemeinsam mit seinen Kolleginnen und Kollegen konstruiert er deshalb zunächst fiktive Welten, in denen geschützte Attribute keine kausalen Effekte mehr auf die Ent-

scheidungen der Systeme haben. Voraussetzung dafür ist ein möglichst realistischer und detaillierter „kausaler Graph“, der die ursächlichen Zusammenhänge zwischen den verschiedenen Eingangsvariablen wie Geschlecht, Einkommen sowie Spareinlagen und der Risikoeinschätzung in Form von Pfeilen darstellt.

Stellt sich dann heraus, dass etwa aufgrund des Gender-Pay-Gaps für eine Frau das Risiko höher ist, einen Kredit nicht zurückzahlen zu können, muss dieser Effekt in der fiktiven, fairen Welt zunächst ausgeglichen werden. „Man könnte der Person dann etwa ein fiktives, an eine vergleichbare Erwerbsbiografie eines Mannes angepasstes Einkommen zuschreiben“, sagt Bothmann. So entsteht ein neuer Satz von Eingangsdaten für das Training eines Entscheidungssystems, in dem das Geschlecht keinen Einfluss mehr auf das Ergebnis hat.

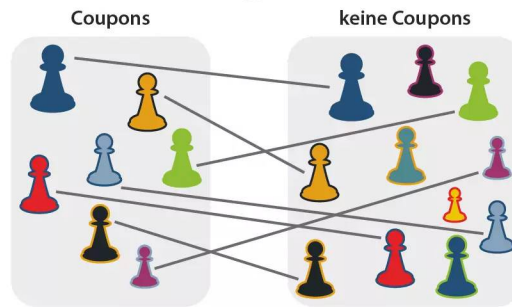
Solche Systeme werden zunächst natürlich schlechtere Vorhersagen treffen, als jene, die mit unveränderten, realen Daten trainiert wurden und davon ausgehen, dass die Welt auch weiterhin unfair bleibt. „Wenn die Gesellschaft mehr Fairness will, ist das aber gewissermaßen der Preis dafür“, sagt Bothmann. „Und wenn der gesellschaftliche Druck groß genug ist, um alle automatischen Entscheidungssysteme derart fair zu gestalten, sollte das letztendlich auf die reale Welt zurückwirken und zu mehr tatsächlicher Fairness führen.“

Bei realen Anwendungen, die künstliche Intelligenz mit Kausalität verbinden, ist in der Regel also menschliche Expertise gefragt. Menschen brauchen nicht erst lange nachzudenken, ob Schokoladenkonsum die Ursache für mehr Nobelpreisträger im



## Simulierte Kaufentscheidung

Die Figuren dieser beiden Gruppen repräsentieren Kunden mit ähnlichem Kaufverhalten, Alter oder Wohnort. Causal ML simuliert ihre Kaufentscheidung nach Erhalt oder Nichterhalt eines Rabatts.



Land sein kann. Schließlich haben wir alle schon Schokolade gegessen und keine Hinweise darauf bemerkt, dass sie uns vielleicht bei einer akademischen Karriere nützlich sein könnte. Maschinen dagegen haben noch keine eigenen Erfahrungen mit der Welt gemacht, weshalb sie bei solchen Fragestellungen meist völlig im Dunkeln tappen.

Zwar ist es möglich, kausale Zusammenhänge zwischen Variablen auch automatisch zu erkennen. Allerdings handelt es sich dabei eher noch um Grundlagenforschung, die nur in einfachen Fällen mit geringer Komplexität funktioniert. In der Realität scheitern solche Methoden meistens daran, dass es entweder zu viele Variablen gibt oder die Ergebnisse mehrdeutig sind, also gleich mehrere verschiedene kausale Graphen das Ergebnis der Unabhängigkeitstests erklären könnten.

## Simulierte Rabattmarken

Vor allem in der Ökonometrie, die seit jeher mit statistischen und mathematischen Methoden versucht, wirtschaftliche Zusammenhänge zu analysieren, ist Causal ML bereits verbreitet. „Big-Tech-Companies wie Amazon oder Google verwenden diese Methoden ständig“, sagt Martin Huber, Professor für Angewandte Ökonometrie und Politik-evaluation an der Universität Freiburg. „Allerdings steht da immer auch ein Mensch dahinter, der für jeden Anwendungsfall ein geeignetes kausales Modell erstellen muss.“ Auf diese Art lassen sich ganz normale, auf Mustererkennung trainierte Machine-Learning-Modelle so geschickt miteinander kombinieren, dass sich aus ihren Ergebnissen beispielsweise die Wirkung eines Gutscheins oder

eines Rabatts auf das Kaufverhalten bestimmter Kundengruppen eruieren lässt.

In einer eigenen, erst kürzlich veröffentlichten wissenschaftlichen Arbeit ist es Huber gemeinsam mit einer Kollegin aus Finnland gelungen, lediglich auf Basis öffentlich zugänglicher Daten eine solche „Impact Evaluation“ durchzuführen. Würde man so ein Experiment mit realen Kunden durchführen, wäre die Aufgabe ziemlich unkompliziert: Man müsste die Probanden lediglich per Zufallsentscheidung in zwei Gruppen unterteilen und in einer Gruppe Rabattkärtchen verteilen, während die Vergleichsgruppe leer ausgeht. Lässt sich schließlich ein Unterschied im Kaufverhalten der beiden Gruppen beobachten, kommt als Ursache dafür nur der Erhalt des Coupons infrage. Der Händler kennt dann den Effekt seiner Aktion und kann feststellen, ob sie sich gelohnt hat.

Solche Experimente kosten aber Zeit und Geld. „Hier kommt das maschinelle Lernen ins Spiel“, sagt Huber. „Es erlaubt uns, aus passiv gesammelten Daten der Vergangenheit gewissermaßen ein solches Experiment zu simulieren.“ Was man dafür benötigt, sind große Mengen an Kundendaten, wie sie etwa durch Treuekarten gewonnen werden. Die Krux dabei: Selbst wenn die Daten sowohl Informationen über Kunden enthalten, die Coupons bekommen haben, und über solche, die keine bekommen haben, ist nicht klar, ob die Verteilung rein zufällig passiert ist. So könnten zum Beispiel auch vorwiegend jene Personen Coupons erhalten haben, die bereits sehr treue Kundinnen und Kunden waren und allein deshalb schon mehr einkaufen. Der Erhalt des Coupons wäre dann also keine Ursache für mehr Konsum, sondern selbst eine Wirkung davon. „Deshalb ist es

entscheidend, zwei Gruppen zu bilden, die sich zwar durch den Erhalt beziehungsweise Nichterhalt eines Coupons unterscheiden, in Bezug auf andere relevante Variablen wie das frühere Kaufverhalten aber vergleichbar sind“, erklärt Huber.

Welche Variablen das genau sind, können zwei einfache Machine-Learning-Schritte klären. In dem vereinfachten Beispiel würde der erste Schritt eine starke Korrelation zwischen dem früheren und dem aktuellen Kaufverhalten feststellen. Und der zweite eine starke Korrelation zwischen früherem Kaufverhalten und dem Erhalt eines Coupons. „Damit wäre das frühere Kaufverhalten als einer der Faktoren identifiziert, in denen die beiden Gruppen unseres simulierten Experiments zur Impact Evaluation vergleichbar sein müssen“, erklärt Huber. „Erst dann kann unterschiedliches Kaufverhalten zwischen den beiden Gruppen auch wirklich als kausaler Effekt des Coupons interpretiert werden.“

## Kunden im kausalen Wald

Für derartige Auswertungen setzen Huber und seine Co-Autorinnen und Co-Autoren auf Programmiersprachen wie R oder Python. „Dazu gibt es auch bereits eine ganze Menge kausaler Machine-Learning-Algorithmen, die sich schon mit ein paar Zeilen Code in die Basisversion implementieren lassen“, sagt Huber. „Das ist mittlerweile sehr benutzerfreundlich und sowohl die Variablenauswahl als auch die Messung des Effekts der Coupons funktionieren im Grunde auf Knopfdruck.“ Der Algorithmus, der dabei zum Einsatz kam, heißt Causal Forest, also kausaler Wald. Er ist einer der meistverwendeten Algorithmen und Huber zufolge bereits sehr gut ausgereift. Wichtig sei es dennoch, die dahinterstehenden Annahmen zu den kausalen Zusammenhängen zu verstehen. Und ein kausaler Effekt lässt sich auch nur dann zuverlässig identifizieren, wenn der Datensatz wirklich alle relevanten Kundeninformationen enthält.

Wenn die Daten genug statistisch relevante Informationen zu den Untergruppen enthalten, können Methoden wie Causal Forest darüber hinaus auch rein datengetrieben herausfinden, für welche Untergruppen die Effekte besonders groß oder besonders klein sind. „Der Algorithmus sagt mir dann beispielsweise, dass ich die Coupons vor allem an Leute abgeben sollte, die früher besonders viel oder besonders wenig gekauft haben, weil sie in diesen Kundengruppen besonders stark wirken“, sagt Huber.

Ähnliche Methoden werden mittlerweile auch in der Medizin angedacht. Zwar gilt auch hier das Experiment in Form randomisierter klinischer Studien mit Placebogruppe immer noch als Goldstandard, um etwa die Wirksamkeit von Medikamenten zu testen. Oft ist das allerdings aus ethischen Gründen gar nicht möglich, weshalb neue Möglichkeiten ausgelotet werden, Patientendaten mittels kausalen maschinellen Lernens auszuwerten.

Eine typische Fragestellung in diesem Bereich wäre etwa, ob tatsächlich ein bestimmtes Medikament Ursache für eine Genesung ist oder ob nicht vielleicht eine dritte Variable wie die finanzielle Situation eines Patienten sowohl Ursache für den Erhalt des Medikaments als auch einen gesünderen Lebensstil sein könnte. Der Medizinforschung kommt dabei zugute, dass aufgrund der fortschreitenden Digitalisierung immer mehr Patientendaten gesammelt und aufbereitet werden, mit denen die Forschenden ihre Machine-Learning-Modelle füttern können.

## Ausblick

Während Causal ML nicht zuletzt auch aufgrund höherer ethischer Standards zum Schutz der Privatsphäre in der Medizin erst am Anfang steht, ist man in der Ökonometrie schon weiter. „Allerdings sind es bisher hauptsächlich große Konzerne, die solche Kausalanalysen durchführen“, schränkt Huber ein. Schließlich könne bei einem weltumspannenden Geschäftsmodell auch ein kleiner Effekt schnell hohe Geldsummen erwirtschaften, die die Ausgaben für Forschung und die Entwicklung der Methoden rechtfertigen.

In mittelständischen Unternehmen gebe es dagegen oft noch nicht einmal das Bewusstsein für die Unterschiede zwischen einer Vorhersage und einer Kausalanalyse. „Inzwischen veröffentlichen die Big-Tech-Unternehmen aber auch wissenschaftliche Papiere zu den von ihnen entwickelten Methoden oder stellen sogar Open-Source-Software zur Verfügung“, sagt Huber. „Das ist eine spannende Entwicklung und wird wohl dazu führen, dass bald auch kleinere Unternehmen diese Methoden einsetzen.“

Das Potenzial der Verknüpfung von maschinellem Lernen mit den Prinzipien von Ursache und Wirkung ist jedenfalls groß. Und wie bereits bei so vielen anderen Spielarten künstlicher Intelligenz zuvor wird wohl auch der Einfluss dieser neuen Methoden auf unser Leben weiter zunehmen. (akr) **ct**



Bild: KI, Collage c't

# KI im Rotlichtviertel

Die Haut sogenannter Sexroboter besteht aus beheizbarem Silikon, Servomotoren bewegen ein Skelett aus Metall. Erste Prototypen zeigen Mimik und führen mittels künstlicher Intelligenz einfache Dialoge. Künftig sollen sie das Gesicht ihres Nutzers erkennen und ein persönliches Profil speichern. Ist das wünschenswert, verwerflich oder gar gefährlich?

Von **André Kramer**

**S**ie heißen „Recreational Androids“. So steht es unter der Fußsohle humanoider Roboter der Marke RealDoll des US-amerikanischen Herstellers Abyss Creations. Zu Deutsch bedeutet das Freizeitandroide, das vorangestellte Adjektiv steht für „der Erholung oder Entspannung dienend“, angelehnt an Begriffe wie „Recreational Vehicle“ (Wohnmobil) oder „Recreational Drug“ (weiche Droge). Die Puppen dienen dem Zweck, sexuelle Bedürfnisse zu befriedigen. Sie sehen im Detail möglichst lebensecht aus. Ihre Proportionen wirken jedoch wie die einer stereotypen Frau der Kategorie „sexy“, die Kur-

ven ausladend, die Wimpern lang, das Make-up kräftig. Die aufreizende Hülle füllen einige Hersteller mit Elektronik, die aus den Puppen mehr oder minder smarte Androiden macht.

Ihre Namen – „Jackie“, „Tiffany“, „Misaki“, „Rocky“ oder „John“ – klingen wie aus einem Pornofilm entnommen. Das ist kein Zufall: Echte Personen verkaufen ihr Aussehen, das die Hersteller auf die Sexpuppen übertragen. Vor allem Pornodarstellerinnen sind bei der Kundschaft beliebt. Die Puppen kommen zudem besser an, wenn sie lebenssechte Attribute wie Muttermale oder Tattoos aufweisen. Ein Nutzer



aus Taiwan habe Scarlett Johansson nachbilden lassen, sagte ein Händler gegenüber c't. Bestätigen lässt sich so eine Aussage nicht.

Obwohl es sich um leblose Produkte aus Kunststoff und Metall handelt, vermenschlicht man die Gegenstände unwillkürlich und spricht von „ihr“ oder „ihm“. Dabei wäre das Pronomen „es“ passender für die lebensgroßen Actionfiguren. Mit ihren billigen, aufblasbaren Verwandten vom Junggesellenabschied haben Tiffany und Rocky allerdings nicht viel gemein und das hat seinen Preis: Sehr günstige Vertreter gehobener Sexpuppen kosten im Sonderangebot etwas über 500 Euro. In der Regel müssen Kunden mehrere Tausend Euro hinblättern. Mit Roboterkopf liegt man schon im fünfstelligen Bereich. Das können und wollen sich nur wenige Kunden leisten.

## Im Berliner Cyberbordell

Das Cybrothel, ein Kofferwort aus Cyber und Brothel (Bordell), in Berlin-Friedrichshain verkauft Rollenspiele mit solchen Sexpuppen als Dienstleistung und kombiniert sie mit VR-Inhalten oder einem KI-gesteuerten Livedialog. Eine Puppe zu buchen, hat nicht nur einen preislichen Vorteil: Fragen wie „Steht die auf mich?“ spielen keine Rolle.



**An Seilen hängen unterschiedliche Puppenkörper im Hinterzimmer des Berliner Cybrothel von der Decke.**

Wir fahren zum Hausbesuch ins Cybrothel. In mehreren Reihen hängen im Hinterzimmer verschiedene Puppenkörper an Seilen von der Decke: kleine, große, üppige, zierliche, männliche, weibliche und sogar ein außerirdischer. Mitarbeiter reinigen und desinfizieren sie, ziehen sie nach Buchungswunsch an und statten sie mit Köpfen sowie Perücken aus. Schauspielerinnen, die „Voice-Queens“, leihen den Puppen derzeit noch ihre Stimmen. Die Schauspielerinnen verkörpern je nach gebuchter Geschichte Rollen wie eine Lehrerin, eine Schülerin oder eine Domina.

Das Programm soll stetig erweitert werden. Vor allem Geschichten aus Fantasy und Science-Fiction möchte das Cybrothel bedienen. Außer Luna, der Vampirin, halten sie eine blaue Asari-Alienpuppe aus der Mass-Effect-Serie für Sci-Fi-Fantasien vorrätig. Der Succubus, ein weiblicher, aus Videospielen bekannter Dämon, ist gefragt, aber auch Meerjungfrauen. Eine Orientierung gibt die Suchstatistik von Videospielcharakteren auf Pornhub. Hoch im Kurs stehen dort Fortnite- und Overwatch-Charaktere sowie Lara Croft aus Tomb Raider und die überlebensgroße Lady Dimitrescu aus dem First-Person-Survival-Horror-Spiel „Resident Evil Village“.

Einen Aufenthalt bucht man im Web, per E-Mail oder WhatsApp. Kunden können sich pseudonym anmelden, bar bezahlen und das Cybrothel anonym besuchen. Ein Mitarbeiter überwacht aus dem Kontrollraum den Eingangsbereich mit einer Kamera. Danach folgen die Kunden einer Stimme aus dem Lautsprecher ins eigentliche Cybrothel. Es ist ausgestattet mit einem Bett, einer Liebeschaukel und einem Beamer. Langfristig soll virtuelle Realität den Raum digital ersetzen. In Kooperation mit der Filmfirma BaDoink hat das Cybrothel einen VR-Porno gedreht: Das intergalaktische Wesen Kokeshi landet auf der Erde und möchte lernen, wie Menschen Sex haben. Alternativ steuern Besucher im auch auf Steam erhältlichen Videospiel VRX Cherry von Polybay über einen Hüft-Controller einen virtuellen Penis.

Eine Stunde im Basic-Tarif mit einer Puppe nach Wahl kostet 89 Euro. Die meisten entscheiden sich für eine Stunde, wenige bleiben zwei Stunden oder über Nacht. Extras lassen sich hinzubuchen: 60 Minuten Mixed-Reality-Sex samt Puppe und VR-Headset kosten 199 Euro. „Bei uns bucht man die Session nicht zwei Monate im Voraus“, sagt Gründer Philipp Fussenegger. Das läuft eher so: „Ich bin jetzt geil, ich komme in zwei Stunden vorbei.“

Die Kunden sind im Durchschnitt etwa 30 Jahre alt, interessieren sich für Technik, Gaming sowie virtuelle Realität und möchten etwas Neues aus-



**Je nach Kundenwunsch stattent Cybrothel-Mitarbeiter ihre Puppen mit verschiedenen Köpfen und Perücken aus.**

probieren. „Das Klischee vom Vater, der seinen Sohn zum 18. Geburtstag mit auf die Reeperbahn nimmt, existiert nicht mehr“, sagt Fussenegger. „Heute ist man auf Onlyfans oder möchte mit seinem Lieblings-spielecharakter Sex haben.“

## Silikon über metallischem Endoskelett

Direkt menschlich fühlen sich die Puppen nicht an. Unter dem Arm aus Silikon spürt man deutlich das Metallskelett. Handgelenk, Schulter, Ellenbogen, Knie oder Hüfte lassen sich bewegen wie bei einer übergroßen Barbie.

Frühe Vertreter hatten eine unnatürlich glatte Haut. Aktuelle hochpreisige Sexpuppen bestehen aus thermoplastischen Elastomeren (TPE) oder Silikon. Puppen aus TPE sind weicher, elastischer, aber auch empfindlich gegenüber Hitze und Flecken. Das Material kann außerdem schwitzen und muss daher mit Babypuder behandelt werden. Silikonpuppen sind leichter zu reinigen und weniger empfindlich, aber teurer und fester.

Das Gerippe aus Metall bildet nicht nur die Gelenke des menschlichen Körpers grob nach, sondern gibt der Puppe auch die nötige Stabilität. Ein Real-Doll-Skelett besteht aus Stahl oder Verbundmetall mit 14 beweglichen Gelenken an Nacken, Schultern, Ellenbogen, Händen und Lenden. Die Finger bestehen aus Metalldraht. Eine Neuerung ist ein weiteres Gelenk, das die Achseln nach oben bewegt, damit

die Puppe ein Stück menschenähnlicher wirkt, wenn man sie am Oberarm anhebt.

Je nach Modell haben die Puppen eine Körperhöhe von 1,40 bis 1,75 Metern. Ihre Größe wirkt sich nicht nur auf den Verkaufspreis aus, sondern auch auf das Gewicht. Sie wiegen zwischen 28 und 40 Kilogramm. Kleine Puppen wären günstiger und leichter zu handhaben, würden aber allein aufgrund ihrer Größe wie Kinder wirken. Nicht zuletzt vermutlich deswegen hat die kleinste Puppe von IronTech Doll besonders üppige Proportionen.

## Frühstück mit Tiffany

Frühe Sexroboter scheiterten an technischen und finanziellen Hürden. Der Prototyp Roxxxy von der Firma TrueCompanion sollte verschiedene, in Software realisierte Persönlichkeiten verkörpern: die offenerzige „Wild Wendy“, die eher zugeknöpfte „Frigid Farrah“ und die unerfahrene „Young Yoko“. TrueCompanion stellte sein Konzept im Jahr 2010 vor und nahm in der Folgezeit 4000 Vorbestellungen zu Preisen ab 9995 US-Dollar entgegen. Das Medien-echo war schon im Vorfeld groß; es hagelte Kritik an den stereotypen Charakteren. Materialisiert hat sich die smarte Sexpuppe indes nicht. Im Jahr 2017 berichtete der Guardian, dass nie ein Produkt der Firma bei einem Kunden ankam. Die Domain der Herstellerseite steht seither wieder zum Verkauf.

Der US-Hersteller Abyss Creations bietet Roboterköpfe der Reihe RealDoll X an, die man im Webshop



zusammen mit einem nichtrobotischen Puppenkörper erwerben kann. Ein Roboterkopf der Figur Harmony, Solana, Serenity, Nova oder Tanya kostet jeweils rund 7500 US-Dollar. Für einen Standardkörper kommen 4000 US-Dollar hinzu. Piercings und ähnliches kosten extra. Die Lieferzeit beträgt laut Hersteller sechs bis acht Monate.

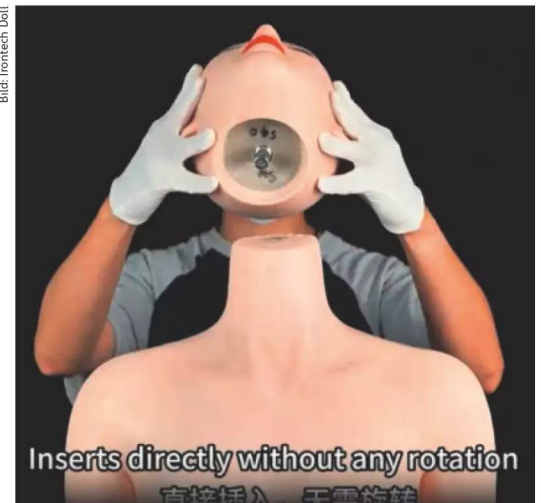
Videos zeigen, wie die Roboterköpfe auf Gesprächsbeiträge reagieren und beim Sprechen Kopf, Augen sowie Lippen bewegen. Die hölzernen Bewegungen erinnern an Kirmesanimatronik und rechtfertigen nur schwer den saftigen Aufpreis. Über die Webseite realdollx.ai kann man die realdollx-App für 20 US-Dollar erwerben und als APK-Datei für Android herunterladen, um mit der KI zu chatten. Eine iOS-Version gibt es gar nicht. Das dürfte sich auch schwer mit den Richtlinien von Apple und Google für den App respektive Play Store decken. Mit uns sprechen wollte der Hersteller trotz mehrfacher Anfragen nicht.

## Von der Puppe zum Roboter

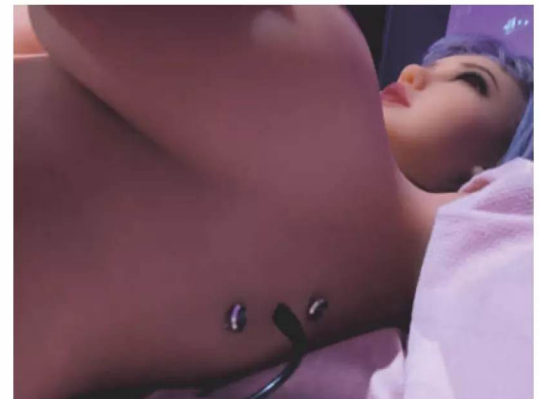
Auch der chinesische Hersteller Irontech Doll fügt seinen Puppen sukzessive elektronische Zusatzfunktionen hinzu. Mit Heizelementen ausgestattet, sollen sich die Körperteile aus Kunststoff etwas lebensechter anfühlen. Mithilfe von Motoren im Hüftbereich lassen sich die schweren Puppen leichter in gewisse Positionen bringen. Sie können den Mund bewegen, allerdings nicht, um Mimik auszudrücken, sondern um eingeführte Gegenstände oder Körperteile anzusaugen. Vergleichbare Funktionen sind in der nachgebildeten Vagina verbaut. Arme und Beine der Puppenkörper lassen sich abgesehen vom Hüftmotor nur passiv bewegen. Die kommerziell verfügbare Puppenelektronik hat mit Robotik bisher nicht viel zu tun. Das soll sich aber ändern.

Wie Abyss Creations bewirbt auch Irontech Doll auf Messen sowie auf Facebook, Instagram und X den Prototyp eines Androidenkopfes, der Mimik zeigt. Er kann die Augen von links nach rechts wandern lassen, kneift beim Lächeln die Lider zusammen und bewegt beim Sprechen den Unterkiefer, was jedoch wie gedankenverlorenes Kaugummikauen aussieht. Der Kopf soll künftig im Kontext der Interaktion mit dem Nutzer lächeln, zwinkern, nicken und die Stirn runzeln. Mit einem Large-Language-Model (LLM) verknüpft, die auch die Basis für KI-Chatbots wie ChatGPT bilden, soll er Anweisungen und Fragen verstehen und darauf reagieren können. Für den Anfang will auch der chinesische Hersteller eine

Bild: Irontech Doll



**Ein Tutorial von Irontech Doll erklärt, wie der separat gelieferte Kopf auf den Körper montiert wird.**



**Eine Buchse für ein Netzteil verbirgt sich in der Haut aus thermoplastischem Elastomer und zerstört die Illusion.**

Smartphone-App veröffentlichen, über die Nutzer mit der KI interagieren können.

Touch-Sensoren erkennen in künftigen Versionen Berührungen und lösen Reaktionen aus. Die Puppen sollen letztlich emotionale Intelligenz simulieren, etwa Empathie und Freude ausdrücken oder Trost spenden. Irontech Doll trainiert dabei nach eigenen Angaben eine integrierte KI. „Die durch Interaktionen





Bild: Irontech Doll

**Der Androidenkopf von Irontech Doll lässt die Augen zur Seite wandern, kneift die Lider zu und bewegt den Mund beim Sprechen.**

gesammelten Daten werden verwendet, um die Reaktionen und Verhaltensweisen der Puppe im Laufe der Zeit zu personalisieren und ihre Persönlichkeit an die Vorlieben und Wünsche des Benutzers anzupassen“, schreibt der Hersteller. Das schließt auch Gesichtserkennung ein.

Über Bluetooth und WLAN greift die Puppe auf das Internet zu, um Informationen abzurufen, Geschichten zu erzählen oder mit dem Smart Home zu interagieren. „Bei unserem Ansatz zur KI-Entwicklung berücksichtigen wir stets die ethischen Standards, Datenschutzbedenken und sozialen Auswirkungen der Nutzung solcher Technik“, versichert Irontech.

Solche Ankündigungen rufen Datenschützer auf den Plan. „Die KI des Roboters speichert biometrische Daten, die der Identifizierung von Personen dienen. Sie sind personenbezogen und somit findet das Datenschutzrecht Anwendung“, sagt Iris Phan, Juristin und Lehrbeauftragte für Roboterethik des Instituts für Philosophie der Leibniz-Universität Hannover im Gespräch mit c't. Neben der Gesichtserkennung werden auch Gesundheitsdaten, Daten zur sexual-

len Orientierung oder zum Sexualleben gesammelt. „Zusammen mit ‚persönlichen Vorlieben‘, die der Nutzer selbst eingeben kann, entsteht hier ein sehr intimes Profil“, so Phan. Die Sexroboter sollen sich mit dem Internet verbinden und es braucht nicht viel Fantasie, sich bei Datendiebstahl den zugehörigen Erpressungsversuch vorzustellen.

## Künstliche Puppenintelligenz

Das Thema Datenschutz treibt auch die Berliner Entwickler des Cybrothel um. Langfristig soll dort künstliche Intelligenz die menschlichen Voice-Queens vollständig ersetzen. Gegenwärtig experimentieren die Berliner mit beiden Ansätzen. „Wir haben die erste sexpositive, queere KI der Welt“, sagt Cybrothel-Entwickler Matthias Smetana im Gespräch mit c't. „Die läuft hier lokal auf einem Gaming-PC.“

Aus Datenschutzgründen laufen das LLM sowie Spracherkennung und -synthese vor Ort im Kontrollraum. „Über APIs ließen sich zwar ChatGPT einbinden und Stimmen generieren, aber das landet dann gleich auf drei, vier verschiedenen Servern“, gibt Smetana zu bedenken. „Wir binden Open-Source-Modelle ein.“ Nicht zuletzt erlaubt der ChatGPT-Entwickler OpenAI auch keine Sexanwendungen. Das Gespräch wäre schnell am Ende, wenn ChatGPT die eindeutigen Avancen der Nutzer zurückweist.

Die Kunden interagieren im Cybrothel über Mikrofon und Lautsprecher mit der KI nebenan. Immer nur der letzte Input und der letzte Output werde gespeichert, versichert Smetana. „Wir haben keine History auf unserem KI-Modell.“ So begegnen die Betreiber Datenschutzbedenken. Sie gehen gleichzeitig davon aus, dass die Leute mehr KI als menschliche Voice-Queens buchen. „Die Erfahrung zeigt, dass sich die Kunden der KI in der Anonymität mehr öffnen als einer menschlichen Gesprächspartnerin“, sagt Cybrothel-Betreiber Fussenegger. Die Cybrothel-Entwickler haben thematische Beschränkungen in ihr Sprachmodell eingebaut. Fragen Kunden nach Sex mit Kindern, Inzest oder Nekrophilie, zeigt die KI Grenzen auf.

## Mensch und Maschine

Ähnlich wie bei der Videospieldiskussion um „Call of Duty“ oder „Counter Strike“ steht auch bei Sexrobotern schnell die Ethikfrage im Raum, geht das Gespräch Richtung Gewalt oder Pädophilie: Dient das Unterhaltungsprodukt lediglich als Ventil oder trainiert es ein unerwünschtes Verhalten? Der Sexual-



Bild: Capcom



**Besonders Geschichten aus Fantasy und Science-Fiction sind im Cybrothel gefragt. Eine der Puppen dafür heißt Blue und stammt aus dem Videospiel-universum von Mass Effect.**

**Das Webportal Pornhub listet die meistgesuchten Videospielcharaktere. Neben Lara Croft aus Tomb Raider liegt Lady Dimitrescu aus dem First-Person-Survival-Horror-Spiel „Resident Evil Village“ hoch in der Gunst der Nutzer.**

therapeut Ulrich Clement äußerte in Zeit online, dass derzeit unklar sei, welche Theorie zutrifft. Ein Händler gibt auf seiner Webseite Hinweise, um die Lebensdauer der Puppen zu maximieren: „Behandeln Sie Sexpuppen nicht grob.“ Da steht nicht „unsachgemäß“, sondern ein vermenschlichendes Attribut, das in der Bedienungsanleitung eines Toasters sicher keinen Platz fände. Der Hinweis hindert dennoch niemanden, Vergewaltigungsszenarien nachzuspielen oder eine Zigarette auf der Silikonhaut auszudrücken. Umgekehrt ist es auch nicht auszuschließen, dass sich jemand in seinen Gegenstand verliebt, statt ihn nur als unbelebtes Objekt der Lust zu benutzen. Es gibt immerhin das Phänomen der Objektophilie. „Sexroboter haben keinen anderen Status – auch nicht rechtlich – als andere Gegenstände“, gibt Juristin Iris Phan außerdem zu bedenken.

In kontrollierten Umgebungen können Sexroboter gleichwohl eine sinnvolle Funktion erfüllen, indem sie pflegebedürftigen, alten oder Menschen mit Behinderung einen Zugang zu Sexualität ermöglichen, den sie sonst nicht bekommen. Auch in Krankenhäusern und in der Psychotherapie gibt es beispielsweise für die Behandlung von Sozialphobien Anwendungsmöglichkeiten. Einer aktuellen Studie im deutschsprachigen Raum zufolge stehen über 80 Prozent der Befragten einer gesellschaftlichen Nutzung unter bestimmten Voraussetzungen positiv gegenüber, insbesondere bei körperlicher Einschränkung, anstelle von Prostitution und um bestimmte sexuelle Fantasien ausleben zu können. Über 80 Pro-

zent der Befragten akzeptierten ferner die Verwendung von Sexrobotern zur Bearbeitung eines sexuellen Problems. Das berichtete das Deutsche Ärzteblatt im Mai 2019.

Allein schon aufgrund des hohen Preises bleiben Sexroboter vorerst aber ein Nischenprodukt. Die Entwicklung humanoider Roboter verschlingt Unsummen. Denn kaum etwas ist schwieriger, als Menschen glaubhaft nachzubilden. Die Evolution hat unsere Gattung in Millionen von Jahren darauf trainiert, die Körpersprache von Artgenossen zu lesen, beispielsweise Aggression, Zuneigung oder Verletzungen zu erkennen. Kleinste Abweichungen fallen sofort auf. Unnatürlich wirkende Mimik irritiert unmittelbar. Der Effekt ist als „Uncanny Valley“ bekannt: Menschen akzeptieren stark verfremdete Charaktere wie gezeichnete Katzen und Mäuse, die einander mit Bratpfannen verprügeln, und verfolgen die Geschichte mit Vergnügen auf der Leinwand. Ein nur fast perfektes Ebenbild lehnen sie aber als unrealistisch ab.

Solange nachgebildete Mimik wie aus der Geisterbahn wirkt, fahren die Hersteller besser damit, die Puppen stumm sowie unbewegt auszuliefern und den Rest der Fantasie der Kunden zu überlassen. Nur eine kleine Gruppe wohlhabender Kunden kann sich derzeit die elektronischen Replikas leisten. Das muss aber nicht ewig so bleiben: Es wäre nicht das erste Mal, dass ein Produkt den Sprung von einer Nische in den Massenmarkt schafft – und von dort in andere Bereiche des Lebens. (akr) **ct**





Bild: Sebastian Steppuhn

# KI fürs Bienenwohl

Varroamilben gelten als einer der Hauptgründe für das Bienensterben. Der Schüler Sebastian Steppuhn aus Pforzheim hat einen cleveren Apparat gebaut, um befallene Bienen zu erkennen und gezielt zu behandeln – noch bevor die Parasiten in den Stock gelangen und das ganze Volk befallen.

Von **Andrea Trinkwalder**

**D**er Milbendetektor BeeAI des Schülers Sebastian Steppuhn ist das, was man eine Vorzeige-KI nennen könnte: Sie nimmt dem Menschen eine für ihn unlösbare Aufgabe ab, dient dem Gemeinwohl – und schafft dabei keine neuen, nur mit noch ausgefeilteren Algorithmen oder dem Einsatz schlecht bezahlter Hilfskräfte lösbarer Probleme. Mit seiner bis ins Detail durchdachten Schleusenkonstruktion für Bienenstöcke, die krankmachende Mil-

ben draußen halten soll, gewann der 17-jährige Hobby-Imker den Hauptpreis des Bundeswettbewerb für künstliche Intelligenz 2024 (BWKI, siehe Kasten „Reflexionsarme Fotofalle“). Doch der Weg dorthin war steinig.

Zunächst schien es, als ob sich die Idee mit überschaubarem Aufwand realisieren ließe, handelt es sich doch um klassische Bilderkennung: Trainingsdaten in Form von milbenbefallenen Bienen besor-



gen und optimieren, damit nach bewährter Deep-Learning-Manier ein neuronales Netz trainieren, auswerten, fertig. Normalerweise findet man für jedes nur erdenkliche Problem Trainingsdaten, mit denen man erst mal loslegen, locker herumexperimentieren und die Erfolgsaussichten eruieren kann. Für dieses gab es jedoch schlichtweg: nichts, außer vereinzelte Bilder von Bienen, denen die Varroamilbe im Pelz sitzt.

So stand erst mal aufwendige Pionierarbeit unter reichlich unklaren Erfolgsaussichten an. Steppuhn ließ sich davon nicht beirren und entwickelte einige clevere Strategien und Skripte, um das mühsame Datenschnürfen sukzessive zu automatisieren. Und nicht nur das: Er entwarf und fertigte auch einen ausgeklügelten Mechanismus, um befallene Bienen sanft in ein separates „Behandlungszimmer“ zu bug-sieren, bevor sie sich wieder zu ihrem Volk gesellen dürfen. Insgesamt arbeitete der Schüler fast zwei Jahre lang an seinem Projekt, in den Aufbau der Trainingsdatenbank investierte er mehrere Monate.

## Die Milbe im Bienenhaufen

Sie ist nur 1,1 Millimeter lang und 1,6 Millimeter breit – aber ein großes Problem: Die Varroamilbe setzt sich auf Brust oder Rücken erwachsener Honig-

## Frisch aus c't Nerdistan

bienen und saugt an deren Fettkörper. Dabei kann sie Viren oder andere Krankheitserreger übertragen. Befallene Bienen werden schwach und sterben. Besonders verheerenden Schaden richten die Parasiten allerdings in den Brutzellen an, in denen die jungen Bienen heranwachsen. Ohne Behandlung stirbt das befallene Volk deshalb binnen weniger Jahre. Um die Milbe zu bekämpfen, muss bislang das gesamte Volk behandelt werden, in der Regel geschieht das mit Ameisensäure.

Weil diese Methode auch viele gesunde Bienen einer unnötigen Behandlung aussetzt und die Säure überdies zu Brutschäden führt, machte sich der Hobby-Imker Sebastian Steppuhn an die Arbeit und konstruierte einen Kasten mit einem schmalen transparenten Gang. Diesen montierte er vor den Eingang seines Bienenstocks, um die Insekten systematisch mithilfe von darüber und darunter angebrachten Kameras beobachten zu können. Sein Plan war, die Parasiten mithilfe eines klassischen Machine-Learning-Algorithmus zu erkennen und dadurch einen Mechanismus auszulösen, der die befallenen Bienen direkt in einen separaten Behälter lotst. Ein geeignetes vortrainiertes Machine-Learning-Modell für solche Echtzeitanwendungen war mit YOLO V8 schnell gefunden. Nun galt es, das knifflige Problem der Trainingsdatengewinnung zu lösen. Mit bloßem Auge eine ausreichende Anzahl von Varroamilben-befallenen Bienen im eigenen Videomaterial zu finden und zu extrahieren, erschien nahezu unmöglich.

## Trainingsdaten: Säen und ernten

Mit dem Sichten seiner Videos hätte Steppuhn sehr, sehr viele mühsame Stunden verbringen müssen. Insbesondere wenn nur wenige Bienen befallen sind, kann es lange dauern, bis das gesuchte Muster auftaucht. Und in dem Moment braucht man wahre Super-Recognizer-Fähigkeiten, um den winzigen, gut getarnten Parasiten schnell genug zu erkennen, denn jede Biene huscht nur Bruchteile von Sekunden lang durchs Bild. Weil Steppuhn weder Super-Recognizer ist noch zu viel Zeit hat, dafür aber gut

Bild: Ale Zee / BWVL



**Mit seinem Milbendetektor BeeAI gewann der 17-jährige Sebastian Steppuhn den Bundeswettbewerb künstliche Intelligenz 2024.**

## Reflexionsarme Fotofalle

Wer in Echtzeit die Milbe im Bienenpelz finden will, kann nicht einfach eine Kamera vor dem Eingang des Bienenstocks platzieren und auf Treffer hoffen. Um der KI das Lernen und anschließende Erkennen zu erleichtern, benötigt man eine möglichst kontrollierte, normierte und witterungsgeschützte Umgebung, die Nahaufnahmen in hoher Qualität garantiert. Im Vorfeld konstruierte und testete Steppuhn also mehrere Varianten für ein solches Kameramodul, bis er mit dem Ergebnis zufrieden war.

Wie eine gut analysierbare und gleichermaßen tierfreundliche Fotofalle aussehen sollte, erarbeitete sich der Nachwuchsimker durch eigene Experimente sowie in Gesprächen mit Experten des Bieneninstituts Hohenheim. Mithilfe von sogenannten Bienenfluchten ermittelte er etwa, dass ein Kanal idealerweise 8 Millimeter Durchmesser haben sollte, damit die Bienen bequem in ihren Stock gelangen. Außerdem fand er heraus, dass sich die Varroamilbe vor allem auf deren Rücken und Brust anheftet. Deshalb filmt sein auf einem Raspberry Pi 4 laufendes Aufzeichnungssystem den Durchgang von oben und von unten, und zwar jeweils mit einer Arducam in 8-Megapixel-Auflösung.

Als Abdeckung für die Ober- und Unterseite des Kanals wählte er schlussendlich Objektträger aus dem Mikroskopiebedarf, weil sie die besten optischen Eigenschaften haben und sich bei Verschmutzung kostengünstig austauschen lassen. Plexiglas und andere durchsichtige Materialien erwiesen sich in Tests als ungeeignet: Sie spiegelten zu stark oder waren nicht gut zu verbauen. Sämtliche Versionen seiner Kameramodule entwickelte er mit dem CAD-Programm FreeCAD. Die Seitenwände bestückte er mit LEDs, die in einem halbtransparenten 3D-Druckteil eingelassen sind. Die Halbtransparenz bewirkt eine gleichmäßigere Ausleuchtung des Gangs.

Allein dieser Aufbau zeigt, wie analytisch und akribisch ein Machine-Learning-Projekt häufig vorbereitet sein will beziehungsweise wie viel fachspezifisches und interdisziplinäres Wissen einfließen muss, um geeignete, hochwertige Daten zu generieren. Selbst kommerziellen oder von Wissenschaftlern durchgeführten KI-Projekten mangelt es mitunter an dieser Sorgfalt oder der notwendigen Expertise – ein Grund, weshalb viele KIs mit Trainings- und Testdaten kaum Fehler machen, aber im praktischen Einsatz scheitern.

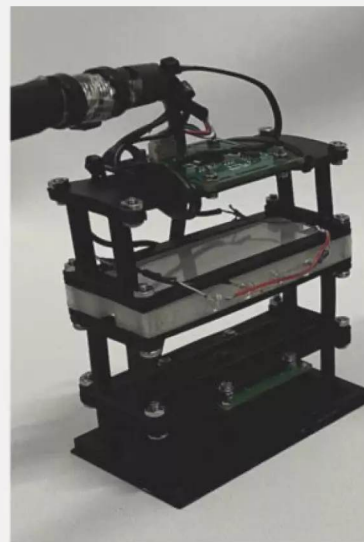
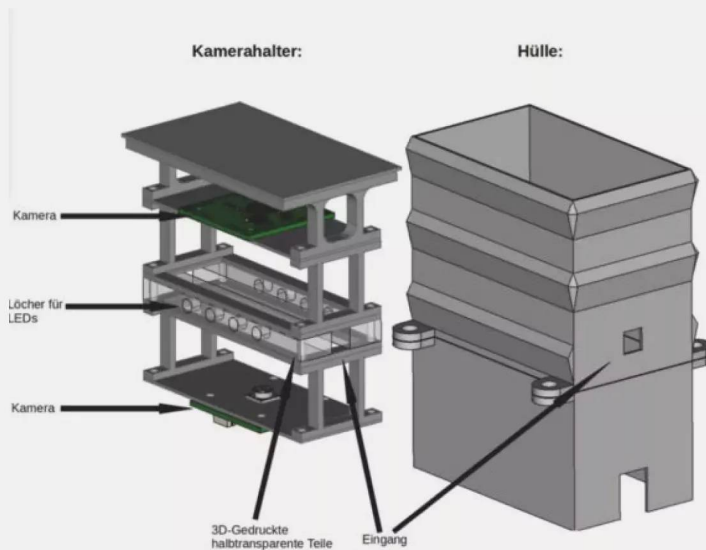
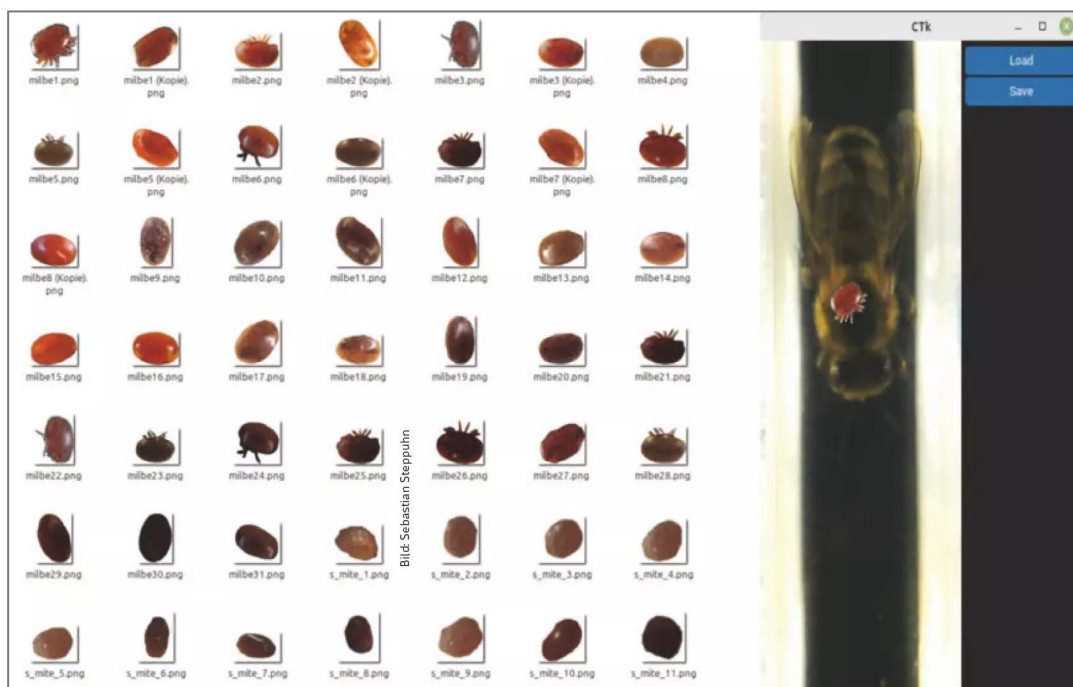


Bild: Sebastian Steppuhn

**Selbst ausgetüftelt, konstruiert, gedruckt und montiert: die perfekt LED-ausgeleuchtete Bienenschleuse mit Kameramodulen.**

**Kleine Milbensammlung für die Fotomontage: Anfangs setzte Steppuhn die Parasiten noch händisch auf die Bienenbilder, rasch ging er per Python-Skript zur Massenproduktion über.**



programmieren kann, bediente er sich eines Tricks: Er generierte zunächst eine synthetische Datensammlung, indem er Milben- auf Bienenbilder montierte und diverse Variationen daraus fertigte. Mit diesen Beispielen wiederum trainierte er das YOLO-System, um sukzessive mehr und mehr Frames mit echten Darstellungen befallener Bienen in seinen Videos zu lokalisieren.

Der erste Part war vergleichsweise einfach: Bienen im Videostream erkennen und markieren, darin die „Region of interest“ (Brust oder Rücken) abgrenzen und mit einem Label versehen, mit diesen Exemplaren das YOLO trainieren und weitere gelabelte Bienenbilder erhalten. Fertig war eine vielfältige Bienensammlung als Basis für die Fotomontage.

Unterschiedliche Abbildungen von Varroamilben fand Steppuhn im Internet und stellte sie mit der Open-Source-Bildbearbeitung GIMP frei. Die isolierten Milbenbilder montierte er zunächst händisch auf Brust und Rücken seiner Basis-Bienenfotos und labelte sie ebenfalls händisch mit dem freien Bild-annotations-Tool Labellmg. Rasch war klar: Auch die Fotomontage muss automatisiert ablaufen, sonst ist der Aufwand nicht zu stemmen. Denn um ein

Milben-Bienen-Gespann unter allen möglichen, auch ungünstigen Bedingungen zu erkennen, benötigt man einen möglichst facettenreichen Trainingsdatensatz. Der muss auch knifflige Fälle enthalten, in denen der winzige Parasit besonders schwer zu erkennen ist, etwa weil er sich kaum vom Hintergrund abhebt, teilweise verdeckt ist oder wegen eines ungünstigen Aufnahmewinkels verzerrt erscheint.

Solche Variationen lassen sich ebenfalls synthetisch erzeugen. Für dieses sogenannte Augmentieren schrieb der Schüler sein eigenes Python-Skript, das die Milben zunächst in Größe und Rotation veränderte sowie zufallsgesteuert auf Brust oder Rücken der Bienen platzierte. Damit erstellte er innerhalb von fünf bis 15 Minuten komplette Datensätze mit durchschnittlich 3000 Trainingsbildern.

Der erste Durchlauf mit 1000 Trainings- und 252 Validierungsbildern verlief allerdings ernüchternd, weil das Modell echte Milben kaum erkennen konnte. Den Durchbruch brachte eine verfeinerte Montagetechnik, wie Steppuhn gegenüber c't erläutert: „Ich habe die Datenpipeline mehrfach überarbeitet und optimiert, bis die Ergebnisse besser wurden. Eine der wichtigsten Verbesserungen war die Variierung



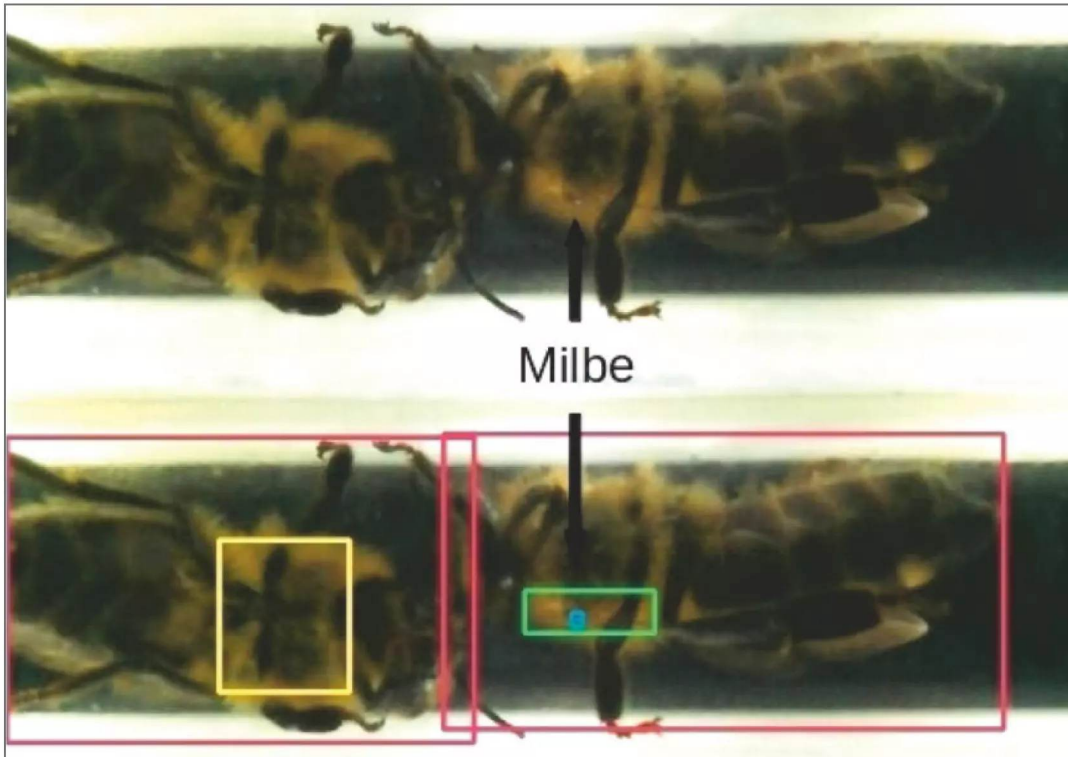


Bild: Sebastian Steppuhn

**Der Lohn für monatelanges Experimentieren und Optimieren: die erste Sichtung einer echten Varroamilbe im Testbienenstock.**

der Transparenz der Milbenbilder, wodurch diese auf dem Bienenbild realistischer aussahen.“

Auf Grundlage dieser wertvollen realistischen Trainings-Samples konnte er weitere Variationen erzeugen und Verfeinerungen vornehmen, bis er schließlich einen umfangreichen Datensatz beisammen hatte, der alle möglichen Winkel und Erscheinungsformen des Milben-Bienen-Gespanns abdeckte. Eine geschickt gewählte Kombination aus Glas und Beleuchtung für den Durchgang gewährleistete störungs- und insbesondere reflexionsarme Aufnahmen.

Der fertig trainierte Milbendetektor geht dreistufig vor: Erst markiert er die Biene, dann deren Rücken oder Bauch, anschließend sucht er die verräterischen Muster innerhalb dieser Region of Interest – ähnlich wie es übrigens auch die Profis bei der Echtzeit-Kennzeichenerkennung machen, wo der Suchbereich zunächst auf das Fahrzeug und dann auf das Nummernschild eingegrenzt wird.

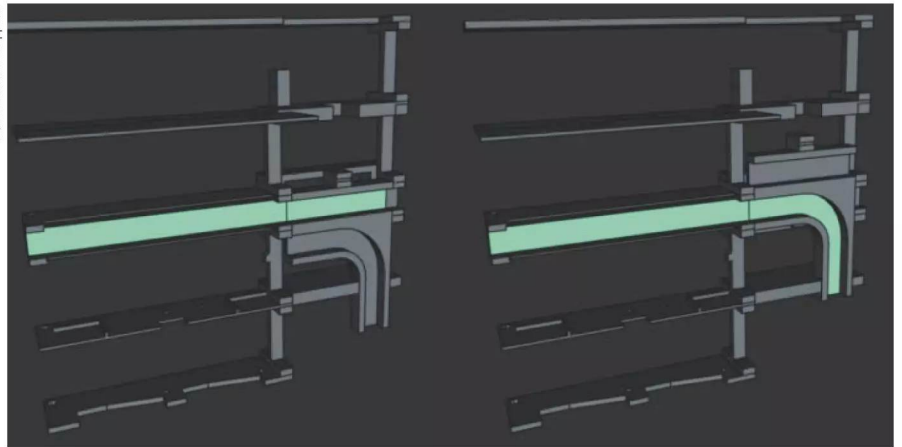
## Sanft ausgepustet

Befallene Bienen zu erkennen, ist prima, aber ohne geeigneten Sortiermechanismus bleibt es eine akademische Übung. Auch bei der Lösung dieses Problems konnte der 17-jährige Schüler nicht auf bewährte Verfahren zurückgreifen. Zwar gibt es bereits Maschinen im industriellen Maßstab, die pilzbefallene Getreidekörner erkennen und über ein starkes Gebläse aussortieren. Solche Hochgeschwindigkeitsverfahren eignen sich für die empfindlichen Insekten aber eher nicht; man will sie ja behandeln und nicht umbringen. Deshalb konzipierte Steppuhn kurzerhand einen Druckluftmechanismus, der die Patienten sanft in einen separaten Behälter umleitet.

Die grün markierten Bereiche sind die Wege, die die Bienen gehen können. Dabei ist der hintere Teil mit den zwei möglichen Wegen beweglich. Sobald das System eine Milbe erkennt, wird zuerst mit einem kurzen Druckluftimpuls der Bereich des beweglichen

Wer milbenbefallene Bienen erkennt, muss sie auch automatisiert separieren: zum Beispiel mit diesem insektenfreundlichen Mechanismus. Gesunde Bienen dürfen den Eingang passieren (links), Milbenträger lösen die Weiche aus und werden zur Behandlung in einen separaten Behälter bugsiert (rechts).

Bild: Sebastian Steppuhn



**KI-gesteuerte Besucherlenkung:**  
Gesunden Bienen bleibt der Weg zum Stock geöffnet, milbenbefallene müssen ins Behandlungszimmer abbiegen.

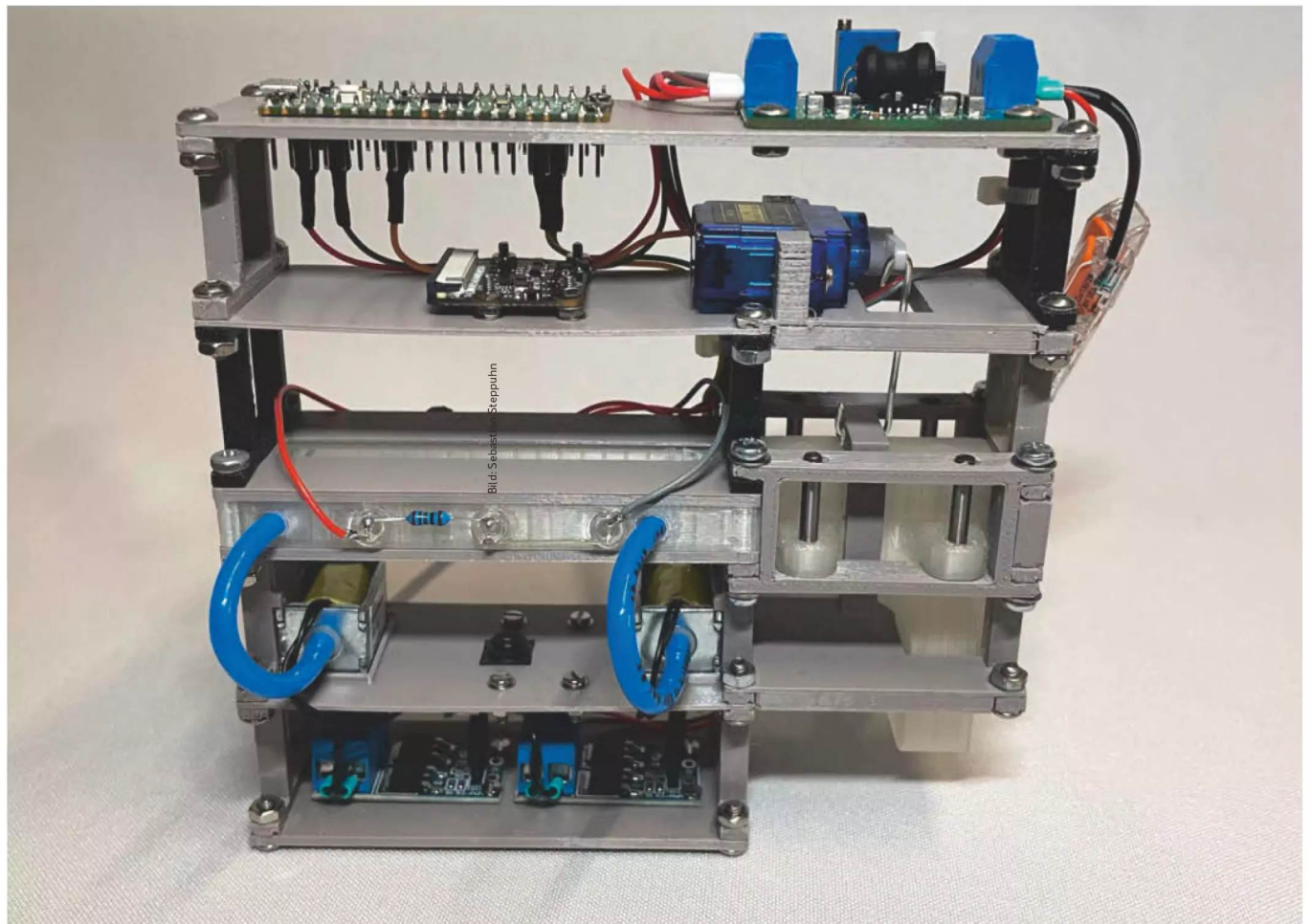


Bild: Sebastian Steppuhn

## Bundeswettbewerb KI: Clevere Ampel, mutige Experimente

Der bereits zum sechsten Mal vom Tübingen AI Center ausgerichtete BWKI stand in diesem Jahr unter dem Motto „Welche Freiheit gibt dir KI?“. Auch die Preisträger der anderen Kategorien punkteten mit funktionierenden und zu Ende gedachten Systemen: Leonie Weiss aus Regensburg wurde für ihre intelligente Ampelschaltung TrafficAid mit dem Preis der Sonderkategorie „AI for Good“ ausgezeichnet. Das System erkennt Fahrzeuge aller Art und schaltet nur dann auf Rot, wenn es wirklich notwendig ist – womit unnötige Wartezeiten vermieden und CO<sub>2</sub>-Emissionen reduziert werden könnten. Bei der Steuerung berücksichtigte die 18-Jährige auch Sicherheits- und Fairness-Aspekte, damit die jeweiligen Wartezeiten nicht allzu weit auseinanderklaffen, wenn aus einer Richtung sehr viele Fahrzeuge kommen und aus der anderen nur wenige.

Geradezu prädestiniert für die Kategorie „No risk, no fun!“ war das Projekt „KI æmazing“ von Anna Perkovic und Nicholas Dahlke aus dem baden-württembergischen Lörrach. Die beiden 17-Jährigen hatten sich zum Ziel gesetzt, ein noch ungeklärtes wissenschaftliches Rätsel mithilfe von KI zu lösen: Warum gefriert heißes Wasser schneller als kaltes? Was in der Projektbeschreibung recht harmlos klingt, entpuppte sich als äußerst ambitioniertes Experiment mit sehr aufwendigem Versuchsaufbau. Auf der einen Seite einer Röhre wurde Wasser mit Öl zusammen eingespritzt,

sodass sich 0,5 Millimeter große, mit einem Ölfilm umhüllte Tröpfchen bildeten. Weiter ging es durch eine minus 40 Grad kalte Kühlanlage, dahinter klassifizierte eine zuvor trainierte Bilderkennung jedes Tröpfchen: gefroren oder nicht. Tatsächlich stellte sich heraus, dass bei der Verwendung von heißem Wasser deutlich mehr Tröpfchen gefroren als bei kaltem Wasser. Zugute kam den Nachwuchsforschern dabei sicherlich die Ausstattung des Schülerforschungszentrums Phaenovum in Lörrach, aber auch ihr Organisationsgeschick: Diverse Spezialgeräte liehen sie sich bei Firmen.

Der Sonderpreis „KI-Forschung“ würdigte die Arbeit von Lorenz Rutkevic aus dem niedersächsischen Leer. Sein System optimiert und analysiert mikroskopische Aufnahmen von Zellen. Zum einen sollen damit auch weniger präzise digitale Mikroskope gestochen scharfe Aufnahmen liefern. Zum anderen unterstützen die von Rutkevic trainierten und angepassten Algorithmen Ärzte bei der Diagnose von Krankheiten, indem sie nach typischen Mustern suchen und diese in der Aufnahme markieren beziehungsweise segmentieren. Auch der Publikumspreis ging nach Niedersachsen, nämlich an den „FolderCopter“ des 16-jährigen Peter Fuchs aus Hannover: ein lokal auf dem Rechner laufendes Sprachmodell, das Fragen zur eigenen Dokumenten- und Wissenssammlung beantwortet.

Teils freigepustet, sodass bei der Bewegung keine Bienen zerquetscht werden und somit keine Bienen-guillotine entsteht. Anschließend fährt der hintere Teil hoch, sodass der Eingang zum Bienenstock verspermt und ausschließlich der Weg nach unten passierbar ist. Nun wird die Biene, die im Kanal sitzt, mit einem zweiten Druckluftstoß nach unten gepustet.

Der auf einer Geforce-RTX-3060-Grafikkarte trainierte und arbeitende Milbendetektor benötigt für seine Entscheidung „etwa 10 Millisekunden pro Bild, sodass der Mechanismus sehr schnell ausgelöst werden kann“, erläutert Steppuhn gegenüber c't. „Das Auspusten läuft anschließend in unter einer Sekunde ab.“ Dafür hat Steppuhn den heimischen Garten mit High-Speed-Internet angebunden. Der



Raspi 4 bietet nicht genügend Leistung für eine Echtzeiterkennung. Für den kommenden Sommer hat er sich vorgenommen, mit einem Gespann aus Raspi 5 und angeschlossener Sony-KI-Kamera als Basis für die Bilderkennung zu experimentieren.

## Fehlertolerantes System

Ob und wie gut ein maschinell trainiertes System funktioniert, verraten diverse Fehlermetriken. Bei der Varroamilbenerkennung handelt es sich um ein klassisches Massen-Screening, dem ein systembedingtes Problem innewohnt: Wenn das gesuchte Phänomen sehr selten auftritt, kommt zwangsläufig auf eine korrekte Erkennung ein Vielfaches an Falsch-Positiven. So auch hier: Das über mehrere Monate in 50 Durchgängen optimierte System konnte am Ende zuverlässig Milben erkennen, die das menschliche Auge nicht wahrnahm. In einem Testlauf mit 12 Stunden Videomaterial sortierte es von 25.671 Bienen, die den Kanal passiert hatten, ganze 83 aus – von denen bei manueller Inspektion aber nur 16 tatsächlich eine Milbe im Pelz hatten.

Das klingt nach einem gescheiterten Experiment. Entscheidend sind aber nicht die nackten Zahlen, sondern die Nutzen-Schaden-Rechnung: Um Schaden vom Bienenvolk abzuwenden, muss die Milbe zuverlässig erkannt werden, damit möglichst wenige von ihnen in den Stock gelangen. Die vielen

Falsch-Positiven hingegen fallen bei dieser Anwendung nicht so sehr ins Gewicht, denn ohne den Detektor müsste man über 25.000 Bienen einer unnötigen Säurebehandlung aussetzen, mit dem Detektor trifft es nur 67 Individuen.

Wie hoch die Erkennungsleistung ist, also wie viele der insgesamt befallenen Bienen der Algorithmus aufspürt, lässt sich derzeit nicht beziffern. Denn dazu müsste der exakte Befall mit etwas aufwendigeren Methoden gemessen werden. Aus Kennzahlen wie der Größe und der Brut seines Volkes, die Steppuhn regelmäßig kontrolliert, konnte er zumindest ableiten, dass es sich um eine gesunde Population mit nicht allzu hoher Befallsquote handelt. Im Fokus stand diese Frage nicht, denn primär ging es dem Tüftler erst einmal darum, nachzuweisen, dass eine KI-Erkennung von Milben möglich ist.

Alles in allem ist es dem 17-jährigen Sebastian Steppuhn gelungen, ein robustes, wetterfestes System mit einer von Anfang bis Ende und ins Detail durchdachten KI-Pipeline auszutüfteln. Wer sich davon oder von den anderen prämierten Projekten nun inspiriert fühlt: Auch 2025 findet der KI-Bundeswettbewerb wieder statt, den Einstieg ins Machine Learning erleichtert der KI-Kurs des Tübinger AI Center.

heise Medien war 2024 Kooperationspartner des BWKI. c't-Redakteurin Andrea Trinkwalder war in diesem Jahr Teil der Jury. (atr) **ct**

Weitere Infos zu BeeAI und zum BWKI:

[ct.de/wktn](https://ct.de/wktn)



**Jetzt absichern!**

**ct SECURITY-EINSTIEG**  
Mit wenigen Handgriffen zu mehr Sicherheit

Passkeys besser als Passwörter  
Schnell und einfach weg vom Passwort  
Phishing und Leaks verhindern

Keine Chance für Cybercrime  
Online-Betrüger brechen und ärgern  
Kriminelle Banden hunderte aushebeln

Extra-Sicherheit für Windows  
Schützliche Software automatisch aus  
Open Security-Tools für noch mehr Schutz

14 Security-Checklisten  
Für alle Fälle  
In wenigen Minuten zu mehr Sicherheit  
Social Media Konten, Homeoffice, Smart  
Windows & Co. absichern

**Vorsicht ist besser als Nachsicht!**

- ▶ 14 Security-Checklisten für alle digitalen Bereiche
- ▶ Extra-Sicherheit für Windows
- ▶ Cybercrime verstehen und bekämpfen
- ▶ Passkeys gegen Phishing und Leaks

... und noch viel mehr hilfreiche Artikel

**Gleich reinschauen**

 [shop.heise.de/ct-securityEinstieg24](https://shop.heise.de/ct-securityEinstieg24)





Bild: OpenAI

# Als die KI-Bilder laufen lernten

Erst vor kurzem revolutionierten generative KI-Dienste wie Dall-E, Midjourney und Stable Diffusion die Art der Bilderzeugung. Jetzt lernen die Bilder laufen. Bisher gibt es nur wenige Anbieter, und die bringen nur kurze Clips zustande. Die Ansätze sind aber vielversprechend. Eine Bestandsaufnahme.

Von **André Kramer**

**D**ie oben gezeigte Frau im roten Kleid erregte im Februar 2024 bei ihrem Spaziergang durch Tokio viel Aufmerksamkeit. Das hier im Screenshot festgehaltene Video wurde nicht etwa gefilmt, sondern entstand durch Eingabe eines Prompts in einen KI-Videogenerator. Das KI-Unternehmen OpenAI stellte auf diese Weise seine experimentelle Software Sora vor. Öffentlich nutzbar ist sie in Deutschland bisher nicht.

Die kanadische Indie-Pop-Band Shy Kids durfte schon vor dem offiziellen Start mit Sora arbeiten. Der dabei entstandene 1:20 Minuten lange Kurzfilm *Air Head* erzählt in mehreren zusammengeschnittenen Clips die Geschichte eines Mannes, der einen gelben Luftballon anstelle seines Kopfes hat. Außer Realfilmaufnahmen enthält der Film auch KI-generierte Clips, siehe [ct.de/wggw](https://ct.de/wggw). Er zeigt, welches Potenzial in den KI-Videogeneratoren steckt.



Mittlerweile gibt es einige Webdienste, die sich kostenpflichtig abonnieren lassen und kommerziell verwertbare Ergebnisse versprechen. Langfristig sollen die KI-Modelle Stockvideos überflüssig machen. Ähnlich exerzieren es Midjourney, OpenAI Dall-E, Stable Diffusion und Adobe Firefly bereits im Fotobereich vor [1]. Denn es führt oft schneller zum Ziel, das gewünschte Bild oder Video detailliert zu beschreiben und von künstlicher Intelligenz generieren zu lassen, als eine Bibliothek mit vorgefertigtem Material in allen möglichen Stilen zu durchstöbern.

Bisher sind KI-Videodienste rar gesät. Sie stecken noch in den Kinderschuhen; wie die Erfahrung mit Bildgeneratoren gezeigt hat, dürften sie sich aber zügig weiterentwickeln. Drei Dienste im Web erzeugen mindestens zehn Sekunden lange Clips aus reinen Textprompts und richten sich an Nutzer, die diese kommerziell verwenden wollen: An unserem Test nehmen der britische Videogenerator Haiper, Luma Dream Machine aus Kalifornien und das Tool des New Yorker Unternehmens Runway ML teil. Das Videoschnittprogramm Adobe Premiere Pro erweitert Video-Clips ebenfalls, jedoch nur um maximal zwei Sekunden. Ferner gibt es Apps, die Selfie-Fotos in lächelnde oder tanzende GIFs wandeln. Sie werden zuhauf auf TikTok und Instagram beworben und erzeugen kurze animierte Schnipsel des geladenen Selfies: zu wenig, um sich als Testkandidat zu qualifizieren.

## Eingabe und Ausgabe

Video-KIs präsentieren dem Anwender ähnlich wie Bild- und Textgeneratoren zunächst nur eine schlichte Promptzeile. Hier beschreibt man in einem kurzen Text oder auch mit einem Bild als visuelle Starthilfe, wie der Videoclip aussehen soll. Ein reiner Textprompt ist die Königsdisziplin für Anwender und KI, denn in diesem Fall muss die Software ohne Referenz bei null beginnen. Der Prompt dazu sollte neben dem dargestellten Motiv und dem bildgebenden Stil auch das Seitenverhältnis enthalten: In der Regel können Nutzer aus verschiedenen Standardformaten wie 16:9, 4:3 und 1:1 sowie den korrespondierenden Hochformaten wählen – im Prompt oder grafisch zum Anklicken.

Alternativ lädt man ein Bild im JPEG- oder PNG-Format hoch. In diesem Fall orientiert sich der Videogenerator an dessen Seitenverhältnis. Weicht es von den unterstützten Maßen ab, passt die Software die Proportionen entsprechend an. Auch beim Motiv muss die KI nicht rätseln. Was mit ihm geschehen soll, kann eine Textzeile spezifizieren – oder man überlässt es der KI, daraus eine Szene zu drehen. Auf diese Weise entstanden zum Beispiel animierte Versionen bekannter Memes wie des „Disaster Girl“ oder des „Distracted Boyfriend“ ([ct.de/wggw](https://ct.de/wggw)). Einige Modelle akzeptieren auch Videos als Eingabequelle oder ergänzen zwischen zwei Bildern die passenden Frames, etwa um Drohnenaufnahmen zu einem Videoclip zu verbinden.

Alle drei Testkandidaten generieren Videos standardmäßig in 720p-Auflösung. Haiper kann sie auf 1080p aufblasen, Runway ML sogar auf 4K. Bei Luma muss man mit der ausgegebenen Auflösung vorliebnehmen. Die fertigen Clips speichern sie in der Cloud des jeweiligen Anbieters. Bei Runway ML ist deren Speichervolumen abhängig vom gebuchten Tarif, bei den übrigen Anbietern nicht. Alle verstehen Prompts in Englisch oder Deutsch; die Oberflächen der jeweiligen Web-Apps sind aber auf Englisch gehalten.

## Stoff für Albträume

Videogeneratoren müssen deutlich mehr leisten als die auf Standbilder spezialisierten Modelle Midjourney, Dall-E und Co, die Grundanforderungen sind aber gleich: Ihre Aufgabe ist keine geringere, als eine korrekte Repräsentation von Dingen und Menschen zu generieren. Dazu gehören beispielsweise anatomische Details wie fünf Finger an jeder Hand und nicht mehr als zwei Arme pro Körper. Die Geh-

Bild: Shy Kids



**Der Kurzfilm Air Head der kanadischen Popband Shy Kids entstand in Teilen mit dem KI-Videogenerator OpenAI Sora.**



## So funktioniert OpenAIs Videogenerator Sora

OpenAIs Videogenerator Sora soll insgesamt eine Minute hochauflösendes Videomaterial generieren können, laut Hersteller in hoher visueller Qualität. Das ist deutlich länger als bei den bereits veröffentlichten Konkurrenten. Den Beweis muss OpenAI allerdings erst erbringen, denn die KI-Szenen im Kurzfilm „Air Head“ dauern jeweils nur einige Sekunden und selbst ausprobieren kann man den Dienst bisher nicht. Neben Videos akzeptiert Sora laut OpenAI auch Fotos, Grafiken oder KI-generierte Bilder als Eingabe. Videodateien kann die KI sowohl vorwärts als auch rückwärts ergänzen. Außerdem soll sie zwei Videos verbinden und dabei einen nahtlosen Übergang schaffen können.

Für Sora trainiert OpenAI Diffusionsmodelle mit Videodaten unterschiedlicher Dauer, Auflösung und mit verschiedenen Seitenverhältnissen. Die Trainingsdaten werden bei Diffusionsmodellen durch Zugabe von Rauschen sukzessive verfremdet. Das Modell lernt anhand solcher Sequenzen, aus reinem Bildrauschen im Idealfall das ursprüngliche, mindestens aber ein ähnliches Bild wiederherzustellen. Dank dieser Technik gelingt es modernen Bild-KIs wie OpenAIs Dall-E, qualitativ hochwertige Bilder aus Textprompts zu generieren.

Der Hersteller selbst bezeichnet Sora als „Diffusion Transformer“. Google nutzte als erster eine Transformer-Architektur, um Text in numerische Darstellungen zu wandeln und auf diese Weise Wörter und Wortschnipsel semantisch zu sortieren. Die Technik verhalf den großen Sprachmodellen zum Durchbruch, insbesondere jenen von OpenAI: dem Basismodell GPT und dem Chatbot ChatGPT. Für Sora wendet das

KI-Unternehmen die Transformer-Methode auf Videokacheln an; die Entwickler bezeichnen sie etwas futuristisch als „Raum-Zeit-Kacheln“. Das sind ähnliche Schnipsel samt Textbeschreibung (Token), wie sie beim Training von Bildgeneratoren verwendet werden. Nur haben die Videoschnipsel außer einer räumlichen auch eine zeitliche Verbindung zu einander. Der Transformer bildet die Pixel-Text-Paare auf Vektoren ab, die einen hochdimensionalen Vektorraum aufspannen, den sogenannten Latent Space. Semantisch ähnliche Kacheln, die in einem Bildmotiv oder einer Videoszene häufig im Kontext auftauchen, liegen dort enger beieinander, artfremde weit auseinander. Diese Ordnung hilft Bild- und Videogeneratoren, stimmige Motive sowie Abläufe zu komponieren.

Den Prompt übersetzt OpenAI auf dieselbe Weise wie bei Dall-E mithilfe seiner eigenen Large-Language-Modelle der GPT-Familie. Kurze Nutzereingaben wandelt das LLM in detailliert ausformulierte Prompts, die es an das Videomodell übergibt.



**Ein Diffusionsmodell wird mit schrittweise verrauschten Bildern trainiert und lernt, aus reinem Bildrauschen Inhalte zu erzeugen.**



Die junge Frau links erbricht eine Bierflasche und transformiert ihren Kopf, während der junge Mann selbst zur Frau wird. „Stoff für Albträume“ nennt das Internet solche KI-generierten Videos. Ihnen mangelt es an Konsistenz.

versuche der Bild-KIs in den letzten zwei Jahren haben bewiesen, dass das alles andere als trivial ist. Das „deutlich mehr“ kommt beim Thema Video durch die zeitliche Dimension zum Tragen. Denn es ist zwar nicht völlig unmöglich, eine Folge von Bildern zu erzeugen, die einander so weit ähneln, dass man sie zu einem Video zusammenkleben kann, ohne dass sichtbare Sprünge entstehen. Foren und Profileiten auf Reddit, Instagram und Facebook sind voll davon. Zusätzlich müssen aber auch die Physik, die Bewegungen und die Logik stimmen.

Dabei geht derzeit noch mit ziemlicher Sicherheit eine Menge schief: Eine Frau sitzt neben einem Mann am Strand, zieht eine Bierflasche aus der Nase, ändert in einer Bewegung Frisur sowie Gesicht und verliert einen Arm, aus dem Mann wird eine Frau, aus seinem Hemd ein Kleid. „Fuel for Nightmares“, Stoff für Alpträume, nennt die Szene das. Und sie hat sichtlich Spaß daran, die generative KI mit einer Abscheulichkeit nach der anderen zu beauftragen. Das Strandvideo sorgt für Schenkelklopper auf Reddit. Kommerziell nutzen, etwa für Werbe- und Imagevideos eines Tourismusverbands ostfriesischer Inseln, lässt sich die Szene nicht.

Unter anderem weil die derzeit verfügbaren Videogeneratoren so schnell thematisch abdriften, also oftmals am Ende völlig andere Personen und Dinge zeigen als am Anfang, erzeugen sie nur vier oder fünf Sekunden lange Clips. Denn bei maximal 150 aufeinanderfolgenden Bildern kann nicht ganz so viel schiefgehen. Die Clips lassen sich zwar um nochmal dieselbe Länge erweitern, aber in der Regel nur einmal. Viele KI-Videos im Netz wurden durch einen simplen Trick deutlich verlängert: Deren Produzenten nutzen einfach das letzte Frame des vorherigen Videos als Eingabe, um auf dessen Grundlage ein weiteres Video zu erzeugen. Die fertigen Clips schneiden sie anschließend zusammen.

## Konsistenz in Raum und Zeit

Für kommerzielle Zwecke müssen die Dienste also in der Lage sein, über einen möglichst langen Zeitraum Abhängigkeiten und Interaktionen zwischen Gegenständen, Menschen und Tieren zu modellieren. Einem Burger, von dem jemand abbeißt, muss anschließend das abgebissene Stück fehlen. Ein Spaziergänger am Strand muss Fußspuren im Sand, ein Maler sichtbare Pinselstriche auf der Leinwand hinterlassen. Bei den Testkandidaten ist es mit der Konsistenz in Raum und Zeit derzeit noch nicht weit her.



**Die Videodienste sollen das Ursprungsbild (oben links) erweitern. Haiper (oben rechts) verzerrt die Gesichtszüge und ergänzt zu viele Finger. Luma (unten links) gibt sich dynamisch, nimmt es aber mit der Anatomie nicht so genau. Runway (unten rechts) animiert das Bild ohne Fehler oder visuelle Dynamik.**

Den Videogeneratoren fehlt komplexes mathematisches Wissen, etwa in Form einer Physiksimulation, die 3D-Programme wie Blender oder die Unreal Engine mitbringen. Was mit einem 3D-Modell spielend gelingt, stellt einen Videogenerator vor nahezu unlösbare Aufgaben: Bei dynamischen Kameraschwenks oder auch der Rotation um eine Szene mit Drohnen- und Helikopterflügen muss der Videogenerator Landschaften, Personen, Gebäude und andere Elemente dreidimensional korrekt rekonstruieren können. Noch anspruchsvoller wird es bei physikalisch komplexen Vorgängen wie zerbrechendem Glas. Das kann derzeit noch keine kommerziell verfügbare generative künstliche Intelligenz abbilden.

Die generative künstliche Intelligenz Kling AI kann laut ihres chinesischen Entwicklers Kuaishou





Das Surfermotiv (oben links) erweitert Haiper (oben rechts) diesmal erfolgreicher als Luma (unten links). Beiden unterlaufen dabei etliche Konsistenzfehler. Runway ML (unten rechts) führt es wiederum stoisch in Zeitlupe weiter.

**„A redneck, wearing a shirt with an american flag, eats a pile of hot dogs in an eating contest“. Von oben: Haiper, Luma, Runway**



Videos mit einer Länge von bis zu zwei Minuten und einer Auflösung von 1080p bei 30 Bildern pro Sekunde generieren. Sie ist derzeit als Demoversion ausschließlich in China verfügbar und verwendet nach Angaben der Entwickler ein 3D-Raum-Zeit-Aufmerksamkeitssystem, um Bewegungen konsistent zu modellieren. Laut Kuaishou simuliert das Modell physikalische Eigenschaften der realen Welt. In einem Beispielvideo beißt ein Junge von einem Burger ab, dem dieser Bissen im weiteren Video fehlt. Auch OpenAI präsentiert solche Videos. Kling AI nutzt dafür wie der Konkurrent von OpenAI einen „Diffusion Transformer“ (siehe Kasten über OpenAI Sora).

## Darstellbare Inhalte

Anders als Sora und Kling AI mit ihren ambitionierten Zielen sind Haiper, Luma und Runway bereits nutzbar. Sie offenbaren dabei unterschiedliche Charaktere. Runway ML setzt ganz auf kommerzielle Anwender. Die Videos erinnern oftmals an Zeitlupeneffekte, die wenig mehr machen als ein hochgeladenes Bild zu animieren oder ein zum Textprompt passendes, animiertes Stillleben zu zeichnen. So bleiben Details, Gesichter, Kleidung und Gegenstände erhalten, aber es passiert auch nicht viel im Video.

Haiper beschreitet mit seiner Betaversion einen ähnlichen Weg. Der Dienst generiert zwar mehr





## Haiper (Beta)

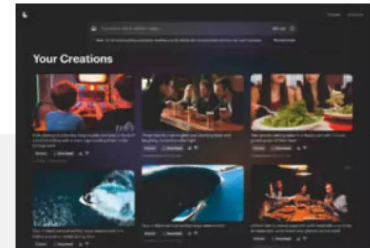
Das Londoner Unternehmen Haiper bietet seinen gleichnamigen Webdienst bisher nur als Betaversion an. Mit einer kostenlosen Testversion lassen sich 300 Videos erstellen, die zügig vorliegen, allerdings Wasserzeichen einblenden und nicht für kommerzielle Nutzung oder als Download freigegeben sind. Alle Videos sind zudem öffentlich einsehbar. Für 20 US-Dollar im Monat bekommt man 5000 Credits, die am Ende des Monats nicht verfallen. Außerdem speichert Haiper die Videos in einem privaten Bereich und entfernt das Wasserzeichen.

Im Einstellungs Menü lassen sich 2, 4 oder 8 Sekunden sowie geringe und hohe Auflösungen wählen. Acht-Sekunden-Clips kosten 40 Credits, viersekündige sind kostenlos. Außerdem kann man das Seitenverhältnis angeben. Haiper setzt hochgeladene Bilder als Video fort oder fertigt neue Frames aus Texteingaben – mit gemischtem Resultat.

Unsere Promptvorgaben setzte Haiper detailreich und mit lebhaften Farben um. Die aktuell in Version 1.5 vorliegende KI erstellt, wenn nicht anders angegeben, ruhige Szenen ohne starke Kamerabewegung. Sie kann auch dynamische Videos generieren, die sich allerdings schon nach wenigen Frames so stark vom Ursprungsmotiv entfernen, dass abgebildete Personen und Objekte nicht mehr wiederzuerkennen sind. Als einziger Dienst im Test kann Haiper auch Videos mit bekannten Charakteren erzeugen, sei es R2-D2 oder Donald Trump. Schwächen zeigt er bei essenden Personen und schnellen Bewegungen etwa von Sportlern.

👉 detailreiche, farbenfrohe Videos

👎 Probleme mit der Anatomie



## Luma Dream Machine

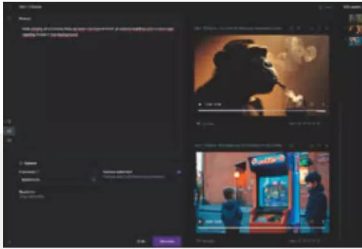
Am 14. Juni 2024 hat die Firma Luma AI den KI-Dienst Dream Machine veröffentlicht, der es schnell zu Popularität im Internet gebracht hat. Das muss aber nichts Gutes heißen: Luma bietet ein kostenloses Konto an, das mit 30 Credits im Monat ebenso viele Videos generiert. Im Standardkonto kommen 120 Credits hinzu. Im Abo darf man Videos ohne Wasserzeichen herunterladen und kommerziell nutzen. Außerdem arbeitet Luma im Bezahlkonto deutlich schneller: Wer einen Gratis-Account hat, muss zuweilen fast einen ganzen Tag auf seinen fertigen Clip warten. Weitere Credits lassen sich hinzubuchen.

Die Eingabemaske ist einfach gehalten: Nutzer dürfen ein Bild hochladen, Text eingeben oder beides. Bei Fotos von Personen bleibt Luma enger am Original als Haiper, animiert allerdings auch Gabelzinken sowie Logos und Schriftzüge auf T-Shirts, Flaschen oder Gläsern. Clever: Luma drückt sich vor Szenen essender Personen, denn das geht auch bei anderen Diensten meist schief. Personen tendieren bei Luma zudem dazu, Flaschen zu erbrechen und Gläser zu verschlucken. Immerhin entgleisen die Gesichtszüge in kurzen Clips weniger stark als bei Haiper. Bei Kamerafahrten um eine Szene erhält die Dream Machine deren Geometrie erstaunlich gut. Die Bilder sind jedoch insgesamt weniger detailreich als bei der Konkurrenz.

Als Technikdemo zeigt Luma Dream Machine, wie weit Videogeneratoren schon sind. Die fertigen Clips machen durch ihre dynamischen Bewegungen Spaß, taugen aber noch nicht für kommerzielle Projekte.

👉 bewegungsreiche Szenen

👎 wenig Details, geringe Konsistenz



## Runway ML

Die New Yorker Firma Runway ML ist in mehrfacher Hinsicht der Pionier der Branche. Im Jahr 2018 von Studenten der New York University gegründet, hat Runway ML bereits im Februar 2023 die ersten kommerziellen Videogeneratoren Gen-1 und Gen-2 veröffentlicht. Die Firma entwickelt aktuell die dritte Generation seiner Video-App, die immerhin im Alpha-Stadium für Tests zur Verfügung steht: Namentlich Runway Gen-3 Alpha. Gen-1 setzt hochgeladene Videos fort, Gen-2 erzeugt bewegte Bilder aus JPEG- und PNG-Dateien und Gen-3 setzt Textkommandos in Videos um.

Der Stil eines Videos lässt sich über die Seed-Nummer kopieren. Im Modell Gen-2 kann man detailliert die gewünschten Kamerabewegungen bestimmen. Ein „Prompting Guide“ klärt über Kamerastile, -fahrten und -geschwindigkeiten sowie Bildstile und -ästhetik auf. Bei 625 Credits im Standard-Abo und 10 Credits für eine Sekunde im Gen-3-Modell leert sich das Konto schnell.

Die von Runway generierten Videos verdienen am ehesten das Prädikat profitauglich. Hochgeladene Porträts verändert Runway in der Animation kaum. Wellen im Surfer-Motiv zeichnet es mit wunderbar sanften Bewegungen. Beim Restaurant-Motiv verzichtet Runway auf mampfende Münder, sondern zeigt ein Close-up des Spaghettitellers wie vom Food-Fotografen. Der Nachteil: Die Videos wirken oft wie in Zeitlupe. Den Runway-Modellen unterlaufen weniger Fehler als den übrigen, ihren Ergebnissen mangelt es aber auch an Dynamik.

➡ ruhige, oft fehlerfreie Szenen

🔴 vergleichsweise teuer

Details, aber zuweilen auch mehr Fehler: Häufig betreffen sie die menschliche Anatomie. Luma verhält sich deutlich weniger zimperlich und erzeugt lebhaftere Szenen, in denen sich allerdings auch Schriftzüge auf Gläsern sowie Zinken von Gabeln bewegen und sich die Gliedmaßen der Personen gegenseitig auffressen.

Den Diensten liegen unterschiedliche sogenannte Ethikrichtlinien zugrunde. Viele Anbieter verpflichten sich selbst, keine Inhalte zu generieren, die sich für Fake-News-Kampagnen nutzen lassen oder Urheberrechte verletzen. Videos, die Präsidenten



„A happy couple holdings hands, walking from the eiffel tower to the statue of liberty“.  
Von oben: Haiper, Luma, Runway

KI-Videogeneratoren			
Produkt	Haiper (Beta)	Luma Dream Machine	Runway ML
Hersteller, URL	Haiper, haiper.ai	Luma, lumalabs.ai	Runway ML, runwayml.com
Systemanforderungen	Webbrowser, iOS ab 15.0	Webbrowser	Webbrowser
Sprache	Englisch	Englisch	Englisch
Funktionen			
Video zu Video	✓	–	✓ (Gen-1)
Bild zu Video	✓	✓	✓ (Gen-2)
Text zu Video	✓	✓	✓ (Gen-3 Alpha)
Auflösung	1280 × 768 Pixel (720p), 1080p-Upscaling (Pro-Abo)	1360 × 752 Pixel	1280 × 768 Pixel (720p), 4K-Upscaling
Seitenverhältnis	16:9, 4:3, 1:1, 3:4, 16:9	16:9, 9:16	21:9, 16:9, 4:3, 1:1, 3:4, 16:9
maximale Videolänge	8 Sekunden (erweiterbar auf 12 Sekunden)	10 Sekunden	15 Sekunden (Gen-1), 16 Sekunden (Gen-2), 10 Sekunden (Gen-3 Alpha)
Inhalte mit Prominenten und Marken	✓ (sofern die KI darauf trainiert ist)	– (Hinweis auf Inhaltsrichtline)	– (Hinweis auf Inhaltsrichtline)
Abo und Anmeldung			
Anmelden und registrieren	Google-Konto	Google-Konto	E-Mail, Apple, Google
kostenloses Angebot	300 Credits, max. 10 pro Tag	30 Generationen pro Monat, max. 5 pro Tag	125 Credits, nur für Modelle Gen-1 und Gen-2
Wasserzeichen	✓ (kostenloses Abo)	✓ (kostenloses Abo)	✓ (kostenloses Abo)
kommerzielle Nutzung	✓ (ab Pro-Abo)	✓ (ab Standard-Abo)	✓
Cloudspeicher	✓ (kein Limit)	✓ (kein Limit)	5 GByte (kostenlos), 100 GByte (Standard), 500 GByte (Pro)
privater Cloudspeicher	✓ (ab Pro-Abo)	✓	✓
Zugriff über API	✓	✓	–
Bewertung			
Bedienung	⊕⊕	⊕⊕	⊕
Funktionsumfang	○	⊖	⊕
Video: Variation / Kontinuität	⊕ / ⊖	⊕ / ⊖	⊖ / ⊕⊕
Realismus / Detailtiefe	⊕ / ⊕	⊖ / ⊖	⊕ / ⊕
Preis	35,70 US-\$ (5000 Credits)	29,99 US-\$ (120 Gen. pro Monat)	15 US-\$ (625 Credits)
⊕⊕ sehr gut   ⊕ gut   ○ zufriedenstellend   ⊖ schlecht   ⊖⊖ sehr schlecht   ✓ vorhanden   – nicht vorhanden   k. A. keine Angabe			

### Literatur

[1] André Kramer, Pixelautomaten, KI-Bildgeneratoren im Test, c't 11/2024, S. 116

Videobeispiele im Web  
[ct.de/wggw](https://ct.de/wggw)

ten der Vereinigten Staaten zeigen, erstellte nur Haiper. Die übrigen Dienste wiesen auf ihre Richtlinien hin. Auch den beliebten Droiden R2-D2 zauberte nur Haiper auf die Leinwand. Luma und Runway lieferten generische, dem Original aber völlig unähnliche Darstellungen. Den Eiffelturm hingegen bauten sie alle nach, wobei nur Haiper unsere Vorgabe annähernd korrekt umsetzte, sowohl ihn als auch die Freiheitsstatue abzubilden.

### Fazit

Der Anfang ist gemacht, doch es ist noch ein weiter Weg, bis KI brauchbare Videos generiert. Viele Dienste scheitern vor allem daran, Personen im Laufe des Videos konsistent darzustellen. Bei Luma Dream Machine morphen Gliedmaßen munter durcheinander. Das sichert dem Dienst breite mediale Auf-

merksamkeit auf Reddit und Co. Als Einspieler für kommerziell genutzte Videoprojekte eignen sich die meisten Clips jedoch nicht.

Haiper und Runway machen ihre Sache zumindest bei Standardaufnahmen wie händchenhaltenen Hochzeitspaaren oder Snowboardern besser. Bei unnatürlichen und aufwendigen Szenen müssen auch sie passen. Der Dienst Haiper steht zwar bisher nur als Betaversion zur Verfügung, besticht aber durch hohe Detailtiefe, lebhaftere Farben und starken Kontrast. Die dargestellten Personen überzeugen allenfalls, wenn man sie von hinten sieht. Landschaften, Gegenstände und Explosionen hingegen gelingen Haiper schon erstaunlich gut. Die besten Videos liefert bisher Runway. In kluger Selbstbeschränkung verzichtet der Dienst auf hektische Action und liefert dabei ruhige, aber authentische und glaubwürdige Szenen ab. (akr) **ct**



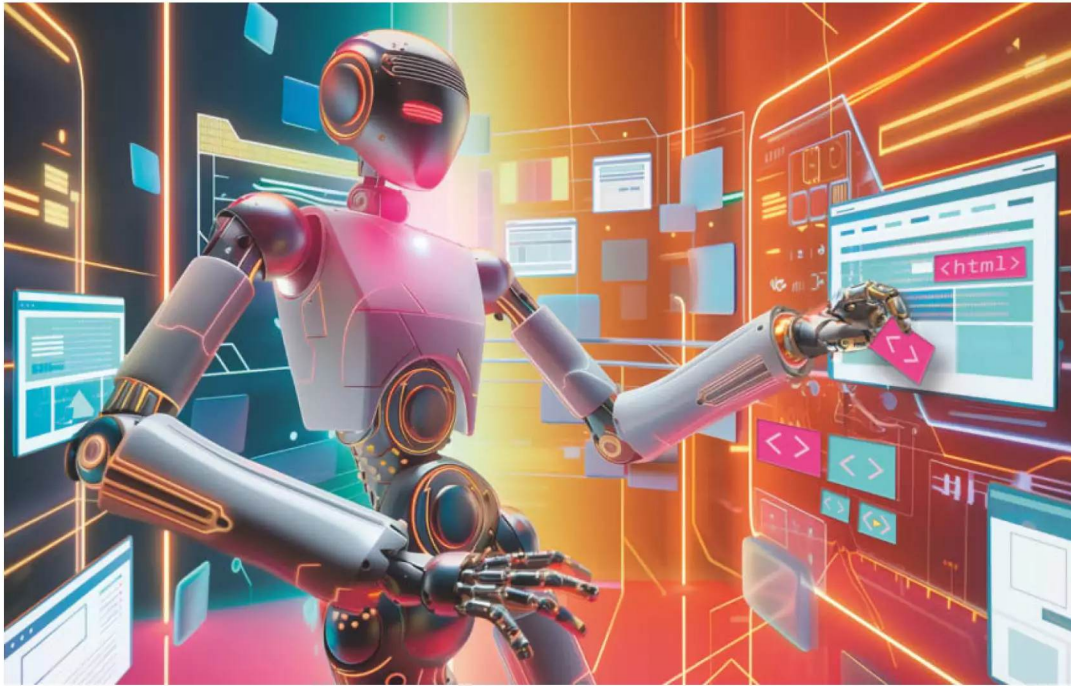


Bild: KI Midjourney | Collage c't

# Browser-Helferlein per KI generieren

Mithilfe von KI können Sie fremde Webseiten beliebig umbauen und ganz an Ihre Bedürfnisse anpassen: Fügen Sie neue Funktionen hinzu, entfernen Sie nervige Elemente, hübschen Sie das Design auf und vieles mehr. Sie formulieren einfach, was Sie möchten, kurz darauf spuckt die KI ein passendes Userscript aus, das lokal in Ihrem Browser ausgeführt wird.

Von **Ronald Eikenberg**

**S**tellen Sie sich vor, Sie könnten Ihre Lieblingswebseiten ganz nach Ihrem Geschmack verändern und zum Beispiel coole Funktionen hinzufügen oder das Design individualisieren. Sie hätten Einblicke hinter die Kulissen und könnten sogar nervige Bugs beseitigen, die Sie seit Jahren

stören. Durch den geschickten Einsatz von künstlicher Intelligenz ist jeder dazu in der Lage – Sie müssen ihr nur Ihre Wünsche mitteilen.

Das technische Fundament bilden die sogenannten Userscripts. Das ist JavaScript-Code, der die gewünschten Modifikationen an der Webseite durch-

führt. Browser-Erweiterungen wie Violentmonkey, Greasemonkey oder Tampermonkey (siehe ct.de/wwnn) injizieren die Skripte beim Surfen automatisch in die Webseite, der Browser setzt die Änderungswünsche anschließend lokal auf Ihrem Rechner um. Empfehlenswert ist vor allem das Open-Source-Projekt Violentmonkey, das Firefox, Chrome und verwandte Browser unterstützt.

Userscripts gibt es seit Jahren. In Katalogen wie Greasy Fork (siehe ct.de/wwnn) findet man etliche vorgefertigte Skripte für viele Webseiten, häufig sind die jedoch veraltet und funktionieren nicht so wie erhofft. Zudem setzt die Nutzung fremder Skripte ein gewisses Vertrauen voraus. Alternativ können Sie selbst welche entwickeln, zum Beispiel über in Violentmonkey integrierten Code-Editor. Das ist keine Raketenwissenschaft, erfordert aber Zeit und zumindest Grundkenntnisse in Webentwicklung. Die Hürde, „mal eben schnell“ etwas auszuprobieren, ist relativ hoch und der Schuh muss schon ordentlich drücken, damit sich der Aufwand lohnt.

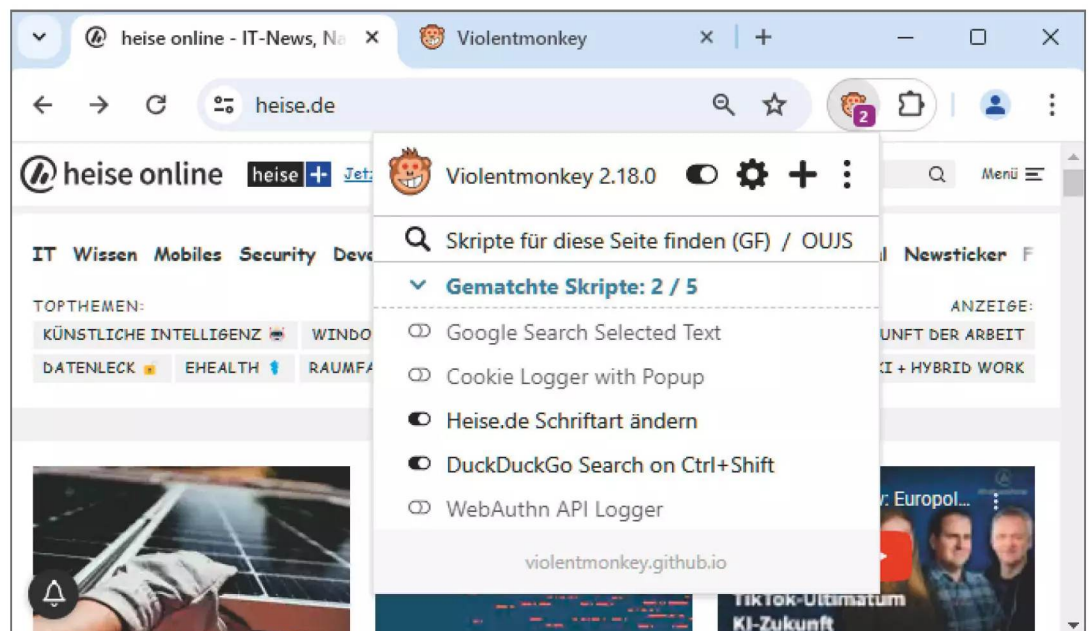
An dieser Stelle eilt die KI zur Hilfe: Moderne Sprachmodelle wie ChatGPT beherrschen die Entwicklung von Userscripts für Violentmonkey & Co. aus dem Effeff. Bei unseren Versuchen entstand oft schon im ersten Anlauf ein funktionstüchtiges

Skript, das eine Webseite wie gewünscht modifiziert hat. Um solche Skripte zu generieren, reicht ein kostenloser OpenAI-Account mit ChatGPT 3.5 aus. Sie können es auch mit anderen Sprachmodellen probieren, JavaScript ist für die meisten keine Herausforderung.

Ein einfaches Beispiel zum Warmwerden: Der Prompt „Erstelle ein Violentmonkey-Skript, das die Schriftart von heise.de in Comic Sans ändert“ liefert bei ChatGPT das abgebildete Userscript. Es beginnt mit den nötigsten Metainformationen (ausführliche Doku siehe ct.de/wwnn), die Erweiterungen wie Violentmonkey & Co. benötigen.

Darin befindet sich der Name des Skripts (@name), der später in der Browser-Erweiterung angezeigt wird. @match gibt vor, auf welchen Webseiten das Skript geladen werden soll. Im gezeigten Beispiel hat ChatGPT korrekt vorgegeben, dass es ausschließlich auf heise.de aktiv ist. Über @grant kann man dem Skript erweiterte Zugriffsrechte einräumen, was in diesem Fall aber nicht nötig ist.

Die eigentliche Magie steckt in der JavaScript-Funktion darunter, mit der das Skript die Schriftart der gesamten Website verstellt. Um mehr über den gewählten Lösungsweg zu erfahren, fragen Sie einfach nach: „Erkläre das Skript“.



**Violentmonkey ergänzt Webseiten um Skripte, die Aussehen und Funktion modifizieren.**



In diesem einfachen Beispiel programmiert ChatGPT ein Userscript, das die Schriftart auf der Website heise.de ändert. Die KI kann auch wesentlich komplexere Skript erstellen.

Daraufhin liefert die KI eine ausführliche Erklärung ab: „Ein neues style-Element wird erstellt, das in den Kopfbereich der Webseite eingefügt wird. Der CSS-Stil innerhalb des style-Elements ändert die Schriftart des body-Elements zu ‚Comic Sans MS‘. !important wird verwendet, um sicherzustellen, dass diese Stilanweisung Vorrang vor anderen Stilen hat. Schließlich wird das style-Element dem head-Element der Webseite hinzugefügt.“

Da immer Zufall im Spiel ist, kann sich der konkrete Lösungsweg bei Ihnen unterscheiden, bekanntlich führen viele Wege nach Rom. Wenn Sie mit dem Ergebnis nicht zufrieden sind, lohnt es sich daher, denselben Prompt zu wiederholen.

## Skripte hinzufügen

Um solche Skripte auszuprobieren, gehen Sie wie folgt vor: Installieren Sie einen Userscript-Manager wie Violentmonkey als Erweiterung im Browser und klicken auf das Symbol der Erweiterung (bei Violentmonkey ist es ein grinsender Affenkopf). Danach klicken Sie auf das Plus, um ein neues Skript im Code-Editor zu öffnen. Dort hinein kopieren Sie schließlich den KI-generierten Code. Das vorgegebene Grundgerüst überschreiben Sie dabei. Zuletzt klicken Sie oben rechts auf „Speichern“.

Wenn Sie jetzt zurück zu dem Tab mit der Webseite wechseln und ihn aktualisieren, ist das Skript



aktiv. Auch nach späteren Änderungen müssen Sie einen Refresh durchführen oder die Website neu aufrufen, damit das geänderte Skript beim Laden injiziert wird. Ob und wie viele Skripte aktiv sind, erkennen Sie an der Ziffer an dem Violentmonkey-Symbol.

In dem simplen Beispiel erstrahlt heise.de von nun an in Comic Sans MS. Wenn Sie genug haben, können Sie das Skript über einen Klick auf den Erweiterungsknopf und den Skriptnamen wieder abschalten. Änderungen sind jederzeit über den Knopf </> neben dem Skriptnamen möglich, der Sie wieder in den Editor befördert.

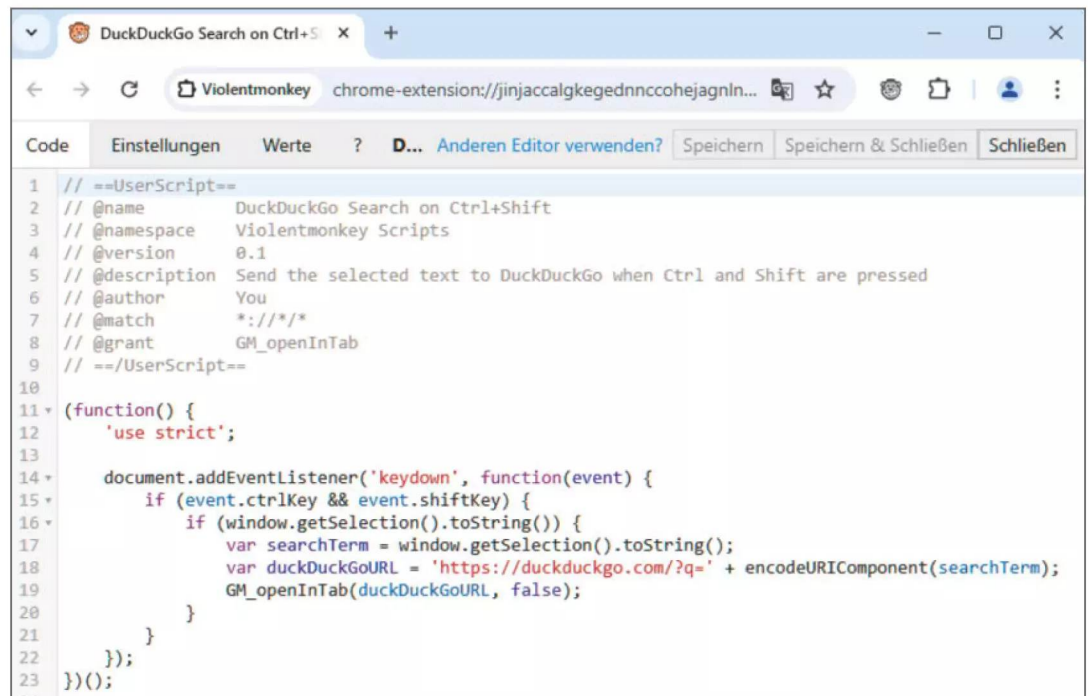
Zugegeben, das Comic-Sans-Beispiel ist weder besonders komplex noch nützlich – aber Sie haben jetzt eine gute Vorstellung davon, wie das Ganze funktioniert. Jetzt ist an der Zeit, die KI mit konkreten und schwierigeren Aufgaben zu betrauen.

Wenn Sie zum Beispiel die YouTube-Oberfläche aufräumen möchten, können Sie die KI bitten, Elemente zu entfernen, die Sie nicht interessieren, etwa die Kommentare anderer Nutzer. In diesem Fall genügt meist ein einzelner, präziser Prompt, um zu

einem brauchbaren Ergebnis zu kommen: „Erstelle ein Violentmonkey-Skript, das bei YouTube die Kommentare ausblendet“. Herauskommen kann ein Skript, das mit der Funktion `document.getElementById()` nach dem Web-Element mit der ID „comments“ sucht und dieses anschließend mit `.style.display = 'none'` unsichtbar macht.

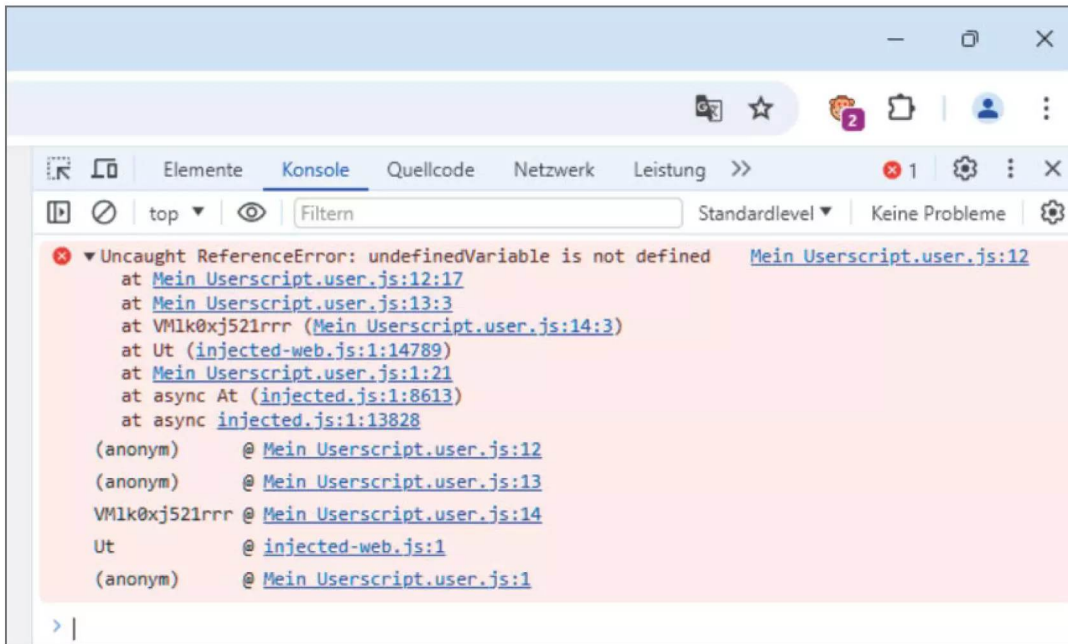
## Fehlersuche

Wenn das KI-generierte Skript nicht funktioniert, probieren Sie es einfach erneut. Oder Sie teilen der KI mit, dass etwas klemmt, um Lösungsvorschläge und Korrekturen zu erhalten. Überprüfen Sie zudem, ob oben in den Metainfos des Skripts bei `@match` die korrekte URL der Webseite steht, etwa `https://www.youtube.com/*` für alle HTTPS-URLs, die mit `www.youtube.com` beginnen. Ohne das Wildcard-Sternchen muss die URL exakt stimmen. Wenn Sie die Zeile weglassen, injiziert Violentmonkey Ihr Skript in alle Webseiten. Das sollten Sie sich jedoch gut überlegen, denn Skripte, die speziell an die Struktur



```
1 // ==UserScript==
2 // @name      DuckDuckGo Search on Ctrl+Shift
3 // @namespace  Violentmonkey Scripts
4 // @version    0.1
5 // @description Send the selected text to DuckDuckGo when Ctrl and Shift are pressed
6 // @author     You
7 // @match      *://**/*
8 // @grant      GM_openInTab
9 // ==/UserScript==
10
11 (function() {
12     'use strict';
13
14     document.addEventListener('keydown', function(event) {
15         if (event.ctrlKey && event.shiftKey) {
16             if (window.getSelection().toString()) {
17                 var searchTerm = window.getSelection().toString();
18                 var duckDuckGoURL = 'https://duckduckgo.com/?q=' + encodeURIComponent(searchTerm);
19                 GM_openInTab(duckDuckGoURL, false);
20             }
21         }
22     });
23 })();
```

Über den Code-Editor fügen Sie Userscripts hinzu und bearbeiten sie.



**Funktioniert ein Skript nicht wie erhofft, füttern Sie die KI mit der Fehlermeldung, die Sie auf der JavaScript-Konsole Ihres Browsers finden.**

einer bestimmten Webseite angepasst wurden, können andernorts Fehler verursachen.

Am besten funktioniert die Fehlersuche, wenn Sie der KI eine konkrete JavaScript-Fehlermeldung liefern. Öffnen Sie über das Browsermenü die Entwicklertools und dort die Konsole, um etwaige Fehler bei der Ausführung des Codes zu inspizieren. Alternativ können Sie unter Windows und Linux auch Strg+Umschalt+I und unter macOS Opt+Cmd+I drücken.

Halten Sie dort nach Fehlern Ausschau, die von Ihrem Userscript ausgehen. Sie erkennen es an dem in den Metainfos festgelegten Skriptnamen (name), gefolgt von .user.js und der Zeilennummer, in der der Fehler aufgetreten ist. Kopieren Sie den Fehler in den KI-Chat, damit der JavaScript-Code korrigiert wird.

Wenn Sie Skripte für weniger prominente Webseiten oder Dienste aus dem Intranet erstellen, müssen Sie davon ausgehen, dass KI nichts über deren Aufbau weiß und ins Blaue rät, wenn bestimmte Elemente oder Funktionen manipuliert werden sollen. Dann sollten Sie einen Kontext liefern, zum Beispiel einen Ausschnitt des Webseiten-

Quellcodes oder Identifizierungsmerkmale von Elementen.

Hierfür sind wieder die Entwicklertools des Browsers hilfreich. Klicken Sie in den Entwicklertools oben links auf den Knopf mit dem Mauszeiger und anschließend auf ein Element der Webseite, um direkt zum dazugehörigen HTML-Code zu gelangen und kopieren Sie ihn in den KI-Chat mit einer Fehlerbeschreibung.

## Backstage-Pass

Mit Userscripts können Sie interne Abläufe der Webseiten beobachten, die Ihnen ansonsten verborgen bleiben. Probieren Sie doch mal den folgenden Prompt: „Erstelle ein Violentmonkey-Skript, das das Setzen von Cookies anzeigt“. Nachdem Sie das Skript importiert und die Webseite aktualisiert haben, finden Sie auf der Konsole Ausgaben wie „Cookie gesetzt: cookiebanner=closed; path=/; expires=Wed, 22 May 2024 13:27:49 GMT“.

Wenn Sie eine Anforderung wie „Das Skript soll die Cookies als Hinweisfenster auf der Webseite melden, das sich nach 5 Sekunden schließt“ nach-

reichen, sollte das Skript das Setzen von Cookies direkt auf der Webseite anzeigen und Sie können sich das Öffnen der Konsole sparen. In unserem Versuch wurde ein kleines `div`-Element als Overlay über die Webseite gelegt und ein Timer sorgte dafür, dass es nach kurzer Zeit wieder verschwand. Wie immer gilt: Der Lösungsweg kann in Ihrem Fall ein anderer sein, möglicherweise sind auch mehrere Versuche nötig, ehe ein Skript funktioniert.

Um sich bemerkbar zu machen, könnte Ihr Userscript auch das Notification-API des Browsers verwenden und Benachrichtigungen erzeugen, die als Mitteilung Ihres Betriebssystems erscheinen. Die bekommen Sie auch dann mit, wenn die Webseite im Hintergrund läuft oder Sie das Browserfenster minimiert haben. Das ist zum Beispiel nützlich, um informiert zu werden, wenn auf der Seite ein zeitintensiver Vorgang abgeschlossen ist oder eine neue Nachricht eintrifft. Auf diese Weise können Sie auch andere Web-APIs einspannen.

## Workflow optimieren

Zudem kann der KI-Helfer Ihrem Browser neue Funktionen ganz nach Ihren Wünschen spendieren, für die normalerweise eine eigene Erweiterung nötig wäre. Um schnell und einfach im Netz zu recherchieren, könnten Sie die künstliche Intelligenz bitten, ein Skript zu schreiben, das den gerade markierten Text an eine Suchmaschine schickt, wenn Sie eine bestimmte Tastenkombination drücken.

Auch das erfordert nicht viel Eigenleistung, Sie sagen ChatGPT einfach, was Sie möchten: „Erstelle ein Violentmonkey-Skript, das den markierten Begriff an DuckDuckGo sendet, wenn ich Strg und Umschalt drücke“. Spätestens nach zwei, drei Versuchen sollten Sie ein funktionsfähiges Skript am Start haben, das eine Suche anwirft, wenn Sie die Tasten Strg+Umschalt gleichzeitig drücken, während Text markiert ist.

Es dürfte darauf hinauslaufen, dass das Skript mit `document.addEventListener` nach Tastatur-Ereignissen lauscht (`ctrlKey` und `shiftKey`), den markierten Text an den URL-Parameter „q“ der DuckDuckGo-URL (<https://duckduckgo.com/?q=>) hängt und diese mit `window.open` in einem neuen Tab öffnet. Bonusaufgabe: Fügen Sie doch mal eine weitere Tastenkombination hinzu, die den markierten Begriff direkt an die ChatGPT-URL <https://chat.openai.com/chat?q=> hängt, damit Sie den KI-Helfer jederzeit im Zugriff haben. Mit dem Wissen über diese Möglichkeiten finden Sie sicher bald eigene Ideen.

Zuletzt noch ein Beispiel aus dem Redaktionsalltag: Für c't 9/2024 haben wir uns mit dem passwortfreien Login-Verfahren Passkeys beschäftigt, das von immer mehr Webdiensten angeboten wird. Zur Implementierung wird das WebAuthn-API verwendet. Webdienste können darüber einige Sicherheitsanforderungen stellen, etwa, wo ein Passkeys gespeichert werden darf und welche Kryptografie zum Einsatz kommt.

Diese Details bekommt man aber nicht zu Gesicht. Um herausfinden zu können, was hinter den Kulissen passiert, baten wir die KI mit folgendem Prompt, uns ein Analyseskript zu erstellen: „Erstelle ein Violentmonkey-Skript, das beim Erstellen eines Passkeys Details über den Einsatz des WebAuthn-API auf der Konsole anzeigt“. Dabei kam ein Skript heraus, dass die an das API übergebenen Daten zur Analyse sichtbar macht.

## Sync und mobile Nutzung

Wenn mit Ihren Skripten zufrieden sind, können Sie Ihre Kreationen über die Violentmonkey-Einstellungen als Backup in eine Zip-Datei sichern. Für die Nutzung mit mehreren Rechnern ist die Sync-Funktion der Browser-Erweiterung besser geeignet, die Dropbox, Google Drive, OneDrive und WebDAV-Server wie Nextcloud einspannen kann.

Auch auf Smartphones und Tablets können Sie Ihre Skripte ausprobieren: Firefox für Android kann die Violentmonkey-Erweiterung einfach aus dem Add-on-Katalog installieren, für iOS und iPadOS gibt's die Safari-Erweiterung Userscripts (siehe [ct.de/wnwn](https://www.ct.de/wnwn)). Die ist auch mit der macOS-Version von Safari nutzbar, für die es kein Violentmonkey gibt.

## Fazit

Mach Dir die Welt, wie sie Dir gefällt: Mit Userscripts passen Sie Webseiten ganz nach Ihren Vorstellungen an. Die technischen Grundlagen gibt es seit vielen Jahren, doch mit Unterstützung von ChatGPT & Co. kommen Sie schneller denn je und auch ohne Programmiererfahrung zum Ziel. Nebenbei ist es interessant zu beobachten, welche Lösungswege die KI für eine bestimmte Aufgabenstellung findet.

Der Kreativität sind dabei kaum Grenzen gesetzt: Rüsten Sie individuelle Funktionen nach, bauen Sie das Design um oder entfernen Sie lästige Elemente, die Sie schon länger stören. Teilen Sie uns gern mit, welche individuellen Probleme Sie mithilfe von Userscripts und KI lösen konnten. (rei) **ct**





Bild: KI Midjourney | Bearbeitung: c't

# ChatGPT als Sprechtrainer

Sie möchten gerne mit einem KI-Sprechtrainer eine Fremdsprache üben, dafür aber nicht gleich ein Abo abschließen? Kein Problem: ChatGPT lässt sich auch in der kostenlosen Version am Rechner und am Smartphone als solcher nutzen.

Von **Nico Jurrán**

**V**om Start weg begeisterte ChatGPT als Chatbot mit der Fähigkeit, „natürliche“ Gespräche führen zu können. Da die generative KI zudem etliche Sprachen beherrscht und von sich aus geeignete Nachfragen stellt, liegt es nahe, sie auch als Konversationstrainerin für Fremdsprachen

einzusetzen. Für Textchats funktioniert das am Rechner ohne weiteres Zutun, für Voice Chats muss man ChatGPT aber erst dazu bringen, zuzuhören (Speech To Text) und zudem selbst zu sprechen (Text To Speech). Wie das genau funktioniert, erfahren Sie in diesem Artikel.

## Lass uns reden!

Um loszulegen, benötigen Sie einen Rechner mit Lautsprecher und Mikrofon oder einem Headset, einen Chrome Browser und Zugang zu ChatGPT. Zwar liefert die kostenpflichtige GPT-Version schnellere und bessere Antworten und ist hinsichtlich der Anfragen weniger limitiert, für ein einfaches Sprechtraining reicht aber das kostenlose ChatGPT seit Version 3.5 aber völlig aus. Um die KI ohne zu viele Einschränkungen nutzen zu können, sollten Sie sich aber bei OpenAI unter <https://chat.openai.com/> mit Ihrer E-Mail-Adresse registrieren.

Als Browser kommt Google Chrome deshalb in Frage, weil dessen Web Store kostenlose Erweiterungen anbietet, die ChatGPT um eine Sprachverarbeitung ergänzt, also um eine Eingabe und Ausgabe. Bei uns brachte „Sprachsteuerung für ChatGPT“ von VoiceWave die besten Ergebnisse (Download unter [ct.de/w9kq](https://ct.de/w9kq)). Ist die Sprachsteuerung aktiviert, erscheint im Eingabefeld von ChatGPT ein Mikrofonsymbol. Wichtig: In einigen Fällen taucht das Mikrofonsymbol erst auf, wenn die Erweiterung im Entwicklermodus läuft. Diesen aktiviert man über das Feld oben rechts neben der Adressleiste (das Puzzle-Symbol) unter „Erweiterungen verwalten“.

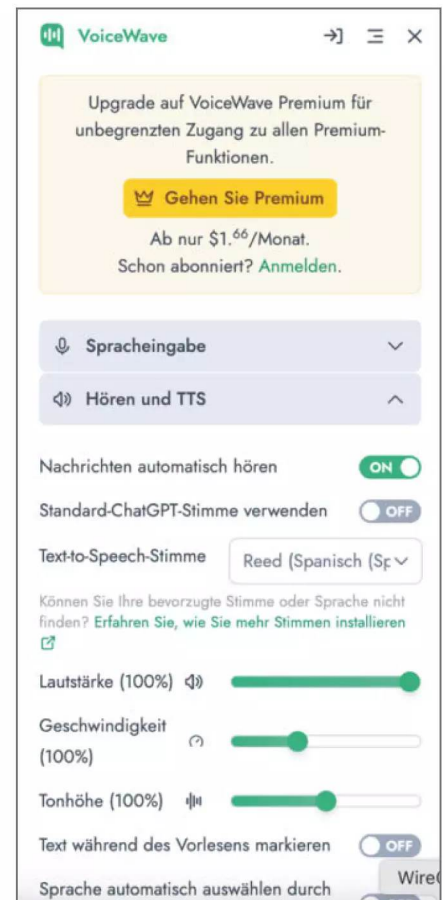
Das Schöne an der Add-On-Lösung ist, dass man den KI-Sprachlehrer nutzen kann, ohne dafür jedes Mal die Systemsprache umzustellen oder sonstige Veränderungen Einfluss auf das gesamte System nehmen. Stattdessen stellt man einfach in der Erweiterung die gewünschte Sprache ein, in der man mit ChatGPT sprechen will.

## Einigung auf eine Sprache

Die Sprache, in der man mit der KI sprechen möchte, legt man unter „Spracheingabe/Sprachsprache“ in den Einstellungen fest, die man über das Rädchen-Icon erreicht, das sich über dem gerade erwähnten Mikrofonsymbol befindet. Dies sollte man auf jeden Fall vor der ersten Konversation einstellen: Ist die falsche Sprache vermerkt, versucht ChatGPT krampfhaft, dem Kauderwelsch einen Sinn zu geben. Aus dem spanischen „hablas“ (Du sprichst) wird dann etwa bei Deutsch als Eingabesprache „Ablass“.

ChatGPT beantwortet Nachfragen auch auf Deutsch, allerdings stellt man diese dann besser über den Textchat, um nicht wieder die Eingabesprache wechseln zu müssen, damit die KI sie auch versteht.

Ebenso ist es wichtig, in den Settings unter dem Punkt „Hören und TTS/Text-To-Speech-Stimme“ die



**Bei der Google-Chrome-Erweiterung muss man für die Ein- und Ausgabe die Lernsprache einstellen, da ChatGPT sonst keine verwertbaren Daten erhält beziehungsweise ungefähr so klingt, als wenn im Auto ein auf Deutsch eingestelltes Navi fremdsprachige Straßennamen ausspricht.**

Stimme in der Sprache auszuwählen, in der die KI sprechen soll. Dies ist üblicherweise dieselbe wie die für die Eingabe. Vergisst man diesen Punkt, hört sich ChatGPT in etwa so an wie Fremdsprachler, die ohne Vorkenntnisse krampfhaft versuchen, Deutsch zu sprechen.

Die Spracherweiterung nutzt die Systemstimmen des Betriebssystems, wobei teilweise für eine Spra-

che bereits verschiedene Sprecher und sogar unterschiedliche Dialekte bereitstehen. Finden Sie keine brauchbare TTS-Stimme, können Sie weitere kostenlos über die Systemeinstellungen installieren. Erklärungen, wie Sie das konkret für die unterschiedlichen Versionen von macOS beziehungsweise MacOS bewerkstelligen und wie sich weitere kostenlose Stimmen in Microsofts Betriebssystem einpflegen lassen, finden Sie unter [ct.de/w9kq](https://ct.de/w9kq).

Redet die KI für Ihren Geschmack zu schnell oder zu langsam, können Sie schließlich in der Erweiterung dem Feld „Geschwindigkeit“ die Ausgabegeschwindigkeit von 50 bis 300 Prozent anpassen.

## Durcheinander vermeiden

Wenn Sie in Ihren Sätzen längere Pausen machen, was zu Beginn ganz normal ist, stellen Sie unter „Spracheingabe“ ein, dass die Spracheingabe nicht automatisch gesendet wird. Andernfalls passiert es, dass Ihre Eingaben bereits an ChatGPT gehen, obwohl Sie noch gar nicht zu Ende gesprochen haben. Schalten Sie die Automatik ein, empfiehlt es sich, unter dem Punkt unter „Erweiterte Einstellungen“ den Haken „Mikrofon ausschalten, wenn Antwort gehört wird“ zu aktivieren. Andernfalls kommt es vor, dass die KI mit sich selbst redet. Generell ist es die beste Lösung, ein Headset zu benutzen, damit solche Probleme gar nicht erst entstehen.

Damit sind die Vorbereitungen abgeschlossen, sodass man mit einer Eingabe (Prompt) à la

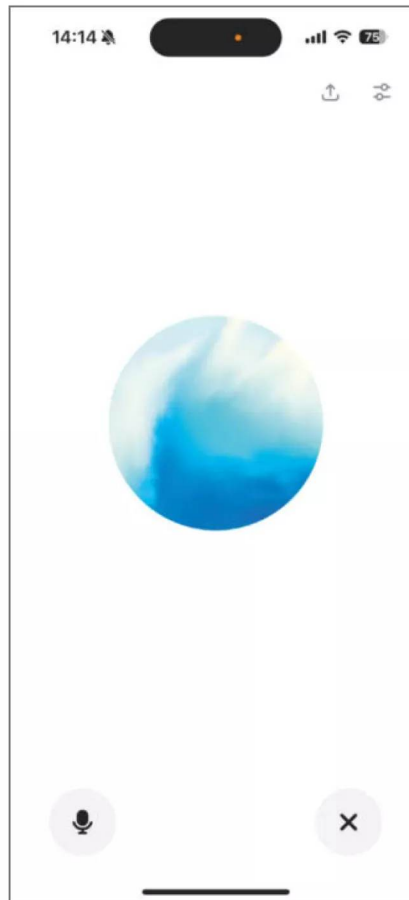
Lass uns ein Gespräch auf Französisch führen

die Konversation starten kann. ChatGPT quatscht darauf munter los, stellt Fragen und geht auf Ihre Antworten ein.

## Korrektur erbeten

Sie werden allerdings schnell merken, dass die KI großzügig ist, was Fehler bei der Eingabe betrifft. ChatGPT geht also einfach darüber hinweg, wenn Sie etwa bei einem Artikel das falsche Geschlecht oder Singular statt Plural benutzen. Die KI macht also genau das, was man als Sprachschüler eigentlich gerade nicht möchte, weil sich ohne Korrekturen und Verbesserungen Fehler manifestieren.

Deshalb ist es sinnvoll, wenn Sie ChatGPT mit auf den Weg geben, dass er auf Fehler achten und diese korrigieren soll. Allerdings sollten Sie ihn auffordern, das Gespräch danach fortzusetzen, damit die Kon-



**In der ChatGPT-App kommt man über das Kopfhörersymbol in den Voice-Chat-Modus. Hier kann man frei drauflosreden, eine Transkription der KI-Ausgaben erhält man an dieser Stelle aber nicht.**

versation nicht bei jedem Fehler stockt. Ein möglicher Prompt hierfür wäre:

Korrigiere mich, wenn ich etwas falsch sage. Führe danach die Konversation fort.

Nach unserer Erfahrung kommt es vor, dass ChatGPT die Anweisung, die angesprochenen Sätze zu korrigieren, nicht gleich umsetzt. Um sicherzustellen, dass die KI die Aufgabe verstanden hat, empfiehlt es sich daher, testweise einen fehlerhaften Satz zu sprechen und gegebenenfalls ChatGPT dann noch



einmal auf seine Pflicht zur Korrektur hinzuweisen.

Um sich auf bestimmte Alltagssituationen vorzubereiten, können Sie ChatGPT dazu auffordern, mit Ihnen ein Rollenspiel in der Zielsprache zu führen. Ein geeigneter Prompt könnte wie folgt lauten:

Lass uns ein Rollenspiel auf Französisch führen, bei dem ich als Gast ein Restaurant besuche und du der Kellner bist.

## Auf dem Smartphone

Shon seit einiger Zeit gibt es in der ChatGPT-App eine Sprachverarbeitung – zunächst nur für die kostenpflichtige Version, später schaltete OpenAI diese Funktion generell frei.

In der Smartphone-App starten Sie den Voice Chat, indem Sie auf das Kopfhörersymbol rechts neben dem Texteingabefeld drücken. Die komplette Steuerung funktioniert in diesem Modus mündlich. So

können Sie ChatGPT etwa auffordern, ein Rollenspiel in einer bestimmten Fremdsprache zu starten. Zwischen Sprachen wechselt die KI nahtlos und beantwortet daher Nachfragen zu ihren Ausführungen auch in anderen Sprachen (wie Deutsch). Der Bitte um Korrekturen kommt ChatGPT auf dem Handy ebenfalls nach.

Haben Sie gerade erst mit einer Fremdsprache begonnen, kann dieser Voice-Chat-Modus aber überfordernd sein. So bekommen Sie hier keine direkte Transkription, sondern sehen sie erst nach der Rückkehr auf den Startbildschirm, weshalb Sie bereits über die nötige Hörkompetenz verfügen sollten. Problematisch wird die automatische Erkennung der Eingaben, wenn man längere Pausen beim Sprechen macht. Dieses Problem ließ sich früher lösen, indem man „Hold for manual control“ beim Sprechen gedrückt hielt. Dann verarbeitete die App die Eingabe erst, wenn Sie den Finger wieder vom Bildschirm nahmen. Diese Funktion hat OpenAI leider entfernt. (nij) **ct**

Download Sprachsteuerungserweiterung

[ct.de/w9kq](https://ct.de/w9kq)

# NEU: c't Desinfec't 2024/25

## Das Notfall-System für den Ernstfall

Komplett auf  
32 GByte USB-Stick.  
Desinfec't startet  
direkt vom Stick.

### Das kann c't Desinfec't:

- ▶ PC-Schädlinge jagen: perfekt geschützt vor Viren & Malware
- ▶ Daten retten: Bedrohungen erkennen, bevor Schaden entsteht
- ▶ Gelöschte Daten wiederherstellen: schnell & sauber

**GRATIS:**  
Signatur-Updates  
bis Oktober 2025



[shop.heise.de/desinfec24](https://shop.heise.de/desinfec24)



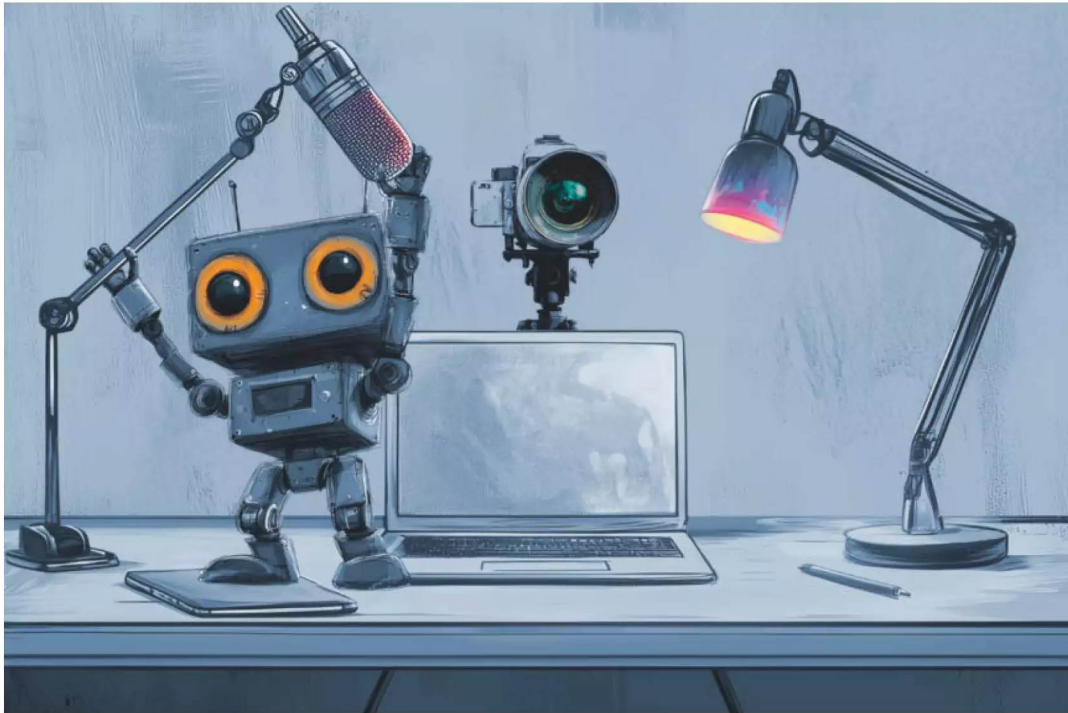


Bild: KI, Collage c't

# Podcast-Helferlein: Acht Tools mit KI

Podcasten macht Spaß, an Kapitelmarken basteln nicht unbedingt. Manche Produktionsschritte sind nicht kreativ, sondern repetitiv, fummelig und anstrengend. Wir stellen acht KI-gestützte Tools vor, die Podcastern dabei helfen können.

Von **Greta Friedrich**

**S**chnitt, Transkript, Kapitelmarken: Bevor die neue Podcast-Folge endlich erscheint, sind einige Arbeitsschritte nötig, die viele als lästig empfinden. Zum Beispiel, weil sie sich beim Transkribieren ein Interview nochmals in voller Länge, langsamer und oft auch mehrfach anhören müssen. Oder weil es ihnen schwerfällt, prägnante Punkte

für Titel und Beschreibung herauszufischen. Oder die vielen Ähms des Gesprächspartners von der Tonspur zu tilgen.

Software kann Podcaster bei diesen Arbeitsschritten unterstützen. Mithilfe künstlicher Intelligenz erzielen die Programme oft sogar ziemlich gute Ergebnisse. Einige haben sich auf bestimmte Stellen

im Podcast-Produktionsprozess spezialisiert, andere Tools decken mehrere Schritte ab oder wollen gleich die gesamte Produktion stemmen.

Die Tool-Auswahl ist riesig. Mehr oder weniger ausgelegene Werkzeuge für diese und jene Aufgabe sprießen nur so aus dem Boden. Wir stellen in dieser Übersicht eine Auswahl von acht Programmen vor, die die Bandbreite des Angebots abbilden: Adobe Podcasts, Auphonic, Mindverse, Podcastle, Riverside, StreamYard, Whisper und Zencastr. StreamYard, Auphonic und Whisper nutzen die c't-Redaktion und die Kollegen von heise online selbst für die haus-eigenen Podcasts.

## Misstrauisch bleiben

Die auf den ersten Blick oft beeindruckenden Arbeitsergebnisse der KI-Tools täuschen leicht darüber hinweg, dass diese keinen blassen Schimmer haben, was sie da bearbeiten. Deshalb erkennen die Programme manche inhaltlichen Zusammenhänge nicht, verwechseln Sprecher oder greifen sich für eine Zusammenfassung die drei langweiligsten Aspekte heraus.

Damit der Podcast und die Shownotes am Ende nicht vor Fehlern strotzen, ist es daher wichtig, stets zu überprüfen, was die Software fabriziert. Idealerweise hat die Person, die die Tools bedient, das Gespräch auf der Tonspur selbst geführt oder war zumindest dabei – und es ist seitdem nicht viel Zeit vergangen. So findet sich der Postproduzent im Material gut zurecht und erkennt KI-Fehler schnell.

## Aufnahme läuft

Um einen Podcast aufzunehmen, braucht es zunächst ein Setting aus Mikrofon und Kamera, das gut austariert und beleuchtet ist. Hier kann KI nur bedingt helfen. Das auch bei heise-Podcasts weitverbreitete Aufnahmewerkzeug StreamYard bietet immerhin ein Teststudio, in dem man Bild, Ton, Bühne und Chat testen, Hintergrundmusik, Logos und Overlays hinzufügen und, wenn alles passt, direkt live streamen oder aufzeichnen kann. Gäste lädt der Host per Link ein. Mithilfe von künstlicher Intelligenz verbessert das browserbasierte System die Audioqualität und filtert zum Beispiel Hintergrundgeräusche heraus. Die heise-Podcaster bearbeiten den Ton trotzdem mit Auphonic nach.

Beim Mikrofon-Setup will Adobe mit seinem neuen Browsertool „Adobe Podcasts“ helfen. In der Funktion „Mikro testen“ sprechen Nutzer einen vor-

gegebenen kurzen Satz in ihr Mikrofon und bekommen einen Augenblick später eine kurze Auswertung der Aspekte Abstand zum Mikrofon, Verstärkung, Hintergrundgeräusche und Echo. Das gibt zumindest eine erste Orientierung, wenn man unsicher ist.

Auch bei der Aufnahme will Adobe dabei sein, die Funktion „Studio“ ist aber noch im Betastatus. Im Vergleich zu den anderen hier vorgestellten Tools sind die Möglichkeiten eher klein, aber für einfache Tonaufnahmen reichen sie. Videoaufnahmen beherrscht Adobe allerdings bisher nicht.

Deutlich mehr Möglichkeiten bietet die Software von Zencastr [1], die es neben der Web- auch als iOS-Version gibt: Auch hier testen Podcaster ihre Aufnahmeumgebung und können dann direkt starten und Gäste einladen. Über Buttons am oberen Rand blenden sie kurze Geräusche (Trommel) oder Musik ein (z. B. Intro/Outro). Ähnlich funktioniert auch Riverside, hier gibt es allerdings neben der Browserversion auch Apps für Mac, iOS und Android. Podcastle zeichnet Gespräche im Browser und der iOS-App auf. Dort lassen sie sich auch bearbeiten.

## Ton verbessern

Manche Tools wie StreamYard polieren den Ton der Aufnahme mithilfe von KI direkt etwas auf. Für tiefergehendende Bearbeitung eignet sich beispielsweise die österreichische Software von Auphonic. Hier lädt man eine Datei hoch und wählt dann zwischen mehreren fertigen Modi (Presets), zum Beispiel „Voice Cleaner (keep music)“, das es auch ohne Musik oder mit Atmer-Entfernung gibt, oder auch „Automatic Chapters, Summary and Transcript“.

Alternativ gestaltet man vorher eigene Presets. Die Bedienoberfläche von Auphonic ist weniger schick und intuitiv als beispielsweise die von Zencastr oder Riverside, dafür bietet das Tool viel mehr Möglichkeiten für die Audioproduktion, etwa einen Leveler, Normalisierung und Geräuschreduktion. In unserem Mini-Test mit den voreingestellten Vorgaben von Auphonic reduzierte das Tool Hintergrundrauschen erfolgreich, die Stimmen wurden viel deutlicher.

## Schnitt im Transkript

Alle hier vorgestellten Tools transkribieren auch die Gespräche. Bei heise hat sich für diesen Arbeitsschritt die Spracherkennung Whisper von OpenAI etabliert, die sich in unserem Vergleich mit anderen Transkriptionsprogrammen im April 2024 ziemlich gut schlug



[2]. Drei der anderen Anbieter werben damit, dass man bei ihnen das Aufnahmematerial sogar direkt im Transkript schneiden könne, also ohne tiefere Kenntnisse im Bereich Audio- und Videoschnitt. Wir haben das im kleinen Rahmen ausprobiert.

Zencastr zeigt links eine Vorschau des Videos und rechts das Transkript, das in Sinnabschnitte unterteilt ist, aber noch einige Fehler enthält. Darin kann man nun Abschnitte löschen und Zencastr schneidet den entsprechenden Teil aus dem Video heraus. In unserem Mini-Test holperte das mal mehr, mal weniger. Man kann das Transkript aber auch einfach bearbeiten, wenn zum Beispiel ein Name falsch verstanden wurde.

Bei Riverside funktioniert das ähnlich. Achten Sie darauf, in Ihrem „Studio“ die gewünschte Sprache einzustellen (Drei-Punkt-Menü des Studios/Settings/Language). Wir hatten das zunächst übersehen und wunderten uns über die unvollständige und englischsprachige Transkription aus einem kurzen deutschsprachigen Video. Nachdem wir die passende Sprache ausgewählt hatten, klappte es auch mit dem Transkript. Riverside zeigt Füllwörter und Pausen im Transkript als Kringel und Punkte an, Nutzer können sie direkt entfernen.

## Sprachverwirrung

Adobe Podcasts transkribiert deutschsprachige Texte grundsätzlich auf Englisch – die Transkriptionssprache lässt sich nicht umstellen. Das Ergebnis liest sich zwar poetisch, ergibt aber keinen Sinn. („Divide getting thrown into too deep, which in turn to anger and disenchantment and conservatism going from this November. Not a good thing. Took nine took a vote and thrown.“ – das Original war ein kurzer Text über die Aktion „Gewalt gegen Frauen kommt nicht in die Tüte.“)

Der automatische „Filter“ von Adobe verbesserte die Sprachqualität tatsächlich. Wir hatten allerdings auch den Eindruck, dass er bewirkte, dass das Wort „Frauen“ (das Adobe für das Wort „thrown“ hielt) mit einem seltsam gerollten „r“ erklang. Adobe Podcasts befindet sich allerdings noch in der Testphase.

Podcastle erstellt klaglos ein deutschsprachiges Transkript, in dem man auch direkt schneiden kann. Sprachverwirrung taucht hier an anderer Stelle auf: Das Tool bietet mehrere Stimmen, die man einen Text sprechen lassen kann – allerdings sind die Stimmen nur auf englischsprachige Texte vorbereitet, Deutsch klingt dementsprechend seltsam. Ein KI-Assistent hilft zahlenden Nutzern, indem er Hinter-

grundgeräusche herausfiltert, den Klang ausbalanciert und längere Pausen entfernt.

## Aufbereiten und werben

Einige Tools liefern lizenzfreie Musik für Hintergrundgedudel oder Jingles mit. Wer sich selbst kreativ austoben möchte, kann mithilfe von KI eigene Musik basteln. Dafür eignen sich zum Beispiel die Dienste Suno und Udio, die die Kollegen von heise+ bereits getestet haben [3].

Bei der Verbreitung der Folgen in sozialen Medien helfen einige der Tools den Podcastern, indem sie automatisch Beschreibungstexte oder Social-Media-Posts generieren oder fertige Podcasts direkt veröffentlichen. In Auphonic können Nutzer selbst Kapitelmarken anlegen oder importieren, für zahlende Nutzer erstellt die App automatisch Shownotes. Automatische Kapitelmarken dagegen beherrschen von den vorgestellten Tools nur Riverside und Zencastr.

Kurze Videoschnipsel für Social Media generieren Nutzer zum Beispiel mit Riverside oder auch mit Zencastr, hier sogar direkt im Transkript-Bearbeitungsmodus. Die KI pickt dafür die „viralsten“ Momente des Videos heraus; alternativ markiert der Social-Media-Verantwortliche eine mindestens 30-sekündige Textpassage im Transkript, die er als Clip extrahieren möchte. Veröffentlichen kann er seinen Podcast direkt über die Plattform. Zencastr erstellt dann einen RSS-Feed und eine Podcast-Website. Ebenso funktioniert die Veröffentlichung bei Riverside und über das sogenannte „Hosting Hub“ von Podcastle.

heise-Podcaster nutzen für Titel, Vorspann- und Beschreibungstext die hauseigene Unternehmenslösung heise I/O. Hier kann sich jeder Podcaster für sein Produkt einen Prompt bauen, mit dem ein gewünschtes KI-Modell dann die Texte generiert – die der Podcaster wiederum überarbeitet.

## Vollmundige Ankündigungen

In dieser Übersicht haben wir bewusst einige Tools außen vor gelassen, die Podcastern auch das Sprechen abnehmen wollen. So heißt es zum Beispiel im Newsletter des deutschen KI-Ökosystems Mindverse: „Mit unserem KI-Agenten kannst du im Handumdrehen ein Dokument in einen spannenden Podcast verwandeln.“ Man müsse dafür nur ein Dokument oder einen Link einfügen, Host und Gast festlegen und Tonalität und Sprache wählen. Die

Funktion ist nur für zahlende Nutzer, im Wust etlicher KI-Engines auf der Mindverse-Website finden sich aber auch kostenlose Funktionen, die Podcastern helfen und beispielsweise ein Skript schreiben oder aus einer Audiodatei einen Blogeintrag generieren, ohne gleich den ganzen Podcast zu übernehmen.

Den Podcast-Generator haben wir der Vollständigkeit halber kurz ausprobiert, man findet ihn unter „KI Agenten“ in der Leiste links. Das Startformular ist so knapp wie im Newsletter versprochen. Aus den Angaben generiert die KI tatsächlich ein kurzes Podcast-Zwiesgespräch über das gewünschte Thema, inklusive Skript und Audiodatei. In unserem Test-Podcast auf Deutsch klangen die Stimmen teils authentisch, meist aber künstlich, eine hatte einen leichten amerikanischen Akzent. Inhaltlich war das Gespräch, das sich um einen sechsseitigen c't-Artikel drehte, in Ordnung, sprachlich aber zu euphorisch.

Ähnliche Versprechen wie Mindverse macht auch der teure Anbieter podcastAI, der in der Pro-Version „MagicPod“ für monatlich 600 Dollar automatisiert Podcasts mit der Stimme des Nutzers produzieren will. NotebookLM von Google dagegen generiert kostenlos einen englischsprachigen Gesprächspodcast mit zwei Sprechern. Dafür nutzt das Tool Quellen, die der Nutzer festlegt [5].

## Fazit

Jeder Nutzer muss für sich und seinen Podcast entscheiden, welches Tool am meisten nützt. So kann es zum Beispiel sein, dass Sendungen über Fachthemen einem KI-Tool Probleme bereiten, weil es viele Begriffe nicht oder falsch versteht. Dann muss der Produzent mühsam Transkript, Beschreibungstexte und Posts durchkämmen und Fehler korrigieren. Diese Arbeit kann es ihm aber wert sein, wenn



# OKR DAY

Mit gemeinsamen Zielen zu messbaren Erfolgen

13. März 2025 • Online

OKR ermöglicht es Unternehmen, klare, fokussierte Ziele (**Objectives**) mit messbaren Erfolgen (**Key Results**) zu verbinden.

### Darum dreht sich die Konferenz:

- Funktionsweise von OKR verstehen
- bewährte Best Practices kennenlernen
- OKR im eigenen Unternehmen implementieren

Für **agile Coaches**, **agiles Management** und **Führungskräfte**, **Scrum Master**, **OKR Master** sowie **OKR-Interessierte**.



Jetzt  
Tickets  
sichern!

[okr.inside-agile.de](https://okr.inside-agile.de)



## Adobe Podcasts

Kreativ-Experte Adobe will im Podcastmarkt mitmischen. Die Werkzeuge dafür sind noch in der Test-, das Aufnahmestudio sogar in der Beta-Phase. Gratis sind nur Audioprojekte, das Premium-Abo umfasst auch Videos. Verglichen mit anderen Tools ist Adobes Angebot bescheiden: Ton bereinigen mithilfe einer KI, die auch Tipps fürs optimale Mikrofon-Setting gibt, außerdem simple Aufnahmen mit Gästen und Transkripte auf Englisch, anhand derer man die Aufnahmen rudimentär bearbeitet.

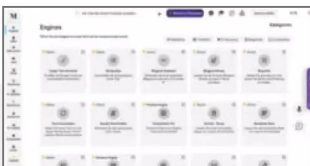
**Preis: Basisversion gratis, Premium 11,90 Euro monatlich**



## Auphonic

Der Dienst aus Österreich hat sich ganz auf die KI-gestützte Audio-Postproduktion spezialisiert und ist in diesem Bereich in der freien Podcast-Szene momentan führend. Schon die vorbereiteten Presets helfen gut und sehr schnell, gleichen beispielsweise Lautstärken an und filtern Störgeräusche weg. Wer mehr will, kann sich eigene Presets zusammenbauen. Mithilfe von Whisper transkribiert Auphonic auch, außerdem kann es Füllwörter und Stille entfernen und unterstützt Nutzer dabei, Kapitelmarken zu erstellen.

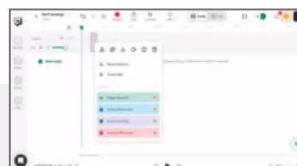
**Preis: monatlich 2 Stunden gratis; gestaffeltes Guthabensystem: 9 Stunden ab 11 Euro pro Monat**



## Mindverse

Das KI-Ökosystem aus Deutschland will sehr viel und wird dadurch unübersichtlich. Die meisten der über 70 KI-Engines arbeiten mit Texten, einige mit Audio, es gibt auch eine KI-Suchmaschine. Nach Hilfen für Podcasts muss man etwas suchen und findet die Engines „Podcast Audio zu Blog“, „Podcast Skript“, „Sprachsynthese – Text to Speech“ und „Audio Transkribieren“; aber auch die Funktionen „Korrigieren“ und „Zusammenfassen“ für Texte und „Social Media Posts“ können beim Podcasten helfen. Mindverse ist nach eigener Aussage DSGVO-konform.

**Preis: kostenlose Basisversion; Abos ab 16 Euro pro Monat**



## Podcastle

Seine kostenlose Basic-Testversion versteckt Podcastle: Wir landeten unerwartet darin, weil wir, statt unsere Zahlungsdaten anzugeben, oben rechts auf „Skip“ klickten. Podcastle lief bei uns in Firefox, nur das „Schließen“-Kreuzchen von Pop-ups war hier unsichtbar – Google Chrome zeigte es korrekt an. Der Dienst beherrscht Audio- und Videoaufnahmen mit Gästen und bearbeitet diese mithilfe von KI rudimentär. Beim Upload wählt man die gewünschte Sprache aus. Deutsche Aufnahmen transkribierte Podcastle in unseren Tests gut.

**Preis: 7-Tage-Testversion; beschränkter, kostenloser Basic-Plan; Abo ab 12 Euro pro Monat**



## Literatur

[1] Holger Bleich, Podcast-Produktion: Hochwertige Audioaufzeichnungen kostenlos mit Zencastr, c't 6/2023, S. 170

[2] Hartmut Gieselmann, Stefan Wischner, Greta Friedrich, Transkriptionsdienste: Whisper V3 im Vergleich mit Online-Diensten, c't 9/2024, S. 92

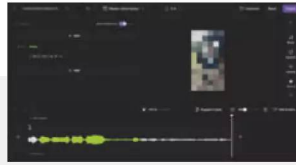
[3] Jan Philipp Wöbbeking, Kostenlos hochwertige KI-Musik erzeugen: Das leisten Suno und Udio, heise.de/-9687928

[4] Immo Junghärtchen, Spracherkennung und Transkription mit KI: Sprache in Text umwandeln mit Whisper, c't 14/2023, S. 140

[5] Jo Bager, Erste Schritte mit Googles KI-Assistent NotebookLM, c't 27/2024, S.148

Links zu den Tools und den erwähnten Artikeln:

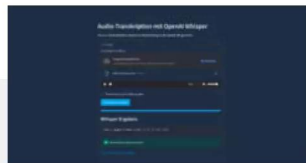
[ct.de/wk52](https://ct.de/wk52)



## Riverside

Als einziger der vorgestellten Dienste bietet Riverside neben der Browserversion auch Apps für iOS und Android an. Das Programm enthält ein kleines Aufnahmestudio, in das man auch Gäste einladen kann. Bevor man Riverside transkribieren lässt, legt man im Studio die gewünschte Sprache fest – dann klappt es zum Beispiel auch mit Deutsch. Füllwörter und Pausen zeigt das Werkzeug direkt im Transkript, mit dessen Hilfe Nutzer die Aufnahme schneiden. Automatische Texte und Clips für Social Media sollen bei der Vermarktung helfen.

**Preis: 2 Stunden gratis, Abo ab 15 US-Dollar pro Monat**



## Whisper

Die Spracherkennungs-KI von OpenAI ist kostenlos und quelloffen, allerdings nicht einfach so nutzbar. Man muss sie entweder selbst hosten (das Bild zeigt die heise-Version) oder über die Programmierschnittstelle von OpenAI verwenden. Letzteres bieten manche Dienste wie zum Beispiel Auphonic oder auch „Whisper Transcription“ von Good Snooze (lokal auf Macs). Alternativ richten Nutzer Whisper via Kommandozeile ein [4]. Das Werkzeug liefert ziemlich gute Transkripte, auch auf Deutsch.

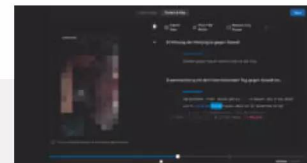
**Preis: kostenlos**



## StreamYard

StreamYard hat sich auf Aufnahmen und Live-streams spezialisiert. Der Webdienst kann mit bis zu zehn Teilnehmenden zugleich umgehen. Die Aufnahme läuft jeweils lokal, damit die Internetqualität nicht die Qualität von Ton und Bild beeinträchtigt. Streamer können mehrere Plattformen gleichzeitig bespielen, darunter Twitch und Facebook. KI-gestützt bearbeitet StreamYard Tonspuren nach und entfernt zum Beispiel Hintergrundgeräusche. Außerdem erzeugt es Clips für Social Media. Seine Abopreise zeigt das Tool nur angemeldeten Nutzern.


**Preis: kostenlose Basisversion; Abo ab 36 US-Dollar pro Monat**




## Zencastr


In der Zencastr-Aufnahmeumgebung kann man Gäste einladen und Videos zum Beispiel mit Soundeffekten würzen. Der Dienst generiert nach der Aufnahme automatisch ein gegliedertes Transkript und erkennt auch die Sprache. Nutzer können den Text bearbeiten und darin die Aufnahme schneiden. Clips für die sozialen Medien generiert Zencastr mithilfe von KI, alternativ markiert der Nutzer die gewünschte Stelle im Skript. Der fertige Podcast geht direkt über Zencastr online, via RSS-Feed und eine eigene Podcast-Seite.


**Preis: 14-Tage-Testversion; Abo ab 18 US-Dollar pro Monat**


 Podcast

Datei

 Lokaler Upload

 Dateilink einfügen

 Testtext.pdf  
PDF • 832.29KB



Tonalität

Humorvoll

Name des Hosts  
(Weiblich)


Isabel Grünewald (weiblich)

Name des Gast  
(Weiblich)

Greta Friedrich (weiblich)

Sprache

Deutsch

 Chat starten

Das KI-Ökosystem Mindverse aus Deutschland umfasst auch einen Podcast-Generator. Dieser bastelt für zahlende Nutzer ein Podcast-Zwiesgespräch zu einem Dokument oder einer Webseite.

Podcast-Tools mit KI				
Name	Adobe Podcasts	Auphonic	Mindverse	Podcastle
Hersteller, URL	Adobe Systems Software Ireland Limited, podcast.adobe.com	Auphonic GmbH, auphonic.com	Relativity GmbH, mind-verse.de	Podcastle Inc., podcastle.ai
Plattformen	Browser („Studio“ nur in Google Chrome; „Sprache verbessern“ auch in Safari, Edge, Firefox)	Browser	Browser (Website und Chrome-Erweiterung), Desktop (Windows, macOS, Linux)	Browser (Chrome <sup>3</sup> , Firefox <sup>4</sup> ), iOS ab 14.0
Datenschutzerklärung auf Deutsch	✓	— (EN)	✓	— (EN)
Registrierung via ...	Google / FB / Apple / Mail	Google / Apple / Mail	Google / Microsoft / Mail	Google / FB / Apple / Mail
Bedienoberfläche auf Deutsch	✓	— (EN)	✓	— (EN)
Funktionen				
unterstützt Audio / Video	✓ / ✓ <sup>4</sup>	✓ / ✓	✓ / ✓	✓ / ✓
Aufnahme	✓	✓ <sup>8</sup>	—	✓
Transkript	✓ <sup>5</sup>	✓	✓	✓
Schnitt	✓ <sup>1</sup>	— <sup>9</sup>	—	✓
Audio-Postproduktion	✓ <sup>1</sup>	✓	—	✓ <sup>1</sup>
Social-Media-Clips	—	— <sup>10</sup>	—	—
Beschreibungstexte generieren	—	✓	✓	—
Kapitelmarken	—	✓	—	—
direkte Veröffentlichung	—	✓	—	✓
Kosten				
Preis	Basisversion gratis, Premium 11,90 €/Monat	monatlich 2 Stunden gratis; gestaffeltes Guthabensystem <sup>6</sup> : 9 Stunden ab 11 €	kostenlose Testversion; Abo ab 16 €/Monat	7-Tage-Testversion; begrenzter <sup>7</sup> , kostenloser Basic-Plan; Abo ab 12 €/Monat
<sup>1</sup> nur eingeschränkt <sup>2</sup> Anzeigefehler <sup>3</sup> vom Anbieter empfohlen <sup>4</sup> im kostenpflichtigen Abo <sup>5</sup> nur auf Englisch <sup>6</sup> Wahl zwischen einmaligen Käufen und Abo <sup>7</sup> begrenzt sind Aufnahmedauer, Transkriptionsmenge, Zeichenzahl Text-to-Speech, Speicherplatz <sup>8</sup> simpler Audio-Recorder <sup>9</sup> Funktion für 2025 geplant <sup>10</sup> keine Videoclips, aber Audiogramme				

110 KI-Praxis im Alltag

c't KI-Wissen 2025

er zum Beispiel die Möglichkeit schätzt, aus dem (korrigierten) Transkript heraus Clips für Social Media zu generieren.

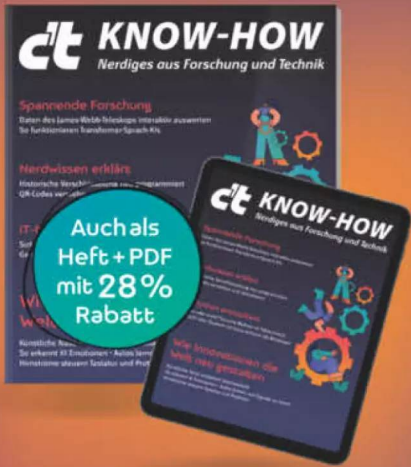
Glücklicherweise haben die hier näher vorgestellten Tools kostenlose Basis- oder zumindest Testversionen. Die günstigsten Abopläne bewegen sich bei fast allen um die 15 Dollar beziehungsweise Euro pro Monat, StreamYard ist mit monatlich 45 Euro teurer. Das kurz erwähnte podcastAI fällt aus dem Rahmen, ist mit monatlich 99 Dollar für die günstigere Variante sehr teuer und bietet keine kostenlose (Test-)Version. Whisper ist dagegen kostenlos und quelloffen, dafür ist die Nutzung komplizierter.

Für manche Anwendungsfälle kann sich auch eine Kombination mehrerer Tools lohnen. So zeichnet zum Beispiel heise viele seiner Podcast-Episoden mit StreamYard auf, überarbeitet die Tonspur mit Audiponic, transkribiert sie mit Whisper und generiert mit der hausinternen KI dann Texte zum Podcast, etwa Titel, Vorspann und Beschreibungstext. (gref) **ct**

	Riverside	StreamYard	Whisper	Zencastr
	RiversideFM, Inc., riverside.fm	StreamYard, Inc., streamyard.com	OpenAI, openai.com/index/whisper	Zencastr, Inc., zencastr.com
	Browser (Chrome, Edge), iOS ab 16.0, Android ab 7.0	Browser (Desktop: Chrome, Firefox, Edge, Opera; Android: Chrome, Firefox; iOS: Safari) <sup>3</sup>	eigener Server, API	Browser (Brave, Chrome, Edge), iOS ab 16.0
	— (EN)	— (EN)	— (EN)	— (EN)
	Google / Apple / Spotify / Mail	Google / Mail	—	Mail
	— (EN)	✓	keine eigene Bedienoberfläche	— (EN)
	✓ / ✓	✓ / ✓	✓ / ✓	✓ / ✓
	✓	✓	—	✓
	✓	✓ <sup>5</sup>	✓	✓
	✓	—	—	✓
	✓ <sup>1</sup>	✓ <sup>1</sup>	—	✓ <sup>1</sup>
	✓	✓	—	✓
	✓	—	✓	—
	✓	—	—	✓
	✓	✓	—	✓
	2 Stunden gratis, Abo ab 15 US-\$/Monat	kostenlose Basisversion; Abo ab 36 US-\$/Monat	kostenlos (Open Source)	14-Tage-Testversion; Abo ab 18 US-\$/Monat
	✓ vorhanden — nicht vorhanden			

# Geballtes Nerd-Wissen!

Spannendes aus Forschung und Technik



- Künstliche Nase analysiert Gasmoleküle
- So erkennt KI Emotionen
- Daten des James-Webb-Teleskops interaktiv auswerten
- Historische Verschlüsselung neu programmiert
- Sicher oder nicht? Security-Mythen im Faktencheck

Heft für 14,90 € • PDF für 12,99 €  
Heft + PDF 19,90 €

 [shop.heise.de/ct-knowhow24](https://shop.heise.de/ct-knowhow24)

Generell portofreie Lieferung für Heise Medien- oder Maker Media Zeitschriften-Abonnenten oder ab einem Einkaufswert von 20 € (innerhalb Deutschlands).  
Nur solange der Vorrat reicht. Preisänderungen vorbehalten.



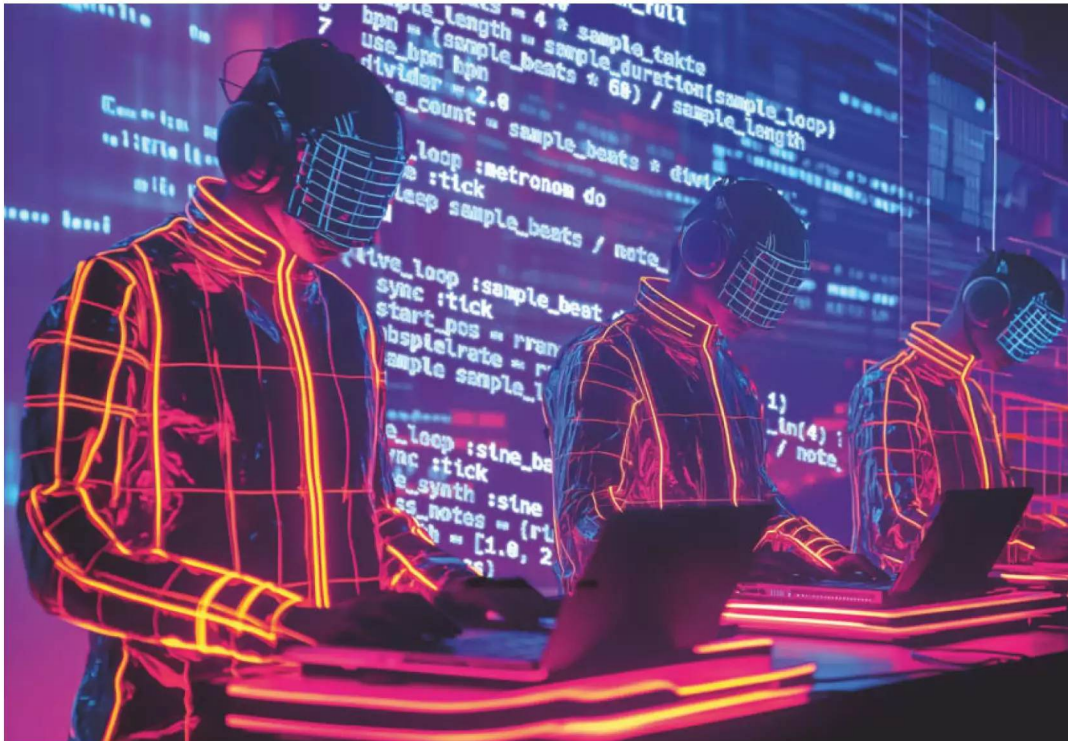


Bild: KI, Collage c't

# Musik live coden mit Sprachmodellen

Mit ein paar Zeilen Code programmieren Sie im kostenlosen Sonic Pi einen coolen Elektro-Track. Dank umfangreicher Tutorials lernen insbesondere Schüler dabei spielend die Grundlagen von Musik und Informatik. Wir zeigen, wo ChatGPT bei aufwendigeren Songs helfen kann.

Von **Hartmut Gieselmann**

**S**eit nunmehr zwölf Jahren können Schüler, DJs und andere Musikbegeisterte in dem kostenlosen Open-Source-Programm Sonic Pi mit Klängen und Songs experimentieren. Nicht nur, dass sich das Programm kinderleicht auf Windows-Rech-

nern, Macs und sogar Raspberry Pis installieren lässt. Dank der eingebauten Tutorials und der Fähigkeit des integrierten Compilers, Änderungen am Code quasi in Echtzeit umzusetzen, kann man seine Musikstücke live verändern und den pumpenden Beats

zuhören, während man eine neue Bassline bastelt. Das macht ungeheuren Spaß und lädt zum Experimentieren ein.

Diese Algorave genannte Musikrichtung hat in den vergangenen Jahren viele Freunde gewonnen (siehe Kasten „Die Algorave-Bewegung“). Um mitzumachen, bietet Sonic Pi den leichtesten Einstieg, denn die Programmierumgebung lässt sich mit ein paar Mausklicks installieren. Einsteiger nimmt Sonic Pi mit ausführlichen Tutorials in Deutsch und anderen Sprachen an die Hand. Man braucht tatsächlich null Vorkenntnisse: weder in Musiktheorie noch im Programmieren der Skriptsprache Ruby, die Sonic Pi verwendet.

Wer nach den ersten Experimenten, wie wir sie in unseren Einstiegsartikeln in [1], [2] und [3] erklärt haben, komplexere Stücke mit mehreren Song-Abschnitten und Instrumenten komponieren will, muss sich beispielsweise mit der Synchronisation der sogenannten Live-Loops auseinandersetzen und Notenwerte in Listen verwalten. Dabei kommen schnell weit über hundert Zeilen Code zusammen. Und genau hier kann ein generatives Sprachmodell wie ChatGPT helfen, den Überblick zu behalten und

die einzelnen Programmteile zu verknüpfen. Dieser Artikel erklärt, wie man dabei grundsätzlich vorgeht, und stellt ein Template vor, das dabei hilft, komplexere Songs zu organisieren.

## Wie musikalisch ist ChatGPT?

Bevor wir darauf eingehen, wie ein Sprachmodell beim Musizieren mit Sonic Pi helfen kann, müssen wir zunächst einmal klarstellen, was ein Sprachmodell nicht kann. Zuerst kann es keine Musik hören und analysieren. Wenn man ChatGPT fragt: „Wie programmiere ich die Kickdrum-Sequenz aus dem Song *Blue Monday* von New Order?“, gibt der Chatbot zwar eine Kickdrum-Sequenz aus, den typischen Rhythmus des besagten Stücks hält diese aber nicht ein. Und auch bei Melodien oder Bass-Sequenzen ist es eher Zufall, wenn das ChatGPT-Ergebnis passende Töne trifft. Das Sprachmodell patzt auch oft, wenn es darum geht, Notenwerte und Takte zusammenzuzählen.

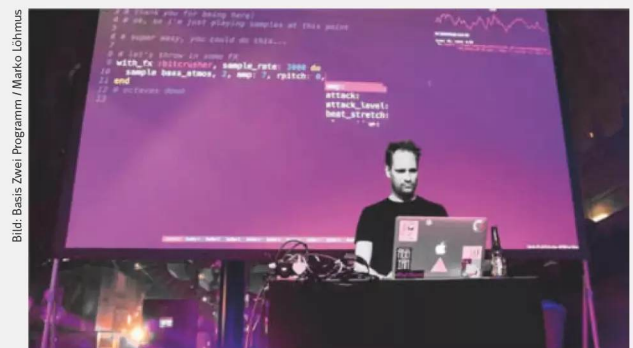
Befiehlt man ChatGPT, es solle einen Hit im Stil von Depeche Mode in Sonic Pi programmieren, stellt sich schnell Ernüchterung ein: Der Code produziert

## Die Algorave-Bewegung

Als sogenannter Algorave hat sich weltweit eine neue Musikszene entwickelt. Unter dem Dach der Toplap-Organisation veranstalten weltweit mehr als 2000 Musiker rund um den Globus Konzerte und Raves, bei denen sie meist keine Schlagzeuge, Gitarren oder Plattenspieler bedienen, sondern vor den Augen des Publikums Code in ihren Laptop hacken. Neben Sonic Pi nutzen sie dazu auch andere Open-Source-Entwicklungsumgebungen wie TidalCycles zusammen mit der Plattform SuperCollider zur Audio-Synthese; oder auch Hydra und Gibber, die zur Musik pulsierende Grafiken und Farbmuster per Beamer auf Wände projizieren.

Alle nur erdenklichen Programmier- und Skriptsprachen kommen dabei zum Einsatz: Ruby, Haskell, ChuckK, Max, Moxie, Nyquist, Pure Data und so weiter. Auch dabei kann eine KI helfen, den richtigen Syntax zu finden. Weil sämtliche Werkzeuge auf Plattformen wie Github kostenlos verfügbar sind, gibt es kaum materielle Hürden, um einen

Algorave zu starten. Die wichtigsten Ressourcen haben wir unter [ct.de/wuj8](https://ct.de/wuj8) verlinkt.



**Sonic-Pi-Erfinder Sam Aaron codet Musik live auf der Bühne.**



So sieht die Programmoberfläche von Sonic Pi aus. Im Fenster oben links geben Sie Ihren Code ein, rechts daneben zeigt Sonic Pi die aktuelle Wellenform und darunter aktuelle Statusmeldungen an. Die umfangreichen Tutorials können Sie im unteren Fenster durcharbeiten.

zwar einen lauffähigen Song, der aber einen lahmen Grundrhythmus hat und mit einer Synth-Stimme immer wieder die gleichen vier Töne als Melodie wiederholt – das langweilt bereits nach wenigen Sekunden.

Die Erklärung ist einfach: ChatGPT wurde zwar mit Millionen Webseiten und Codebeispielen trainiert, kann aber nicht zwischen einer guten und einer schlechten Komposition unterscheiden. Daher generiert es nur sehr einfache Abschnitte, die zwar von der Syntax her oft stimmen, aber niemanden auf die Tanzfläche locken.

Die musikalische Idee und das Rhythmusgefühl müssen also vom Musiker kommen, der seinen Track in Sonic Pi programmiert – egal ob mit KI oder ohne. Für welche Zwecke er die KI einspannt, hat gewaltige Auswirkungen darauf, wem der fertige Song

schließlich gehört: An den Melodien aus Text-to-music-Generatoren wie Udio oder Suno, die Hit-Kopien auf einfache Textanweisungen hin fabrizieren, können allenfalls die Künstler, mit deren Material diese Generatoren trainiert wurden, eine Urhebererschaft geltend machen – jedoch nicht der Nutzer, wenn er lediglich einen kurzen Prompt als Auftragsanweisung geschrieben hat [4]. Wenn sich die KI-Hilfe jedoch nur auf die technische Umsetzung einer Idee beschränkt, bleiben die Rechte am Song bei dem, der ihn mit Codezeilen komponiert hat.

## Programmierhilfen

Wenn Sie die mitgelieferten Tutorials von Sonic Pi durchgearbeitet haben, kennen Sie zwar die wichtigsten Grundelemente. Sie wissen jedoch noch



nicht, wie Sie daraus einen umfangreicheren Song bauen. Und genau hier kann ChatGPT helfen. Starten Sie mit einer eigenen kleinen Idee und schreiben Sie die ersten Codezeilen selbst. Sobald Sie nicht mehr weiterkommen, weil Sie die nötigen Befehle oder deren Syntax noch nicht kennen, formulieren Sie das, was Sie tun wollen, auf Deutsch und geben diese Beschreibung zusammen mit Ihrem Anfangscode an ChatGPT weiter.

Da ChatGPT auch beim Verfassen von Code zum Halluzinieren neigt, gehen Sie am besten Schritt für Schritt vor. Dann können Sie den Code nachvollziehen und eventuelle Fehler schneller finden.

Dazu ein Beispiel:

Unser Song soll das in Sonic Pi mitgelieferte Sample vom Amen-Break immer wieder zufällig neu zusammensetzen und so einen sich verändernden Drum-&Bass-Rhythmus als Grundlage schaffen.

Das genannte Amen-Break können Sie leicht selbst einbauen:

```
live_loop :amen_brother do
  sample :loop_amen_full
  sleep 16
end
```

Wenn Sie das Amen-Break in voller Länge mitzählen, kommen Sie auf vier 4/4-Takte. Da eine Pause mit der Länge 1 genau eine Viertelnote lang ist, wartet der Befehl `sleep 16` also 16 Schläge, bevor die `live_`

loop wieder zum Anfang springt und das Sample erneut abspielt. Von Haus aus arbeitet Sonic Pi aber mit dem Tempo von 60 bpm (beats per minute, Schläge pro Minute); das ist deutlich langsamer, als das Amen-Break eingespielt wurde. Deshalb ist auch eine lange Pause zwischen den wiederkehrenden Sample-Wiedergaben zu hören.

## bpm berechnen

Als Erstes sollte ChatGPT also zeigen, wie man die genaue Länge des Sample-Loops herausfindet und entsprechend das Tempo des Songs anpasst. Wenn wir ChatGPT bitten, genau dies zu tun, spuckt es nicht nur den erweiterten Code aus, sondern erklärt ihn auch Zeile für Zeile. Nach ein paar Iterationen, bei denen wir die Vorschläge ausprobierten und ChatGPT immer wieder mitteilten, was noch nicht funktionierte, hatten wir uns folgenden Code erarbeitet:

```
sample_loop = :loop_amen_full
sample_takte = 4.0
sample_beats = 4 * sample_takte
sample_length = sample_duration(sample_loop)
bpm = (sample_beats*60)/sample_length
use_bpm bpm
```

```
live_loop :sample_beat do
  sample sample_loop
```

# WILLKOMMEN IM NEUEN IOT-ÖKOSYSTEM

Mit LoRaWAN und C-Programmierung über lange Distanzen messen und steuern



Heft +  
LoRaWAN-  
Set

Make Special LoRaWAN  
inkl. Experimentierset für 64,90 €



[shop.heise.de/make-lorawan24](https://shop.heise.de/make-lorawan24)

JETZT  
BESTELLEN!



```
sleep sample_beats
end
```

Der Befehl `sample_duration` ermittelt die Länge des Samples in Sekunden. Zusätzlich geben wir an, wie viele Takte das Sample dauert (4.0), sodass das Skript die Anzahl der beats per minute (bpm) berechnet und als Tempo festlegt (140 bpm beim Amen-Break). In der `live_loop` wartet der `sleep`-Befehl dann immer die Anzahl der Viertelnotenschläge des Samples ab, bevor die Schleife von Neuem startet. Nach demselben Prinzip können Sie auch andere Samples als Drum-Loop einsetzen. Dazu müssen Sie nur den Namen von `:loop_amen_full` austauschen und dessen Taktlänge in `sample_takte` eintragen.

Als Nächstes soll ChatGPT den Code ergänzen. Es soll das Sample bei der Wiedergabe in Achtelnoten zerhacken und diese in einer zufälligen Reihenfolge wiedergeben. Eine Achtelnote dauert genau den halben Schlag einer Viertelnote. Die `live_loop` aus dem obigen Beispiel ändert sich dann folgendermaßen:

```
live_loop :sample_beat do
  divider = 2.0
  note_count = sample_beats * divider
  start_pos = rand_i(0, note_count-1)
  sample sample_loop, start: start_pos
  start_pos / note_count.to_f,
  finish: (start_pos + 1) / note_count.to_f
  sleep 1 / divider
end
```

Als `divider` ist hier 2.0 eingetragen, um die Viertelnoten in Achtelnoten zu zerlegen. Anschließend berechnet das Skript die Anzahl der Achtelnoten des Samples und daraus eine zufällige Startposition. Der `sample`-Befehl spielt dann den Sample-Teil zwischen der zufälligen Startposition und der Endposition ab. Deren Werte werden durch den Teiler `note_count.to_f` auf Werte zwischen 0 und 1 normalisiert. Schließlich wartet der `sleep`-Befehl eine Achtelnote lang und springt wieder an den Anfang. Probieren Sie aus, was passiert, wenn Sie den `divider` auf 0.5, 1.0 oder 4.0 setzen.

Schließlich soll ChatGPT die Achtelnoten in zufälligen Geschwindigkeiten im Bereich von 0.7 bis 1.3 vorwärts oder rückwärts abspielen. Dazu berechnet ChatGPT eine zufällige Abspielrate, die entweder positiv oder negativ ausfällt, und ergänzt den `sample`-Befehl um den Parameter `rate`:

```
abspielrate = rand(0.7, 1.3) *
  (one_in(4) ? -1 : 1)
sample sample_loop, start: start_pos
  / note_count.to_f, finish:
  (start_pos + 1) / note_count.to_f,
  rate: abspielrate
```

Die 4 in `one_in(4)` gibt an, dass nur jeder vierte Abschnitt mit einer negativen Rate, also rückwärts abgespielt werden soll. Sie können dort auch höhere Integer-Werte eintragen, um den Anteil der rückwärts abgespielten Abschnitte zu verringern.

## Metronom für alle

Damit steht die `live_loop` mit dem rhythmischen Grundgerüst. Diese soll anschließend einen Basslauf bekommen, der synchron zu der Sample-Loop läuft. Sonic Pi kann zwar von sich aus mehrere Live-Loops starten. Da es aber unterschiedlich lange dauern kann, verschiedene Befehle abzuarbeiten, divergieren die Loops bald, wenn sie nicht von vornherein synchronisiert sind. Auf langsamen Systemen reduzieren Sie die Rechenlast etwas, indem Sie dem Skript den Befehl `use_debug false` voranstellen.

Zudem benötigen Sie eine `live_loop`, die ein Synchronisationssignal ausgibt. Das umgestellte Skript mit der Synchronisations-Loop sieht dann so aus:

```
use_debug false
sample_loop = :loop_amen_full
sample_takte = 4.0
sample_beats = 4 * sample_takte
sample_length = sample_duration(sample_loop)
bpm = (sample_beats * 60) / sample_length
use_bpm bpm
divider = 2.0
note_count = sample_beats * divider

live_loop :metronom do
  cue :tick
  sleep 1 / divider
end

live_loop :sample_beat do
  sync :tick
  start_pos = rand_i(0, note_count-1)
  abspielrate = rand(0.7, 1.3) *
    (one_in(4) ? -1 : 1)
  sample sample_loop, start: start_pos
    / note_count.to_f, finish:
    (start_pos+1) / note_count.to_f,
```



```

rate: abspielrate

end

```

Die :metronom-Loop gibt per cue einen :tick aus. Die :sample\_beat-Loop wartet per sync auf diesen :tick, bevor es das Sample abspielt.

Nun kann man weitere Loops ergänzen, beispielsweise mit einem Basslauf. Für Drum & Bass eignet sich etwa der Sinus-Synth. Dazu gibt man die Noten an (etwa e2, g2, a2) und packt diese in eine Ring-Liste, damit man sie leicht austauschen kann. Schließlich sollen die Bassnoten nicht immer nur den Achtelnoten der Metronom-Schleife entsprechen. Deshalb definieren wir eine Liste für die Variable length, die die Länge der Noten in Viertelschlägen angibt. Bei jedem Aufruf von length wird ein zufälliger Eintrag aus 1.0, 2.0 und 4.0 ausgewählt. Die Bass-Loop kann man dann einfach hinter die Beat-Loop hängen:

```

live_loop :sine_bass do
  sync :tick
  use_synth :sine

```

```

bass_notes = (ring :e2, :g2, :a2)
length = [1.0, 2.0, 4.0].choose
tick(:bass)
if (look(:bass) % (divider*length)).zero?
  play bass_notes.tick, sustain: length
end
end

```

tick(:bass) lässt einen weiteren Tick in der Loop namens :bass hochzählen. Erst wenn der Rest aus der Division des aktuellen Ticks (abgefragt mit look(:bass)) und (divider\*length) gleich Null ist, erklingt eine neue Note. Deren Tonhöhe schaut Sonic Pi in der Ring-Liste bass\_notes nach, zählt mit jedem tick eine Stelle weiter und beginnt am Ende wieder von vorn.

Wer die Bassnoten zufallsgesteuert variieren will, kann die Noten ähnlich wie die Länge zufällig wählen lassen, zum Beispiel:

```
bass_notes = [:e2, :g2, :b2].choose
```

Die zugehörige play-Zeile lautet dann:



**JETZT IM ABO GÜNSTIGER LESEN**

**GRATIS!**

2x Make testen mit über 30 % Rabatt

#### Ihre Vorteile im Plus-Paket:

- ✓ Als Heft und
- ✓ Digital im Browser, als PDF oder in der App
- ✓ Zugriff auf **Online-Artikel-Archiv**
- ✓ **Geschenk**, z. B. Make: Tasse

Für nur 19,90 € statt 29 €

Jetzt bestellen:

**make-magazin.de/miniabo**





```
play bass_notes, sustain: length
```

Spielen Sie damit ein wenig herum. Probieren Sie andere Noten, Längen und Synthesizer abseits des :sine aus.

## Weitere Instrumente

Nach dem Prinzip dieser Bass-Loop können Sie weitere Instrumente in neuen Loops hinzufügen. Fangen Sie dazu die Loop immer mit `sync :tick` an, definieren Sie dann den Synthesizer, listen die Töne in einer Liste auf und lassen die Noten spielen.

Wenn Sie mit einem Dividervon 4.0 arbeiten, gibt die :metronom-Schleife das Tempo von Sechzehntelnoten aus. Dann können Sie in Listen beispielsweise Pattern aus 16 Noten notieren. Wenn Sie Pausen einbauen wollen, schreiben Sie in die Liste ein :r und ergänzen eine if-Bedingung bei der Wiedergabe, beispielsweise:

```
lead_notes = (ring :e4, :e4, :e4, :r, ♩
               ♩ :g4, :g4, :g4, :r, :b4, :b4, :b4, ♩
               ♩ :r, :c5, :c5, :c5, :r)

tick(:lead)
if lead_notes != :r
  play lead_notes.tick
end
```

Auch hier können Sie für lead\_notes wie zuvor alternativ eine Liste mit Notenwerten in [ ] angeben, aus der Sonic Pi mit einem angehängten .choose einen zufälligen Wert wählt.

Mit wachsender Anzahl von Live-Loops wird der Code schnell unübersichtlich. Definieren Sie die Notenfolgen dann besser im vorderen Teil des Skripts. So können Sie diese bei einem Live-Konzert schneller ändern, ohne zu den jeweiligen Loops scrollen zu müssen. Um den Überblick zu behalten, sollten Sie darauf achten, dass die Variablen für die Notenlisten auf den Namen des Instruments hinweisen (etwa bass\_notes, lead\_notes und so weiter).

## Song-Strukturen aufbauen

Etwas größer werden die Umbaumaßnahmen, wenn Sie die einzelnen Loops nicht manuell von Hand starten und stoppen sowie Noten austauschen wollen, sondern eine feste Song-Struktur mit Intro, Strophe, Refrain und Outro abspielen lassen. Dazu müssen Sie einige Variablen und Listen definieren.

Als Erstes eine Liste, in der Sie festlegen, welches Instrument in welchem Songabschnitt spielen soll:

```
sections = {
  intro: { drums: false, bass: true },
  verse: { drums: true, bass: true },
  outro: { drums: false, bass: true }
}
```

Dann legen Sie fest, wie viele Takte jeder Abschnitt dauern soll:

```
section_lengths = {
  intro: 4.0,
  verse: 8.0,
  outro: 4.0
}
```

Jetzt definieren Sie, mit welcher Sektion der Song starten soll:

```
current_section = :intro
```

Nach der :metronom-Loop bauen Sie eine Loop ein, die die Variable current\_section von :intro auf :verse und schließlich auf :outro setzt, sobald die jeweiligen Sektionslängen durchlaufen wurden:

```
define :play_section do |section|
  current_section = section
  takte = section_lengths[section] * 4 * divider
  takte.times do
    sync :tick
  end
end
```

Den gesamten Songablauf legen Sie schließlich über einen Thread-Abschnitt fest. Dieser stoppt die :metronom-Loop, sobald das Outro durchgespielt wurde:

```
in_thread do
  play_section(:intro)
  play_section(:verse)
  play_section(:outro)
  live_loop :metronom do
    stop
  end
end
```

Damit jedes Instrument in seiner Live-Loop auf die richtigen Noten der aktuellen Song-Sektion zugreift,

## Literatur

[1] Gerhard Völkl, Vom Programmierer zum DJ, Mit Python und Sonic Pi elektronische Dance Music komponieren, c't 17/2016, S. 164

[2] Pit Noack, Programmier deinen Song, Erste Schritte mit der Musik-Software Sonic Pi, c't 12/2017, S. 144

[3] Pit Noack, Beats basteln wie die Großen, mit Sonic Pi eine Rhythmusspur programmieren, c't 13/2017, S. 132

[4] Nik Akinci, Joerg Heidrich, Im rechtsfreien Raum, Urheberrecht und Datenschutz bei KI-generierten Werken, c't 5/2023, S. 70

Template und  
Algorave-Quellen

[ct.de/wuj8](https://ct.de/wuj8)

```

# Reich Phase

# Steve Reich's Piano Phase
# See: https://en.wikipedia.org/wiki/Piano_Phase

# use_synth :piano
notes = (ring :E4, :Fs4, :B4, :Cs5, :D5, :Fs4, :E4, :Cs5, :B4, :Fs4, :D5, :Cs5)

live_loop :slow do
  play notes.tick, release: 0.1
  sleep 0.3
end

live_loop :faster do
  play notes.tick, release: 0.1
  sleep 0.295
end

```

**Der Komponist Steve Reich experimentiert mit leicht abweichenden Tempi und Rhythmen. Hier erzeugt er mit zwei leicht unterschiedlich ablaufenden Notensequenzen ständig neue Muster. Versuchen Sie, den Code um eine dritte Stimme zu erweitern und verändern Sie die Wartezeiten des sleep-Befehls.**

bauen Sie die live\_loop-Routinen nach folgendem Muster auf (hier am Beispiel des Basses mit Sechzehntelnoten):

```

live_loop :bass_line do
  sync :tick
  use_synth :fm
  tick(:bass)
  if sections[current_section][:bass]
    note = bass_notes[current_section].tick
    if note != :r
      play note
    end
  end
end

```

Zunächst prüft die if-Bedingung, ob das Instrument in der aktuellen Sektion überhaupt spielen soll. Falls ja, dann übergibt sie der Variable note den aktuellen Wert aus der Ring-Liste bass\_notes mit dem Namen der aktuellen Sektion. Schließlich prüft eine weitere if-Bedingung, ob nicht eine Pause (:r) gemacht werden soll, und gibt anderenfalls mit play die Note wieder.

Wir haben den gesamten Ablauf mit einigen Zusätzen in einem ausführlich kommentierten Beispiel-Template zusammengestellt, mit dem Sie Ihre Songs strukturieren können (siehe [ct.de/wuj8](https://ct.de/wuj8)). Betrachten Sie die dort eingesetzten Noten als Platzhalter und fangen Sie an zu experimentieren.

## Jenseits der Sequencer

Diese Template-Struktur orientiert sich an typischen Sequencern, die Sechzehntel- oder Achtelnoten im festen 4/4-Takt abspulen. So arbeiten für gewöhnlich Digital Audio Workstations (DAW) oder auch MIDI-Tracker. Der Vorteil ist, dass man die Strukturen leicht nachvollziehen kann. Der Nachteil: Es entsteht ein sehr enges Korsett für den Rhythmus, die Melodien und die Song-Struktur. Das passt gut zu Dance-Tracks mit festen Beats, weniger aber zu Klangexperimenten, die solche Strukturen aufbrechen.

Deshalb sollten Sie sich nicht damit zufriedengeben, mit Sonic Pi einen Tracker oder eine DAW nachzuahmen. Dabei stoßen Sie nämlich schnell an die Grenzen der Entwicklungsumgebung. Nutzen Sie stattdessen die Vorteile und Besonderheiten der algorithmischen Generierung und Manipulation.

Dazu gibt es verschiedene Methoden. Sie können zum Beispiel die Länge der Ring-Listen von 16 Schritten auf weniger Schritte mit einer ungeraden Schrittanzahl verkürzen. Wenn beispielsweise eine Kick-Drum vier Schritte nutzt, die Hi-hat nur 3 und die Snare 5 Schritte, ergeben sich mit sehr kurzen Zeilen Code ungewöhnliche rhythmische Muster, die sich erst nach langer Zeit wiederholen.

Statt Melodien vorzugeben, können Sie Sonic Pi auch zufällige Töne aus einer Liste wählen lassen.

```
note = [:a2, :c2, :d2, :e2].choose
```

## # Monday Blues

```
# Coded by Sam Aaron

use_debug false
load_samples [:drum_heavy_kick, :drum_snare_soft]

live_loop :drums do
  puts "slow drums"
  6.times do
    sample :drum_heavy_kick, rate: 0.8
    sleep 0.5
  end

  puts "fast drums"
  8.times do
    sample :drum_heavy_kick, rate: 0.8
    sleep 0.125
  end
end

live_loop :synths, delay: 6 do
  puts "how does it feel?"
  use_synth :mod_saw
  use_synth_defaults amp: 0.5, attack: 0, sustain: 1, release: 0.25, mod_range: 12, mod_phase: 0.5, mod_invert_wave: 1
  notes = (ring :F, :C, :D, :D, :G, :C, :D, :D)
  notes.each do |n|
    tick
    play note(n, octave: 1), cutoff: (line 90, 130, steps: 16).look
    play note(n, octave: 2), cutoff: (line 90, 130, steps: 32).look
    sleep 1
  end
end

live_loop :snare, delay: 12.5 do
  sample :drum_snare_soft
  sleep 1
end
```

In den Beispielen von Sonic Pi finden Sie den „Monday Blues“. Probieren Sie, den Code auf die im Text erklärte Loop-Synchronisation umzustellen, und ergänzen Sie weitere Stimmen aus dem Hit von New Order.

Da Sonic Pi mit Pseudo-Zufallszahlen arbeitet, sind die Ergebnisse reproduzierbar, die ungewöhnlichen Stücke klingen also bei jedem Neustart gleich.

Der Sonic-Pi Entwickler Sam Aaron hat dazu zahlreiche Beispiele verfasst, die Sie zusätzlich zum Tutorial in Sonic Pi finden. Erstaunlich ist etwa „Haunted“, das mit nur vier Zeilen Code eine endlose Reihe von Glockentönen erzeugt:

```
live_loop :haunted do
  sample :perc_bell, rate: rand(-1.5, 1.5)
  sleep rand(0.1, 2)
end
```

Oder folgen Sie dem Komponisten Steve Reich, der zwei leicht voneinander abweichende Synchronisationsspuren kombiniert, um sich fortwährend ändernde Polyrhythmen zu erzeugen (Beispiel Reich Phase).

Sie können ChatGPT mit solchen Beispielen füttern und das Modell auffordern, diese zu variieren. Dabei unterlaufen dem Sprachmodell unweigerlich Fehler, aber auch glückliche Unfälle, die Sie auf neue musikalische Ideen bringen. So wandeln Sie den Nachteil, dass Sprachmodelle oft halluzinieren und Fehler machen, in einen kreativen Vorteil.

## Fehler finden

Wenn ein vom Sprachmodell vorgeschlagener Code nicht läuft, können Sie den Fehler mit ein paar Hilfsmitteln einkreisen und korrigieren. Zunächst sollten Sie Ihren Code von einer kleinen Idee oder einem Motiv ausgehend immer schrittweise langsam erweitern. Nehmen Sie als Keimzelle ein paar Zeilen funktionierenden Code und bitten Sie ChatGPT um genau eine einzelne Erweiterung. Dann finden Sie mögliche Fehler schneller, als wenn das Sprach-



modell gleich mehrere Ergänzungen auf einmal hinzuprogrammieren soll.

Wenn Sie bestimmte Codeblöcke nicht verstehen, bitten Sie ChatGPT, sie Ihnen Zeile für Zeile zu erklären. Zudem können Sie das aus Performancegründen deaktivierte Debugging mit `use_debug true` wieder einschalten. Dann zeigt Sonic Pi die Zeilen im Code an, in denen es klemmt, und gibt dazu eine Fehlermeldung aus.

Manchmal sitzt der eigentliche Fehler aber woanders. Oft enthalten beispielsweise Variablen andere Zahlen als gedacht. Die möglichen Zustände von Variablen können Sie mit dem Befehl `puts variablenname` ausgeben lassen, den Sie an beliebigen Stellen im Code einfügen. Um die Fehlerquelle weiter einzugrenzen, empfiehlt es sich zudem, einzelne Code-Zeilen mit dem `#`-Zeichen auszukommentieren. Das sollten Sie aber mit Bedacht tun, denn dadurch können neue Fehler entstehen.

Während unserer Testphase reagierte Sonic Pi manchmal nicht mehr, obwohl der Code in Ordnung war. In solchen Fällen half ein simpler Neustart, und die Musik lief weiter.

## Fazit

Sonic Pi ist ein wunderbarer Spielplatz, auf dem man nicht nur neue musikalische Ideen ausprobieren, sondern auch gefahrlos mit den Programmierfähigkeiten von Sprachmodellen nach Lust und Laune experimentieren kann. Hier passiert viel Überraschendes, von der komplett missratenen Melodie bis hin zum coolen Zufallsrhythmus; dank Live-Coding fällt beides schnell auf, sodass man Fehler sofort ausbessern und schnell von Version zu Version iterieren und so den Sound verbessern kann.

Solche Experimente schulen nicht nur das musikalische Verständnis, sondern helfen auch, die Möglichkeiten und Grenzen von Sprachmodellen zu erkennen. Am besten kommt voran, wer dem Chatbot Schritt für Schritt genau vorgibt, was er tun soll. Die Sprach-KIs sind gut darin, die passende Syntax zu finden; die musikalischen Ideen müssen jedoch vom Anwender kommen. Deshalb eignet sich Sonic Pi ideal für den Schulunterricht, um Musik und Informatik zusammenzubringen und Schüler spielerisch fürs Coden und für Musiktheorie zu begeistern. (hag) **ct**

# Für alles gerüstet! Tests, Tipps und Tools



Das Sonderheft richtet sich vor allem an Privatanutzer, Freelancer und kleinere Unternehmen und enthält Kaufberatungen, Tests und Praxisanleitungen zu typischen Büroprogrammen, auch abseits von Microsoft Office.

**JETZT BESTELLEN!**



Heft für 14,90 € • PDF für 12,99 € • Heft + PDF 19,90 €



**shop.heise.de/ct-homeoffice24**

# Interview: Kann KI grün sein?

**Wasser, Energie, seltene Erden: Der Einsatz von Sprachmodellen erfordert einige Ressourcen. Ein Team des Instituts für ökologische Wirtschaftsforschung hat untersucht, wie man KI-Systeme nachhaltiger gestalten kann. Wir haben mit der Hauptautorin des Grundsatzpapiers gesprochen, Dr. Friederike Rohde.**

Von **Greta Friedrich und Isabel Grünewald**

**S**prachmodelle nutzen große Energiemengen zum Training und große Wassermengen, um Server zu kühlen, in denen wiederum Ressourcen wie seltene Erden verbaut sind. Und auch in der Nutzungsphase ist viel Energie nötig, um Anfragen zu verarbeiten und zu beantworten (siehe Artikel „Fünf Sprachmodelle gegen ChatGPT“).

Das Institut für ökologische Wirtschaftsforschung (IÖW) in Berlin hat im Leuchtturmprojekt SustAIIn über drei Jahre ausgearbeitet, wie KI-Anwendungen nachhaltiger werden können. Das Bundesumweltministerium förderte das Projekt. Wir haben mit Dr. Friederike Rohde, Hauptautorin des Grundsatzpapiers zum Projekt, über ihre Erkenntnisse, nachhaltige KI-Systeme und deren Regulierung gesprochen.

**c't: Frau Rohde, was ist für Sie die wichtigste Erkenntnis aus der SustAIIn-Studie?**

**Friederike Rohde:** Dass es sehr viele Aspekte der Nachhaltigkeit gibt. Für eine ernsthafte Nachhaltigkeitsperspektive muss man KI entlang des Lebenszyklus betrachten. Da sind dann einerseits ökologische Auswirkungen, aber auch der indirekte Ressourcenverbrauch und alles, was in der ethischen Debatte diskutiert wird: Verteilungsgerechtigkeit, Transparenz, menschliche Aufsicht, partizipatives Design. All diese Punkte spielen eine wichtige Rolle bei der Entwicklung von KI-Systemen.

Wir konnten mit unserer Arbeit sicherlich einen Beitrag dazu leisten, dass die Politik nicht nur gesellschaftliche, sondern künftig auch ökologische Auswirkungen von KI-Systemen berücksichtigt. Zum

Beispiel bei der europäischen KI-Verordnung. Auch Entwicklerinnen und Entwickler machen sich Gedanken, wie man den Energieverbrauch messen und reduzieren kann, wie man Modelle anders gestaltet, sie komprimiert oder kleinere Modelle nutzt. Es gibt beispielsweise eine Bewegung hin zu Tiny Models, also sehr kleinen Modellen.



Bild: Institut für ökologische Wirtschaftsforschung

**Friederike Rohde, Hauptautorin des Synthesepapiers zum SustAIIn-Projekt: „Viele KI-Systeme haben ganz andere Dinge zum Ziel, als irgendetwas Positives für die Umwelt oder die Gesellschaft zu tun.“**

**c't: Warum geben Anbieter von Sprachmodellen oder KI-Anwendungen bisher nicht an, wie viele und welche Ressourcen ihre Systeme verbrauchen?**

**Rohde:** Warum die Unternehmen das nicht veröffentlichen, ist nochmal ein anderes Thema. Es gibt eine ganze Reihe von Tools, mit denen man den Ressourcenbedarf fürs Training relativ einfach abbilden oder einschätzen kann, indem man die Trainingszeit angibt und dazu sagt, wo die Datencenter stehen, wie lange trainiert wurde, auf welcher Hardware und so weiter. Da eine Größenordnung anzugeben ist also nicht das Problem. Es gibt auch schon einige Studien, die Zahlen geliefert und Sprachmodelle miteinander verglichen haben.

Schwieriger einzuschätzen ist, was dem Training vor- und nachgelagert ist und ja auch Energie und Ressourcen verbraucht. Vorher passiert die ganze Modellentwicklung. Also welches Modell nehme ich überhaupt, wie groß soll es sein? Hinterher kommt noch der Energieverbrauch in der Nutzungsphase hinzu. Dann zieht das trainierte Modell aus neuen Daten Schlussfolgerungen, man nennt das auch Inferenz.

Eine relativ neue Studie aus den USA hat festgestellt, dass gerade bei generativer KI der Energie- und Ressourcenverbrauch in der Nutzungsphase wahnsinnig viel höher ist als bei anderen, spezialisierten Modellen. Die Entwicklung hin zu sogenannten Foundation Models, die sehr viele unterschiedliche Dinge tun können, geht damit einher, dass sehr große Modelle entstehen, die mehr Energie und Ressourcen verbrauchen.

**c't: Welche Ressourcen verbrauchen KI-Modelle denn?**

**Rohde:** Der Energieverbrauch ist ein Thema, ein zweites ist der wahnsinnig hohe Wasserverbrauch der Rechenzentren, ein indirekter Ressourcenverbrauch. Um die Rechenzentren zu kühlen, wird in der Regel Trinkwasser genutzt und es gibt Verteilungskonflikte um diese Ressource. Es gibt Beispiele in Ecuador und Chile, wo die Bevölkerung oder auch lokale Behörden versucht haben zu verhindern, dass Datenzentren gebaut werden. Wasser ist einfach eine viel knappere und endlichere Ressource als Energie.

Entlang des ganzen Lebenszyklus kommen noch viele andere Dinge dazu. Auch die Hardware muss ja produziert werden, mit Rohstoffen wie seltenen Erden. Da ist nicht nur die Frage, welche Ressourcen

verwendet werden, sondern vor allem auch, wo sie gewonnen werden. Häufig passiert das unter sehr problematischen Arbeitsbedingungen, das ist dann die soziale Komponente. Da kann man sich schon die Frage stellen: Ist KI eine Lösung, die im ersten Moment nachhaltiger aussieht als sie ist?

**c't: Wie sind Sie in Ihrer Studie vorgegangen, um den Energie- und Ressourcenbedarf von LLMs oder KI einzuschätzen?**

**Rohde:** Wir haben den Stand der Wissenschaft zusammengetragen, die Hauptleistung war aber, ein umfassendes Set an Nachhaltigkeitskriterien zu entwickeln. Damit haben wir Fallstudien durchgeführt und uns den Energie- und Ressourcenverbrauch von einzelnen Systemen angeguckt – aber auch geschaut, was bewirkt die Technologie in diesem oder jenem konkreten Fall? Welche Potenziale im Hinblick auf Nachhaltigkeit sind damit verbunden?

Ein Fallbeispiel war die energetische Optimierung für ein Wohnquartier durch KI. Tatsächlich war der Energieverbrauch des betrachteten Machine-Learning-Modells nicht so hoch. Da konnte man sagen, okay, in diesem Fall wiegen die ökologischen Vorteile schwerer als der Energieverbrauch fürs Training. Weil hier fossile Energieträger eingespart und die erneuerbare Energieeinspeisung optimiert wurden.

In der Debatte heißt es oft: „Wir können KI für ganz viele Dinge einsetzen, wo wir Energie und Ressourcen einsparen, und am Ende sparen wir dann mehr ein, als wir nutzen.“ Aber die Anwendungen, die wirklich in diesen Bereich fallen, sind nur ein Bruchteil aller KI-Anwendungen. Wir merken es ja alle selbst. Wir benutzen ChatGPT oder irgendwelche anderen Tools, die gar nichts damit zu tun haben, irgendwo irgendetwas einzusparen. Es gibt einfach sehr viele KI-Systeme, die ganz andere Dinge zum Ziel haben, als irgendetwas Positives für die Umwelt oder für die Gesellschaft zu erreichen.

**c't: Sie erwähnten eben ein KI-System, das den Energieverbrauch eines Wohnviertels reguliert. In welchen Bereichen kann man KI so nachhaltig einsetzen, dass sie ihren Ressourcenverbrauch ausgleicht oder rechtfertigt?**

**Rohde:** Ob sie den Ressourcenverbrauch ausgleichen kann, lässt sich so pauschal nicht beantworten. Man kann zumindest sagen: Ja, es gibt Anwendungsbereiche, in denen es auf jeden Fall Sinn ergibt. Oftmals sind es Optimierungsanwendungen. Wenn



beispielsweise die Landwirtschaft Machine-Learning-Modelle verwendet, um den Pestizideinsatz zu verringern; es gibt Studien, die zeigen: Ja, der verringert sich tatsächlich um 20 oder 30 Prozent. Oder KI kann zum Beispiel die Materialeffizienz in Unternehmen optimieren. Oder denken Sie an den ganzen Bereich Umwelt- und Klimaschutz, ob es jetzt große Klimamodelle sind oder Waldmonitoring oder überhaupt Biodiversitätsmonitoring.

Es gibt also viele Einsatzbereiche für KI, bei denen man davon ausgehen kann, dass damit ein Nutzen für Gesellschaft und Umwelt verbunden ist. Ob am Ende die Energiebilanz positiv ist, ist schwierig zu sagen, und ich weiß auch nicht, ob das der einzige Maßstab sein sollte. Ich glaube, man sollte grundsätzlich sagen: Es ist ein sehr machtvolleres Werkzeug und wir können es für Zwecke einsetzen, die wir als Gesellschaft als sinnvoll erachten und die dem Gemeinwohl dienen.

Häufig haben wir aber das Problem, dass die positiven Effekte im globalen Norden auftreten und die negativen Effekte in anderen Ländern der Welt. Also da, wo beispielsweise die Rechenzentren stehen, oder da, wo die Materialien, die ganzen Sensoren und Chips und so weiter, hergestellt werden. Es ist dann die Frage, für wen ist diese KI-Anwendung jetzt eigentlich Netto-Positiv? Nur aus unserer Perspektive oder auch global betrachtet?

**c't: Wenn man den ganzen Lebenszyklus eines Modells betrachtet, was verbraucht mehr Ressourcen, das Training oder die Nutzung?**

**Rohde:** Das ist unterschiedlich, aber tendenziell verbraucht die Nutzung mehr Ressourcen. Das Verhältnis von Training und Inferenz hängt natürlich immer vom konkreten Anwendungsfall ab. Bei einem System, das nur einmal am Tag die Inferenz macht, wie zum Beispiel die Energieoptimierung für Häuser, ist natürlich auch der Energieverbrauch dafür nicht so hoch. Aber bei einem System, das Millionen Menschen verwenden, Empfehlungsalgorithmen zum Beispiel, da passiert die Inferenz sehr, sehr oft und der Energieverbrauch dafür ist sehr viel höher als der fürs Training. Ich würde sagen, dass dieses Szenario sehr häufig ist. Eben weil in der Regel das Training nur einmal passiert und die Nutzung ganz häufig.

**c't: Ist es dann überhaupt vertretbar, dass man KI privat nutzt, zum Zeitvertreib oder für Dinge, die man auch ohne sie schaffen würde?**

**Rohde:** Das ist tatsächlich eine große Frage: Was ist ein gesellschaftlich sinnvoller Einsatz solcher Machine-Learning-Modelle? Macht es wirklich Sinn, dass das frei verfügbar ist, dass alle das benutzen können? Es gibt sehr viele Anwendungsfälle, bei denen man aus unserer Sicht durchaus kritisch hinterfragen kann, ob man dafür wirklich die knappen Ressourcen nutzen sollte. Oder ob man sich nicht eher darauf konzentrieren sollte, in welchen Bereichen wir die Technologie so einsetzen können, dass sie einen gesellschaftlichen Nutzen hat und nicht nur einen ökonomischen.

Wir meinen, es braucht eine größere gesellschaftliche Reflexion darüber. Es ist eine total schwierige Frage, was eine gesellschaftlich sinnvolle Nutzung ist. Man kann natürlich sagen, okay, wo eine Gemeinwohlorientierung erkennbar ist, da macht es Sinn und ist legitim. Aber man kann jetzt nicht Akteuren, Organisationen, Unternehmen oder privaten Nutzerinnen und Nutzern verbieten, KI-Modelle einzusetzen. Wir können das Rad nicht zurückdrehen.

**c't: Kann der steigende Energiebedarf von immer mehr Rechenzentren denn überhaupt mit erneuerbaren Energien gedeckt werden?**

**Rohde:** Es heißt immer, wir haben ja dann erneuerbare Energie, dann ist das mit dem Energieverbrauch alles überhaupt kein Thema mehr. Aber das stimmt nicht. Zum einen elektrifizieren wir ganz viele Sektoren, zum Beispiel den Mobilitätssektor. Da wird der Bedarf an Strom steigen. Und wir haben auch eine Flächenproblematik. Es wird teils um jedes Windrad gerungen, um jede Stromtrasse, es gibt gesellschaftliche Konflikte.

Eigentlich basiert das Ziel, irgendwann mal bei 100 Prozent erneuerbaren Energien zu sein, darauf, dass wir den Energieverbrauch insgesamt senken. Und wenn das nicht passiert, dann wird es ganz schwierig, jemals bei diesen 100 Prozent anzukommen.

**c't: In Ihren Empfehlungen fordern Sie dazu auf, die Entwicklung und Nutzung von KI-Systemen zu regulieren. Wie könnte das aussehen?**

**Rohde:** Wir haben ja die europäische KI-Verordnung, die jetzt in Kraft tritt. Darin sind schon sehr viele Dinge verankert, in Bezug auf Dokumentation, Risikomanagement und so weiter, immer risikobasiert abgestuft. Es gibt zum Beispiel Dokumentationspflichten für die Energie- und Ressourcennutzung

## In voller Länge ...



... können Sie dieses Interview im heise-Podcast KI-Update (siehe [ct.de/w345](https://ct.de/w345)) hören.

von High-Risk-Systemen und auch von der sogenannten General Purpose AI (siehe Artikel „EU-Regeln für KI“). Aber es geht auch um andere Themen, die mit Nachhaltigkeit zu tun haben, zum Beispiel: Wer sind vulnerable Gruppen, welche Datensets werden genutzt, und so weiter und so fort.

Die Regulationslage ist eigentlich gut und ausreichend. Aber wer kontrolliert eigentlich die Umsetzung? Das ist der Elephant in the Room, genauso wie bei der Datenschutz-Grundverordnung. Wichtig ist auch, wie die Gremien aufgestellt sind, die die Umsetzung kontrollieren sollen. Welche Akteure bestimmen jetzt bei der Normierung der Dokumentation mit? Welche Rolle spielt die Zivilgesellschaft?

Weitere Informationen zum SustAIIn-Projekt

[ct.de/w345](https://ct.de/w345)

**ct:** Welche Maßnahmen wünschen Sie sich denn von der Politik?

**Rohde:** Zum einen eine konsequente Umsetzung der KI-Verordnung und Aufsichtsgremien, die keine

zahnlosen Tiger sind, sondern die Umsetzung wirklich begleiten können. Auch Normierungsgremien, in denen die Zivilgesellschaft in angemessenem Maße repräsentiert ist, würden wir uns wünschen. Zum anderen wünsche ich mir in der Technologie- und Innovationspolitik einen Fokus auf Technologien, die wirklich das Gemeinwohl fördern und nicht einen rein wirtschaftlichen Nutzen haben.

Insgesamt wünsche ich mir für die Zukunft eine kritischere Reflexion darüber, was es eigentlich bedeutet, KI-Systeme zu entwickeln und einzusetzen. Wichtig ist eine gesellschaftliche Debatte, in der es auch möglich ist, vorherrschende Technikvisionen, Erzählungen oder Narrative zu hinterfragen.

**ct:** Und was müsste in der Wirtschaft passieren?

**Rohde:** Bei der Wirtschaft geht es aus meiner Sicht sehr viel um Capacity Building. Also darum, die Akteure – sowohl Entwickler als auch Unternehmen, die KI nutzen, – zu befähigen, solche Technologien verantwortungsvoll zu entwickeln. Ich glaube, da braucht es noch sehr viel Aufklärungs- und Bildungsarbeit.

Es gibt in vielen Bereichen schon ein Bewusstsein dafür und ein großes Interesse, sich stärker mit dem Thema nachhaltige KI zu beschäftigen. Aber oftmals fehlen dann die Kapazitäten oder die Ressourcen oder die Programmierer arbeiten unter großem Zeitdruck. Man muss Ansatzpunkte suchen, das organisational so zu verankern, dass es für die Unternehmen gut funktioniert. Dass sie es nicht als Hürde empfinden, Technik verantwortungsvoll zu gestalten, sondern als Bereicherung. (gref **ct**)

# Endlich *gute* Fotos!



**JETZT  
LOSLEGEN!**



**+ TIPPS  
VON PROFI-  
FOTOGRAFEN**



**shop.heise.de/ct-fotoeinsteiger25**



Bild: KI, Collage c't

# Warum Big-Tech auf Atomenergie setzt

**Jahrelang feierten sich die großen Tech-Konzerne als Ökostrom-Champions, plötzlich kaufen sie Atomstrom und finanzieren neue Kernkraftwerke. Was steckt dahinter?**

Von **Christian Wölbart**

**A**m 28. März 1979 um 4 Uhr morgens fiel im Block 2 des Kernkraftwerks Three Mile Island in Pennsylvania die Kühlung aus. Aufgrund einer Reihe weiterer technischer Probleme und Bedienungsfehler überhitzte der Reaktorkern, die Brennstäbe begannen zu schmelzen und der Kern wurde weitgehend zerstört. Weil die Mannschaft im Kontrollraum die Kühlung wieder in Gang setzen konnte, blieb die Bevölkerung verschont – anders als 1986 in Tschernobyl und 2011 in Fukushima. Die USA waren an einer Katastrophe vorbeigeschrammt.

Nun, 45 Jahre später, macht Three Mile Island wieder Schlagzeilen: Ende September verkündete

Microsoft den Abschluss eines zwanzigjährigen Stromlieferungsvertrags mit dem aktuellen Betreiber des Kraftwerks, Constellation. Der Abmachung zufolge wird Constellation den unbeschädigten, 2019 aus wirtschaftlichen Gründen abgeschalteten Block 1 des Kraftwerks im Jahr 2028 wieder in Betrieb nehmen, um Microsoft-Rechenzentren in der Region mit Strom zu versorgen.

Der Neustart von Three Mile Island ist symbolträchtig, aber nicht der einzige Atomdeal von Big Tech: Amazon kaufte bereits im Frühjahr 2024 ein Rechenzentrum, das direkt neben einem anderen Kernkraftwerk in Pennsylvania steht, und will sich



Strom aus dem Kraftwerk sichern. Im September verkündete Amazon zudem mehrere Investitionen in Projekte zum Bau kleiner, modular aufgebauter Atomreaktoren in verschiedenen US-Bundesstaaten. Im Oktober folgte der dritte Tech-Gigant: Google meldete, ab 2030 ebenfalls Atomstrom aus solchen Mini-AKW beziehen zu wollen.

Die Atom-Deals von Big Tech irritieren, weil die Konzerne traditionell stark auf erneuerbare Energien setzen und auch lautstark damit werben. Google zum Beispiel trommelt seit 2017 damit, seine Rechenzentren und Büros zu 100 Prozent mit „Renewables“ zu betreiben. Amazon hat dieses Ziel nach eigenen Angaben 2023 erreicht und kauft laut einer Statistik von Bloomberg mehr erneuerbare Energie ein als jedes andere Unternehmen auf der Welt.

## Geografie und Zeit

Warum nun also plötzlich Atomstrom, wenn Big Tech augenscheinlich genügend Ökostrom bezieht? Weil die Konzerne „100 Prozent erneuerbare Energie“ großzügig definieren: Sie meinen damit lediglich, dass sie ihren jährlichen Energiebedarf auf dem Papier mit diversen Ökostromprodukten ausgleichen. Darunter sind langfristige Abnahmegarantien, die den Bau von Wind- und Solarparks erst ermöglichen, aber in vielen Fällen auch umstrittene Zerti-

fikate von Altanlagen, die kaum dazu beitragen, dass mehr Grünstrom ins Netz fließt [1].

Die Tech-Konzerne nutzen den Trick mit dem jährlichen Ausgleich vor allem aus zwei Gründen. Der erste ist geografischer Natur: Ihre Rechenzentren bauen sie meist in der Nähe der Kundschaft oder wichtiger Netzknoten und nicht unbedingt dort, wo es Ökostrom im Überfluss gibt.

Der zweite Grund hat etwas mit Zeit zu tun: Selbst dort, wo im Mittel relativ viel Ökostrom im Netz ist, ist dieser nicht rund um die Uhr verfügbar, weil nun einmal die Sonne nicht immer scheint und der Wind nicht immer weht. Um diese Flauten abzudecken, müssen die Tech-Konzerne also auch konventionellen Strom einkaufen.

Eine Rund-um-die-Uhr-Versorgung von Rechenzentren mit erneuerbarer Energie wäre an den meisten Standorten kaum machbar oder zumindest extrem teuer: Die Tech-Konzerne müssten noch viel mehr erneuerbare Energie einkaufen, den Überschuss in gigantischen Batterien speichern und auch noch die Übertragungsnetze so ausbauen, dass all das funktioniert.

## Schwenk von „erneuerbar“ zu „CO<sub>2</sub>-frei“

Statt dieses vielleicht unmögliche Unterfangen anzugehen, verkündete Google schon im Jahr 2020

Bild: Kairos Power



**Mini-AKW des Start-ups Kairos Power sollen ab dem Jahr 2030 Google-Rechenzentren mit Strom beliefern.**

ein anderes Ziel: Man werden alle Rechenzentren und Büros ab dem Jahr 2030 „24/7“ ausschließlich mit „CO<sub>2</sub>-freier Energie“ versorgen, kündigte der Suchmaschinenkonzern an. Das heißt, dass zu jeder Betriebsstunde die nötige Menge CO<sub>2</sub>-freier Energie in dasselbe nationale oder regionale Netz eingespeist werden soll. Zu CO<sub>2</sub>-freier Energie zählt Google neben den Erneuerbaren explizit auch Atomstrom. Microsoft fährt inzwischen dieselbe Strategie.

Atomenergie soll also das Problem lösen, dass der Wind nicht immer bläst und die Sonne nicht immer scheint, wenn die Server laufen. Wenn die Konzerne genügend Atomstrom beschaffen, um ihre Rechenzentren bei Licht- oder Windmangel zu versorgen, können sie auf klimaschädlichen Strom aus Kohle- und Gaskraftwerken verzichten und ihre Clouddienste als klimafreundlich anpreisen. Anders formuliert: Um beim Klimaschutz voranzukommen, nehmen sie die Risiken der Atomkraft in Kauf.

Wie weit Google auf dem Weg zur CO<sub>2</sub>-freien Rund-um-die-Uhr-Versorgung ist, zeigt eine lange Tabelle im Umweltbericht des Konzerns. Demnach machte Google zuletzt Rückschritte: 2020 lag der Anteil der CO<sub>2</sub>-freien Energie noch bei 67 Prozent, 2022 und 2023 waren es 64 Prozent. Das hat vor allem mit der Inbetriebnahme von Rechenzentren in Regionen zu tun, in denen es kaum erneuerbare Energien und keinen Atomstrom gibt. Laut dem Bericht lag der Anteil CO<sub>2</sub>-freier Energie in den Google-Rechenzentren in Katar und Saudi-Arabien 2023 bei 0 Prozent, in Kanada waren es – dank Wasserkraft – 100 Prozent, in Finnland 98.

Die Werte stehen wohlgerne nicht nur für die Leistung, die in die Rechenzentren fließt, sondern für die Summe aus lokal bezogener und im gleichen regionalen oder nationalen Netz zugekaufter CO<sub>2</sub>-freier Energie. Schließlich stehen die Windparks und Solarfelder selten direkt neben den Rechenzentren.

## Auftritt ChatGPT

Ein weiteres Problem auf dem Weg zur CO<sub>2</sub>-freien Zukunft könnte der KI-Boom werden, der 2022 mit dem Start von ChatGPT so richtig losging. Solche großen KI-Modelle benötigen zu ihrem Training und im Betrieb gigantische Mengen an Rechenleistung, weitaus mehr als bisher übliche Cloud- und Serveranwendungen. Experten fürchten, dass der Energiebedarf der Tech-Branche künftig noch viel schneller wächst als bisher gedacht. Dazu kursieren zahlreiche Hochrechnungen, die mit Vorsicht zu genießen sind, weil viele schwer abschätzbare Faktoren im Spiel

sind. Dazu zählen etwa die Fragen, wie lang der KI-Hype noch anhält und wie stark die Effizienz von Hard- und Software steigt.

Aber auch konservative Schätzungen zeigen ein dynamisches Wachstum. Die US-Firma Uptime Institute, die Rechenzentren zertifiziert, schätzt den jährlichen Energiebedarf der Anfang 2024 installierten KI-Beschleuniger von Nvidia auf 5,8 Terawattstunden beziehungsweise gut zwei Prozent des gesamten Strombedarfs aller Rechenzentren weltweit. Sollte Nvidia seine Lieferziele einhalten, läge der jährliche KI-Energiebedarf der Anfang 2025 installierten KI-Chips schon bei 21,9 Terawattstunden, also beinahe dem Vierfachen.

Da überrascht es nicht, dass Nvidia-Chef Jensen Huang Ende September Atomenergie als „wunderbare Zukunftsoption“ bezeichnete. „Wir werden Energie aus allen Quellen brauchen und müssen die Verfügbarkeit, die Kosten und die Nachhaltigkeit über die Zeit ausbalancieren“, sagte er gegenüber Bloomberg.

## Mini-AKW ...

Doch selbst, wenn der Energiehunger der KI nicht so stark wächst wie befürchtet: Für die Tech-Konzerne dürfte es so oder so äußerst schwierig werden, an allen Standorten genügend klimafreundlich produzierte Energie zu beschaffen. Wohl auch deshalb begnügen die Konzerne sich nicht damit, einfach nur Strom aus bestehenden Atommeilern einzukaufen. Google und Amazon finanzieren auch die Entwicklung und den Bau neuartiger Kernkraftwerke. Amazon investiert dazu in das US-Start-up X-Energy, Google schloss einen Vertrag mit dem Konkurrenten Kairos Power.

X-Energy und Kairos entwickeln gleich in zweifacher Hinsicht besondere Atomkraftwerke: Die erste Besonderheit ist, dass sie kleine, modular aufgebaute Reaktoren planen (Small Modular Reactors, SMR). Der Kairos-Entwurf zum Beispiel besteht aus einem Tandem aus zwei Reaktoren mit jeweils 75 Megawatt Leistung. Zum Vergleich: Block 1 von Three Mile Island lieferte bis zu seiner Abschaltung rund 830 Megawatt, richtig große Reaktoren leisten über 1600 Megawatt.

Zu den Vorteilen der geringen Größe gehört, zumindest in der Theorie, eine einfachere Standortfindung. Denkbar ist zum Beispiel, dass die Reaktoren künftig direkt neben Rechenzentren platziert werden, was Kosten für den Netzausbau sparen würde. Die modulare Bauweise wiederum soll eine effiziente



**So stellt sich  
X-Energy sein aus  
vier Mini-Reaktoren  
bestehendes Atom-  
kraftwerk vor.**

Bild: X-Energy



Serienfertigung eines Standarddesigns ermöglichen. Die Reaktoren könnten dann in Fabriken vormontiert und zum Einsatzort transportiert werden.

Ob SMR unter dem Strich wirklich günstiger sind als die altbekannten Riesenreaktoren, lässt sich jedoch noch nicht sagen. „Es muss sich zeigen, ob der Ansatz ‚klein und modular‘ mit Anlieferung werksseitig vormontierter Reaktoren mit Typenzulassungen anstatt Einzelgenehmigungen sich wirklich gegen die Ökonomie der Größe durchsetzen kann“, sagt Horst-Michael Prasser, Kernkraftexperte und emeritierter Professor der ETH Zürich, gegenüber c’t.

### ... mit neuer Kühltechnik

X-Energy und Kairos haben allerdings noch ehrgeizigere Pläne als die meisten anderen SMR-Entwickler: Sie setzen nicht auf die bei großen AKW gängige Wasserkühlung, sondern auf neue Kühltechnik. Besonders wenig Praxiserfahrung gibt es mit dem Kairos-Konzept eines mit einer Salzsammelze gekühlten Reaktors mit kugelförmigen Brennelementen. „Es gibt keine Vorläuferanlagen“, betont Prasser.

Kairos hat gerade in Tennessee mit dem Bau seines ersten Testreaktors mit 35 Megawatt thermi-

scher Leistung begonnen, der bis 2027 in Betrieb gehen soll. Schon 2030 will Kairos dann den ersten Reaktor in Betrieb nehmen, der Google mit Strom versorgt.

Etwas mehr Erfahrung gibt es mit heliumgeköhlten Kugelhaufenreaktoren, wie X-Energy sie entwickelt. Zwei AKW mit im Grundsatz ähnlicher Technik waren bis zu den Achtzigerjahren in Deutschland in Betrieb, auch China hat Erfahrung damit. X-Energy betreibt aber noch keine Testanlage. Amazon hofft, dass X-Energy-Reaktoren „Anfang der 2030er Jahre“ im US-Bundesstaat Washington in Betrieb gehen. Vier der Blöcke sollen dort zusammen 320 Megawatt leisten.

Die Zeitpläne sind allerdings sportlich: In den vergangenen Jahrzehnten dauerten die meisten AKW-Projekte selbst beim Einsatz etablierter Technik viele Jahre länger als geplant und wurden deutlich teurer, wenn sie nicht komplett scheiterten. Die Finanzierung von Neuentwicklungen wie bei Kairos und X-Energy habe den „Charakter von Risikokapitalinvestitionen mit riesigen Renditen, wenn es klappt“, urteilt Prasser. Sollten die Projekte scheitern, wären die Verluste für Konzerne wie Amazon oder Google aber eher „unbedeutend“.

(cwo) **ct**

### Literatur

[1] Christian Wölbert, Redmonds CO<sub>2</sub>-Problem, Warum Microsoft beim Klimaschutz Rückschritte macht, c’t 14/2024, S. 112



# KI-Hardware für PCs und Notebooks

AMD, Apple, Intel, Nvidia und Qualcomm werben mit speziellen KI-Funktionen ihrer Chips. Apple und Microsoft wiederum preisen die KI-Funktionen ihrer Betriebssysteme an, sprich Apple Intelligence und Copilot+. Bei Künstlicher Intelligenz geht vieles durcheinander, wir dröseln es auf.



Von **Christof Windeck**

## Was bringen KI-Rechenwerke?

**?** Was genau können KI-Beschleuniger besser als normale CPU-Kerne?

**!** Viele KI-Algorithmen benötigen sehr hohe Rechenleistung, kommen aber größtenteils mit zwei sehr speziellen Funktionen aus: der Multiplikation großer Matrizen und einer anschließenden Addition, genannt Matrix Multiply Accumulate (MMA). Eine gezielt für diese Rechenschritte optimierte Baugruppe verarbeitet die Daten viel schneller und gleichzeitig energetisch effizienter als ein Allzweck-Prozessorkern. Das klappt in der Praxis jedoch nur, wenn mehrere Bedingungen erfüllt sind.

## Kompatibilität von KI-Apps und -Beschleunigern

**?** Kann jede KI-Software jeden KI-Beschleuniger nutzen?

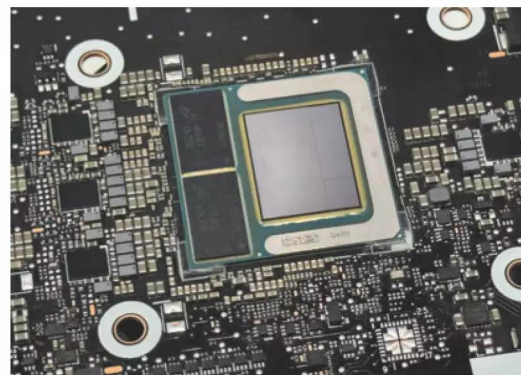
**!** Nein, und hier liegt der Hase im Pfeffer: Um ein bestimmtes KI-Rechenwerk nutzen zu können, muss eine KI-App gezielt dafür programmiert sein. Denn dummerweise sind die unterschiedlichen KI-Rechenwerke in Chips von AMD, Apple, Intel, Nvidia und Qualcomm nicht zueinander binärkompatibel. Schlimmer noch: Manche Chips enthalten zwei oder drei verschiedene KI-Rechenwerke. Sie sind zwar über Treiber ins jeweilige Betriebssystem eingebunden und es gibt dafür auch standardisierte Programmierschnittstellen (APIs) sowie dazu kompatible

KI-Frameworks. Doch manche KI-App kooperiert nur mit bestimmten KI-Einheiten oder bestimmten APIs, andere sind für sie nutzlos.

## KI-Programmierschnittstellen

**?** Wie bekomme ich heraus, welche KI-Software meine Hardware optimal ausreizt?

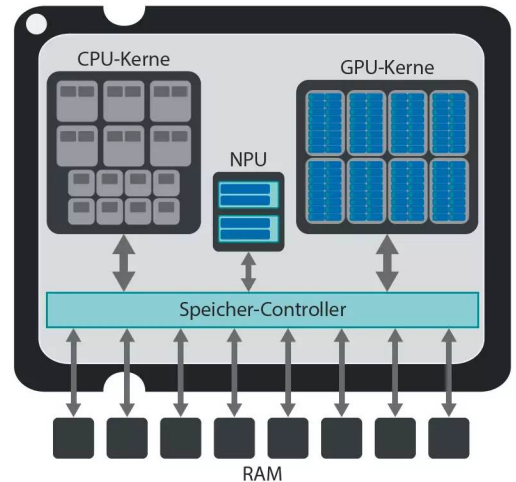
**!** Das ist schwierig. Viele Softwarefirmen verraten nicht einmal, welche Softwareschnittstelle (Application Programming Interface, API) ihre KI-App



**Aktuelle Mobilprozessoren wie der abgebildete Intel Core Ultra 200V (Lunar Lake) enthalten außer CPU-Kernen und einem integrierten Grafikprozessor (IGP) auch eine Neural Processing Unit (NPU) für KI-Apps.**

## KI-taugliche Rechenwerke in x86-Prozessoren

Sowohl die CPU-Kerne, als auch die GPU-Kerne und die NPU verarbeiten KI-Datenformate wie Int8 und BF16. In der CPU sitzen AVX2-VNNI- oder AVX-512-VNNI-Einheiten. Viele GPU-Typen haben außer Shaderkernen auch spezielle KI-Rechenwerke. Die NPU ist nur für Matrix Multiply Accumulate (MMA) ausgelegt.



überhaupt nutzt. Und selbst mit dieser Information ist es schwer einzuschätzen, wie gut die App auf einem bestimmten Rechner laufen wird. Denn je nach Kombination aus KI-Framework, KI-API, Treibern und Hardware unterscheidet sich die Performance gewaltig.

KI-Apps verwenden typischerweise sogenannte KI-Frameworks wie TensorFlow, Caffe, PyTorch oder Keras. Diese wiederum nutzen je nach den vorhandenen KI-Rechenwerken, Treibern und Betriebssystemen verschiedene Programmierschnittstellen wie Microsoft DirectML (Windows ML), Apple CoreML, Nvidia TensorRT, AMD AI Engine, Intel OpenVINO, Qualcomm AI Engine Direct oder auch die generische Schnittstelle Vulkan.

Microsoft wirbt dafür, dass Programmierer unter Windows DirectML nutzen, weil sich darüber KI-Einheiten unterschiedlicher Chipfirmen ansteuern lassen. Doch Benchmarks zeigen, dass DirectML oft deutlich weniger Rechenleistung herausholt als mit dem vom jeweiligen Hardwarehersteller selbst gepflegten API. Das gilt insbesondere für Nvidia TensorRT und Intels OpenVINO.

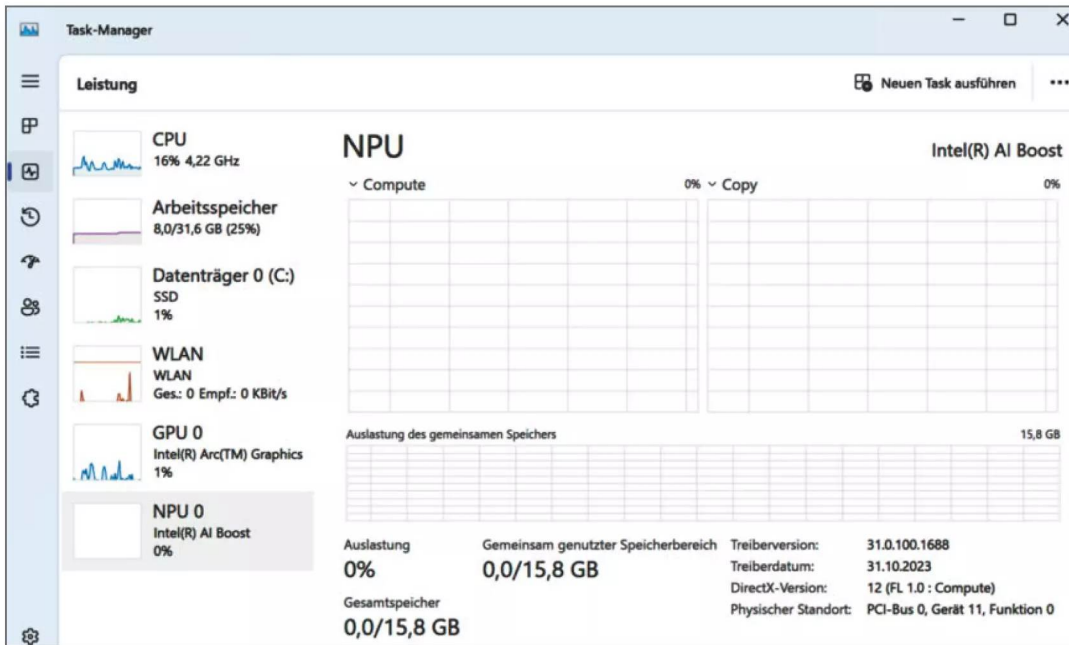
### Was sind Tops?

**?** Immer wieder stoße ich auf die Maßeinheit „Tops“. Ich habe schon begriffen, dass damit

**Tera-Operationen pro Sekunde gemeint sind. Aber was bedeutet das?**

**!** Für den KI-Siegeszug war die Erkenntnis wichtig, dass viele KI-Algorithmen auch dann gute Resultate liefern, wenn sie mit stark vereinfachten Werten arbeiten. Daher rechnen viele KI-Apps mit sogenannten quantisierten Daten. Beispielsweise verwenden sie anstelle von Gleitkommawerten mit je 32 Bit (32-Bit Floating Point, FP32) nur FP16 oder gleich ganze Zahlen (Integer) mit 8, 6 oder gar nur 4 Bit. Ein solcher Int8-Wert belegt 1 Byte, also nur ein Viertel dessen, was ein FP32-Wert braucht. Und moderne CPU-Rechenwerke wie die Advanced Vector Extensions (AVX) verarbeiten pro Taktzyklus von den „schmalen“ Daten viel mehr. AVX-VNNI hantiert beispielsweise mit 256-Bit-Vektoren und verrechnet anstelle von acht FP32-Werten alternativ auch 32 Int8-Zahlen in einem Rutsch. Von den kleineren Daten passen auch mehr ins RAM und in die Caches.

Für die maximale Anzahl an Datenwerten, die ein Rechenwerk pro Sekunde verarbeiten kann, hat sich die Maßeinheit „Operations per Second“ (ops, Rechenschritte pro Sekunde) eingebürgert. Bei KI-Rechenwerken ist damit üblicherweise die Anzahl der bei Matrixmultiplikationen pro Sekunde verarbeiteten Int8-Werte gemeint, die längst in die Billionen geht: Tera-ops, kurz Tops. Viele KI-Rechenwerke –



Windows 11 zeigt im Geräte-Manager die aktuelle Auslastung der NPU an.

aber längst nicht alle – verarbeiten auch FP16-Werte, jedoch nur halb so schnell wie Int8; typischerweise nennen Hersteller dann den höheren Wert. Geht es um Gleitkommazahlen (Floating Point), schreibt man eher Flops: Floating Point Operations per Second. Ohne konkrete Angabe der Datenformate, die das jeweilige KI-Rechenwerk verarbeiten kann, ermöglichen Tops-Werte nur sehr grobe Leistungsvergleiche.

## Welche KI-Beschleuniger gibt es?

**?** Welche Typen von KI-Beschleunigern stecken in aktuellen Windows- und macOS-Rechnern?

**!** Bei den meisten aktuellen x86- und ARM-Prozessoren sind bereits die gewöhnlichen CPU-Kerne für KI-Algorithmen optimiert und daher deutlich schneller als ihre Vorgänger. Denn AMD, Intel und auch ARM haben ihre jeweiligen Vektor-Rechen-einheiten (Advanced Vector Extensions, AVX, und Scalable Vector Extensions, SVE) überarbeitet, so dass sie nun auch KI-Datenformate wie BF16, FP16

oder Int8 verarbeiten. Im Idealfall sind sie bei gleicher Taktfrequenz um den Faktor zwei bis acht schneller als ältere Prozessoren. Die 16 CPU-Kerne eines AMD Ryzen 9 9950X leisten beispielsweise zusammen rund 10 Tops bei Int8.

Die meisten aktuellen Prozessoren enthalten zudem integrierte Grafikprozessoren (iGPU, IGP). Die sind zwar deutlich schwächer als die GPU einer teuren Grafikkarte, enthalten aber ähnliche Rechenwerke, die außer 3D-Berechnungen und Raytracing mittlerweile auch KI-Datenformate kennen. Die IGP des Mobilprozessors Intel Core Ultra 9 288V liefert beispielsweise 67 Tops. Eine 300 Euro teure Nvidia GeForce RTX 4060 liefert hingegen satte 242 Tops.

Zusätzlich enthalten alle aktuellen Mobilprozessoren von AMD, Apple, Intel und Qualcomm separate KI-Rechenwerke, sogenannte Neural Processing Units (NPUs). Die meisten verarbeiten nur Int8- und FP16-Werte und sind dabei oft auch schwächer als die ebenfalls eingebaute GPU: Intels NPU im Core Ultra 100 liegt bei 13 Tops, im Core Ultra 200V immerhin bei 45 Tops. Für das Logo „Copilot+“ verlangt Microsoft eine NPU mit mindestens 40 Tops.



Der Kniff an den NPUs: Sie sind besonders effizient und fressen wenig Strom. Sie sind vor allem für kontinuierlich laufende KI-Anwendungen gedacht, die den Notebook-Akku nicht schnell leeren sollen: Spracherkennung, Optimierung von Audio- und Videostreams.

## Wie viel RAM braucht ein KI-PC?

**?** Braucht ein KI-PC besonders viel Arbeitsspeicher?

**!** Das lässt sich nicht pauschal sagen. Lokal ausgeführte KI-Modelle können deutlich mehr RAM belegen als beispielsweise Office-Apps. Deshalb verlangt Microsoft für Windows-11-Rechner mit dem Logo Copilot+ mindestens 16 GByte RAM. Apple hat vor der Einführung von Apple Intelligence die Mindestbestückung seiner Macs ebenfalls auf 16 GByte erhöht.

Die bisherigen Copilot+Notebooks sowie auch Apple-Rechner haben Prozessoren mit eingebauten GPUs und NPUs, bei denen sich alle drei Rechenwerke das vorhandene RAM teilen. Soll das KI-Modell auf einer separaten Grafikkarte laufen, muss wiederum deren lokaler Speicher groß genug sein.

## KI-Beschleuniger nachrüsten

**?** Kann ich einen KI-Beschleuniger bei meinem Notebook oder PC nachrüsten?

**!** Das klappt sehr gut bei Desktop-PCs mit einem freien PCI-Express-x16-Steckplatz (PCIe x16): Dort können Sie eine moderne Grafikkarte einbauen. Wie stark sie sein darf, hängt außer von Ihrem Geldbeutel auch vom vorhandenen Netzteil ab, denn viele Grafikkarten brauchen zusätzliche Stromkabel, speziell solche mit mehr als 75 Watt Leistungsaufnahme.

Empfehlenswert sind bislang vor allem Nvidia-RTX-Grafikkarten, weil sie nicht nur relativ viel KI-Rechenleistung bereitstellen, sondern weil Nvidia auch Treiber und Programmierschnittstellen gut pflegt. Je nach KI-App kommen aber auch Karten von AMD oder Intel infrage.

Nur sehr wenige Notebooks und Mini-PCs haben Steckfassungen für Grafikkarten. Es gibt aber für unter 100 Euro KI-Beschleuniger in M.2-Bauform wie den Hailo 8L mit immerhin 13 Tops. Welche KI-Apps diesen unter Linux oder Windows jedoch nutzen, konnten wir bisher nicht in Erfahrung bringen. (ciw) **ct**

# NIX VON DER STANGE!

## Wunsch-PC selber bauen oder aufrüsten



Heft für 14,90 € • PDF für 12,99 €

• Heft + PDF 19,90 €



[shop.heise.de/ct-hardwaretipps24](https://shop.heise.de/ct-hardwaretipps24)

Generell portofreie Lieferung für Heise Medien- oder Maker Media Zeitschriften-Abonnenten oder ab einem Einkaufswert von 20 € (innerhalb Deutschlands). Nur solange der Vorrat reicht. Preisänderungen vorbehalten. E-Books können einem DRM-Schutz unterliegen.

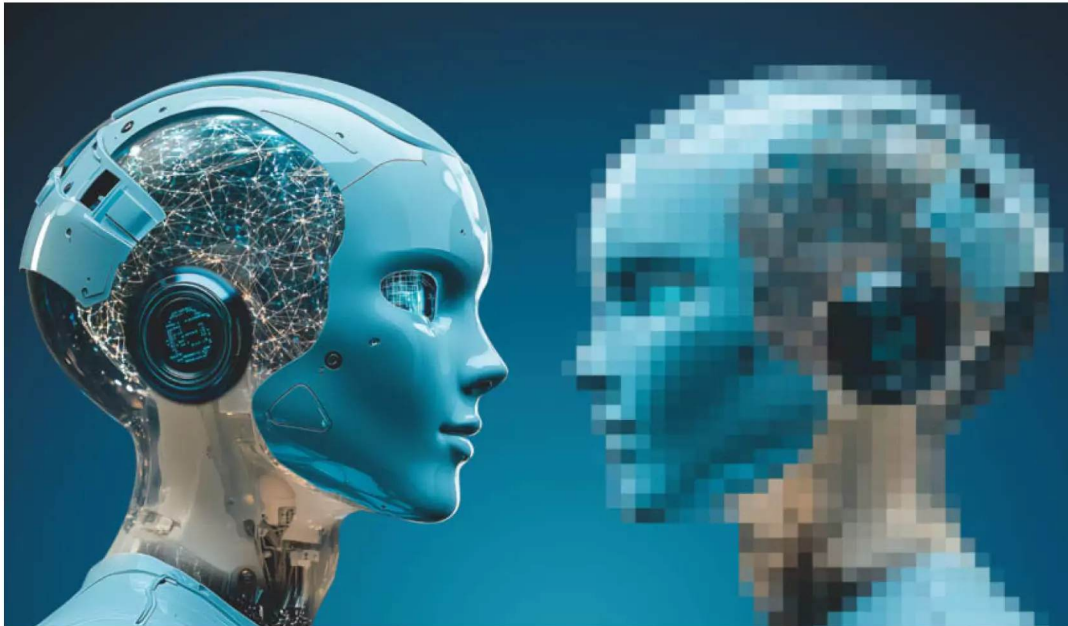


Bild: KI, Collage c't

# Wie Quantisierung Ressourcen spart

Große Sprachmodelle wie ChatGPT benötigen große und teure Server und viel Energie. Man kann sie aber quantisieren, sodass sie mit viel weniger Speicher und Strom auskommen und sogar lokal auf einem Smartphone laufen. Wir erklären, warum quantisierte Modelle viel schneller antworten und trotzdem fast so schlau sind wie die großen Originale.

Von **René Peinl**

**W**ährend der Hype etwas abgeflaut ist und sich Ernüchterung breit macht, steigt das Interesse an selbst betriebener generativer KI. Frei verfügbare Modelle holen gegenüber kommerziellen KI-Angeboten immer mehr auf. Es spricht also einiges dafür, ein eigenes Large Language Model (LLM) zu Hause oder im Unternehmen zu hosten.

Dumm nur, dass die besten Modelle in der Regel auch die höchsten Anforderungen an die Hardware stellen und die meiste Energie verbrauchen. Ein Llama 3.1 mit 405 Milliarden Parametern wird wohl selbst für Freaks mit Power-Workstation im Arbeitszimmer oder Server im Hobbykeller jenseits der Möglichkeiten bleiben. Und auch ein Mittelständler wird sich zweimal überlegen, ob sich die Investition

in zwölf H100-GPUs mit je 80 GByte RAM lohnt, denn schon eine davon kostet schlappe 30.000 Euro.

Zum Glück gibt es kleinere Modelle wie das Mistral Large 2 mit 124 Milliarden Parametern oder das Nvidia Nemotron mit 70 Milliarden Parametern, die dem großen Llama kaum nachstehen. Mit einem kleinen Trick lässt sich Mistral Large 2 auf einer einzigen 80-GByte-Grafikkarte (A100 oder H100) und Nemotron auf einer 48-GByte-GPU (zum Beispiel einer A6000 für rund 6.000 Euro pro Stück) betreiben. Der Trick heißt Quantisierung. Damit ist es inzwischen sogar möglich, Modelle mit der Leistung eines ChatGPT 3.5 lokal auf Smartphones zu betreiben. Dieser Artikel gibt einen Überblick über verschiedene Quantisierungsmethoden und welche Frameworks und Hardwareplattformen sie besonders effizient einsetzen.

### Training mit hoher Genauigkeit

Quantisierung bedeutet, dass man sich Zahlen, die mit vielen Bits in hoher Genauigkeit gespeichert sind, mit weniger Bits in geringerer Genauigkeit nähert. Für Sprachmodelle ist das deshalb interessant, weil es inzwischen empirische Belege dafür gibt, dass die Genauigkeit der Ausgaben von trainierten großen KI-Modellen auch dann kaum abnimmt, wenn sie mit weniger genauen Zahlen arbeiten.

Um das zu verstehen, müssen wir ein wenig ausholen. Tiefe neuronale Netze bestehen aus sehr vielen einzelnen Neuronen. Jedes dieser Neuronen multipliziert die eingehenden Zahlenwerte mit einem sogenannten Gewicht. Die Ergebnisse aller Eingaben werden miteinander verrechnet, im einfachsten Fall addiert. Überschreiten sie einen bestimmten Schwellenwert, werden sie an nachfolgende Neuronen weitergeleitet. Weicht die so erzeugte Ausgabe des gesamten Netzes von der er-

warteten Antwort ab, wird ein Fehler im Ergebnis bestimmt. Dieser Fehler ist ein großer Vektor mit vielen Zahlen.

Die Abweichung zwischen den erwünschten Zahlen und den tatsächlichen Werten wird dann durch die einzelnen Neuronenschichten wieder zurück übermittelt und deren Gewichte ein kleines Stück in die gewünschte Richtung verändert. Wie weit das geht, gibt die sogenannte Lernrate mit Werten zwischen 0 und 1 an. Ist sie hoch (nahe 1), dann lernt das neuronale Netz diesen einen Trainingsfall sehr schnell und produziert das gewünschte Ergebnis.

Allerdings bedeutet eine hohe Lernrate auch, dass das Netz in vielen anderen Trainingsfällen schlechtere Ergebnisse erzielt. Daher wird in der Regel mit einer sehr niedrigen Lernrate (zum Beispiel  $2 \times 10^{-5}$ ) gearbeitet. Um mit vielen Trainingsfällen gute Ergebnisse zu erzielen, benötigt man sehr viele Trainingsläufe, bei denen die Fehlervektoren von Schritt zu Schritt kleiner werden.

Hier kommt die Genauigkeit ins Spiel. Würden die Gewichte ungenau gespeichert, könnten sie bei kleinen Lernraten aufgrund von Rundungsfehlern kaum nachjustiert werden. Weitere Trainingsläufe könnten das Modell dann nicht weiter verbessern. Deshalb wird beim Training meist mit voller Genauigkeit gerechnet. Das sind Gleitkommazahlen (floating point) mit 32 Bit (FP32).

Intelligente Trainingsbibliotheken können selbst entscheiden, wo sie genau rechnen müssen und wo nicht und mit gemischter Genauigkeit arbeiten. Das heißt, sie verwenden dort, wo es notwendig ist, Zahlen mit FP32, ansonsten Zahlen mit halber Präzision (FP16). Das spart nicht nur Speicherplatz, sondern vor allem Rechenzeit. Denn grob gesagt geht eine Multiplikation doppelt so schnell, wenn man die Genauigkeit halbiert.

### Leistungsvergleich von quantisierten LLMs

	Milliarden Parameter	Quantisierung	Dateigröße	MMLU	ARC	HumanEval	IFEval
SmolLM v2 1.7B	1,7	16 Bit	3,4 GByte	35,97	61,55	20,76	18,48
Qwen 2.5 1.5B	1,5	16 Bit	3,1 GByte	60,90	54,70	37,20	29,00
Llama 3.2 3B	3,2	4 Bit	2,0 GByte	63,4	78,6	52	77,4
Phi 3.5 mini 3.8B	3,8	4 Bit	2,4 GByte	69,00	84,60	62,80	57,75

Die Ergebnisse der vier Benchmarks MMLU, ARC, HumanEval und IFEval liegen auf einer Skala von 1 bis 100, größere Werte sind besser. Im Vergleich schneiden die beiden großen, auf 4 Bit quantisierten Modelle Llama 3.2 und Phi 3.5 mini deutlich besser ab als die kleineren Modelle SmolLM und Qwen, die mit nur halb so vielen Parametern ohne starke Quantisierung arbeiten und deutlich mehr Speicherplatz benötigen. (Quelle: Huggingface)



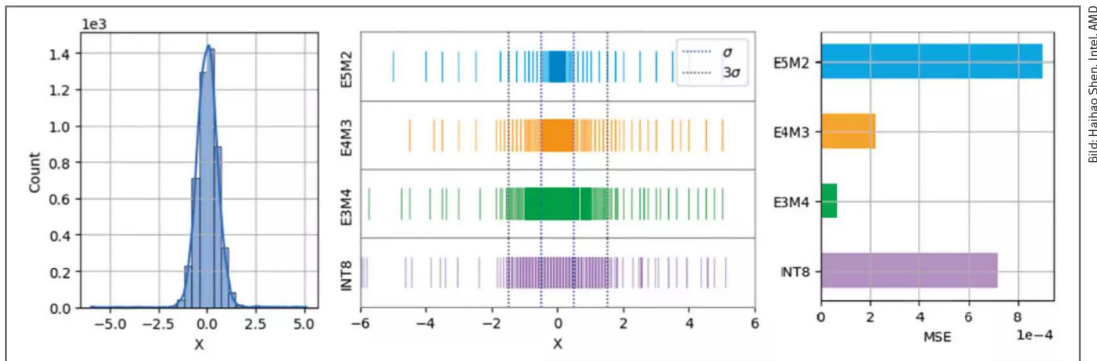


Bild: Haihao Shen, Intel, AMD

Die Bilder zeigen die statistische Verteilung von Fehlerwerten, wie sie bei der Quantisierung von KI-Modellen auftreten. Sie liegen meist im Intervall von  $-2$  bis  $+2$ , das sich mit 8 Bit als Gleitkommazahlen besonders gut abbilden lässt, die das Format E3M4 oder E4M3 statt E5M2 oder Integer (INT8) nutzen. Rechts sieht man die mittlere quadratische Abweichung, die beim grünen E3M4 unter 1 liegt.

## Nach dem Training

Nach Abschluss des Trainings benötigen die Modelle im Betrieb (dem sogenannten Inferencing) keine hohe Genauigkeit. Deshalb wird zum Beispiel beim Inferencing immer nur mit FP16-Werten gerechnet, weil es praktisch keinen Unterschied zu FP32-Werten macht.

Um den Einfluss weiterer Quantisierungen zu untersuchen, veröffentlichte ein Forscherteam um Elias Frantar im März 2023 ein Konferenzpapier, in dem sie den damals neuen Quantisierungsalgorithmus GPTQ (Post Training Quantization for Generative Pre-trained Transformers) vorstellten. Mit ihm gelang es, verschiedene große Sprachmodelle wie BLOOM und OPT mit 175 Milliarden Parametern in nur wenigen Rechenstunden von 16 Bit auf 4 Bit zu reduzieren. Der Speicherbedarf sank um den Faktor 3,25, wodurch sie erstmals ein Modell mit 175 Millionen Parametern auf nur einer GPU betreiben konnten, was Anschaffungskosten, Stromverbrauch und Rechenzeit spart.

Als Maß für die Qualität der Antworten der Sprachmodelle verwendeten sie die Perplexität der Ausgaben. Sie misst die Unsicherheit eines Sprachmodells, einen Satz korrekt zu vervollständigen. Sie kann auch als „Überraschung“ des Modells interpretiert werden: Ist die Perplexität hoch, enthalten die Trainingsdaten kaum ähnliche Formulierungen wie die Ausgabe oder das Modell hat sie nicht ge-

lernt. Die Perplexität (PPL) ist daher ein einfaches Maß zur Abschätzung des Qualitätsverlustes einer Quantisierung. Je kleiner der PPL-Wert, desto besser das Ergebnis.

Die Untersuchungen der Forscher zu GPTQ zeigten, dass sich die Perplexität bei einer 4-Bit-Quantisierung im Vergleich zu einer 16-Bit-Quantisierung bei großen Modellen mit 175 Milliarden Parametern um weniger als 0,4 Prozent verschlechtert.

Waren die Ausgangsmodelle dagegen kleiner und hatten nur 350 Millionen Parameter, verschlechterte sich die Ausgabeperplexität bei einer Quantisierung von 16 auf 4 Bit um über 10 Prozent. Bei kleinen Modellen geht eine so starke Quantisierung also mit erheblichen Qualitätsverlusten einher.

## Neue Zahlenformate

Für eine gute Quantisierung ist es wichtig, dass die Abweichungen zwischen den Gewichten mit hoher Genauigkeit und den quantisierten Werten nicht zu groß werden. Die meisten Werte bewegen sich in einem Intervall von  $-2$  bis  $+2$ , daher ist in diesem Bereich eine besonders hohe Genauigkeit gefragt. Um diesen Zahlenbereich mit hoher Genauigkeit abzubilden, stellte ein Forscherteam um Haihao Shen von Intel und AMD im Mai 2024 angepasste Formate für 8-Bit-Gleitkommazahlen (FP8) vor.

FP8 verwenden acht Bits zur Darstellung von Gleitkommazahlen mit Vorzeichen (s), normalisier-

ter Mantisse (M) und Exponent (E) nach der vereinfachten Formel  $(-1)^s \times M \times 2^{(E - \text{Bias})}$ , wobei die Mantisse Werte zwischen 1 und 2 annimmt. Je nach Zahlenraum, den man damit abbilden will, kann man abseits des Vorzeichenbits die übrigen sieben Bits unterschiedlich auf den Exponenten und die Mantisse verteilen. Der Bias ist ein je nach Anzahl der Exponentenbits festgelegter Wert, um positive wie negative Exponenten ohne weiteres Vorzeichenbit zu erhalten.

Da man keine hohen Werte benötigt, eignen sich die Formate E4M3 und E3M4 deutlich besser als E5M2, da sie Absolutwerte bis 448 beziehungsweise 30 mit höherer Genauigkeit darstellen. Gewichte im Intervall von -2 bis +2 können sie deutlich genauer speichern und die Quantisierungsfehler drastisch reduzieren. Nvidia-Grafikkarten ab der Ada-Lovelace-Generation (RTX 4xxx, H100 und andere) unterstützen diese FP8-Formate nativ und können damit quantisierte Modelle sehr schnell ausführen.

## Quantisierung für mobile Geräte

Da GPTQ bei kleinen bis mittelgroßen Modellen mit relativ hohen Qualitätsverlusten einhergeht, stellte ein internationales Forscherteam um Ji Jin im Frühjahr 2024 ein verbessertes Verfahren namens AWQ (Activation-Aware Weight Quantization) vor. Zusammen mit AWQ entwickelten die Forscher eine Laufzeitumgebung namens TinyChat, mit der sie ein Llama-2-Modell auf 13 Milliarden Parameter verkleinerten, sodass es erstmals auf einer einzigen mobilen GPU mit 8 GByte Speicher und einer Leistungsaufnahme von nur 15 Watt lief. Bei einer Quantisierung auf 4 Bit verschlechterte sich die Perplexität der Ausgabe gegenüber der großen 16-Bit-Version auch bei den relativ kleinen Modellen nur um wenige Prozentpunkte.

Um dies zu erreichen, bedienten sich die Forscher eines Tricks. Sie beließen etwa 0,1 bis 1,0 Prozent der besonders häufig aktivierten Gewichte im Modell bei der vollen 16-Bit-Genauigkeit und reduzierten den großen Rest auf 4 Bit. Es gibt nämlich nicht nur einen Gewichtstyp in neuronalen Netzen, sondern verschiedene Parametertypen. Bei einigen spielt die Genauigkeit eine größere Rolle, bei anderen eine geringere. Dementsprechend werden die Parameter mit mehr oder weniger Bits gespeichert. Auf dem Weg dorthin mussten die Forscher allerdings eine Reihe von Hürden überwinden, die sie in ihrem Forschungspapier näher beschreiben. Sie finden alle hier aufgeführten Studien unter [ct.de/wcfh](https://ct.de/wcfh).

Bild: Ji Jin, MIT



**Während eine plumpe Rundungsmethode das quantisierte Modell OpenFlamingo-9B zu der Beschreibung verleitete, der Mann würde einen Baby-Elefanten auf dem Arm halten, erklärte das mit AWQ quantisierte Modell, dass der Mann und seine Tochter mit einem Elefanten posieren.**

Ein guter Indikator für die benötigte Speichergröße im GPU-RAM ist die jeweilige Dateigröße der Modelle. Für das Inferencing auf der GPU sollte man 10 bis 20 Prozent Puffer zur Dateigröße hinzurechnen. Dies hat sich als Faustregel bewährt. Beispielsweise läuft das Modell Mistral Large 2 mit 65 GByte Dateigröße (4 Bit AWQ) auf einer 80-GByte-GPU halbwegs ordentlich.

## Tempomacher

Die Quantisierung soll aber nicht nur den Speicherbedarf reduzieren, sondern auch die Berechnungen beschleunigen. Bei Sprachmodellen misst man die Ausgabegeschwindigkeit in Token pro Sekunde (T/s). Ein Token ist durchschnittlich etwa vier Zeichen lang.

AWQ und GPTQ haben gemeinsam, dass sie nicht besonders schnell auf CPUs laufen. Ihr volles Potenzial entfalten sie auf CUDA-Kernen, die für schnelle Berechnungen auf Nvidia-Grafikkarten optimiert sind. Dabei muss die ganze Kette aus Hardware, Treiber, CUDA Framework, Kernel, Programmbibliothek (zum Beispiel PyTorch) und neuronalem Netz stimmen, damit die Beschleunigung greift. Passt der Kernel oder die PyTorch-Variante nicht zur CUDA-Version, wird es sofort deutlich langsamer.

Die AWQ-Entwickler konnten ihre geschrumpften Modelle mithilfe optimiertem CUDA-Kernel gegen-



## Perplexity-Vergleich von quantisierten LLMs

LLM-Größe	350M	1.3B	2.7B	6.7B	13B	30B	66B	175B
Original (16 Bit)	22,00	14,63	12,47	10,86	10,13	9,56	9,34	8,34
Quantisiert (4 Bit)	24,24	15,47	12,87	11,39	10,31	9,63	9,55	8,37
Unterschied (absolut)	2,24	0,84	0,4	0,53	0,18	0,07	0,21	0,03
Unterschied (prozentual)	10,2 %	5,7 %	3,2 %	4,9 %	1,8 %	0,7 %	2,3 %	0,4 %

Die Tabelle vergleicht die Perplexity-Werte von LLMs mit verschiedenen Parametergrößen von 350 Millionen (350M) bis 175 Milliarden (175B). Angegeben sind die Perplexity-Werte für 16-Bit-Genauigkeit und 4-Bit-Quantisierung. Je geringer die Perplexität, desto besser. (Quelle: <https://arxiv.org/pdf/2210.17323>)

über GPTQ sogar deutlich beschleunigen. So katalysierten sie laut ihrer Untersuchung beispielsweise das frei verfügbare Modell Mistral 7B mit einer AWQ-Quantisierung auf 4 Bit in ihrer eigens entwickelten Laufzeitumgebung TinyChat auf eine Ausgabe von 38 T/s. Im Vergleich dazu lief eine mit AutoGPTQ quantisierte Version mit der Transformer-Bibliothek von Huggingface nur mit 14 T/s.

Ein Teil des Unterschieds ist auf die Optimierung des Codes in den beiden Laufzeitumgebungen zurückzuführen. So erreichte die 16-Bit-Version des Sprachmodells Falcon 7B mit der Transformer-Bibliothek von Huggingface in einem weiteren Test mit einigen Codeoptimierungen immerhin 53 T/s. Sie blieb jedoch hinter der AWQ-Quantisierung zurück, die mit 124 T/s mehr als doppelt so schnell war.

Derzeit ist das Zahlenverhältnis der auf der Internetplattform Huggingface veröffentlichten Modelle, die mit GPTQ und AWQ quantisiert wurden, noch relativ ausgeglichen. Im November fanden wir unter den insgesamt 1,1 Millionen Modellen etwa 5000 GPTQ-quantifizierte und 4000 AWQ-quantifizierte Modelle. Das Verhältnis wird sich aber in Zukunft aufgrund der Performancevorteile wahrscheinlich zugunsten von AWQ verschieben.

## Quantisierung für CPUs

Während GPTQ und AWQ vor allem für Nvidias GPUs optimiert sind, nutzt das Quantisierungsverfahren GGUF (Georgi Gerganov Unified Format) Optimierungen, um auch auf CPUs akzeptable Geschwindigkeiten zu erreichen.

Der KI-Entwickler Georgi Gerganov nutzt mit seinem auf GitHub veröffentlichten Projekt llama.cpp, in dessen Kontext GGUF und dessen Vorgänger GGML entstanden sind, die Fähigkeiten moderner CPUs aus. Dazu gehören die Erweiterungen AVX, AVX2, AVX512 und AMX. Die Advanced Matrix Extensions (AMX) sind Intels neueste Befehlssatzerweiterungen. Sie beschleunigen die Ausführung von Ma-

trixmultiplikationen, wie sie für neuronale Netze entscheidend sind.

Es ist aber auch möglich, eine oder mehrere Schichten des neuronalen Netzes auf die GPU auszulagern (Offloading) und damit die Performance weiter zu steigern. Denn GPU-Recheneinheiten sind für Matrixmultiplikationen noch weitaus besser optimiert als CPUs mit Vektorerweiterungen. So schafft ein moderner AMD Ryzen 7 7800x3D oder ein Intel Core i7 14700k mit dem Modell Mistral 7B nur knapp 10 T/s. Ein Apple M3 Pro schafft mit gut 20 T/s etwa die doppelte Geschwindigkeit. Eine mobile Geforce 4060 GPU im Notebook verdoppelt die Geschwindigkeit auf 38 T/s. Noch schneller sind jedoch Desktop-GPUs wie die AMD RX 6800XT mit 16 GByte oder die Nvidia RTX 4070 Super. Sie rechnen mit 52 T/s beziehungsweise 58 T/s etwa fünf- bis sechsmal so schnell wie eine moderne CPU.

Aufgrund der guten CPU-Optimierung durch die populäre LLM-Laufzeitumgebung ollama hat sich GGUF zum heimlichen Star der KI-Szene entwickelt. Wenn man bei Huggingface nach Modellen mit GGUF-Quantisierung sucht, findet man nicht weniger als 46.000 Treffer. Ein weiterer Grund für die Beliebtheit ist, dass die Modelle in der Regel nicht nur in einer 4-Bit-Variante vorliegen, wie es bei AWQ und GPTQ der Fall ist. Vielmehr gibt es meist eine ganze Reihe verschiedener Versionen, die von 8 Bit bis 2 Bit reichen.

## Weitere Alternativen

Als weiteres Quantisierungsformat ist noch Bits&Bytes (bnb) zu erwähnen, das auf Huggingface mit etwa 4000 Modellen vertreten ist. Diese sind allerdings im Vergleich schnarchlangsam, weshalb sich das bnb-Format nur für erste Experimente mit Modellen eignet, die in keinem anderen Quantisierungsformat verfügbar sind.

Im Februar 2024 stellten KI-Forscher um Shumin Ma von Microsoft Research ein Verfahren vor, das



Gewichte mit jeweils nur einem sogenannten Trit speichert. Dieses kann in einem ternären System die Werte 1, 0 und -1 annehmen. Ein Trit benötigt etwa den 1,585-fachen Speicherplatz ( $\log_2(3)$ ) eines Bits. Damit hat Microsoft im Vergleich zu einem Modell mit FP16 eine Verkleinerung um den Faktor 3,55 erreicht. Das ist nur geringfügig besser als bei einer 4-Bit-Quantisierung mit AWQ oder GGUF (Faktor 3,2). Das Trit-Verfahren könnte trotzdem in Zukunft eine größere Rolle spielen. Zumal Microsoft auch ein entsprechendes Framework vorgestellt hat, das die Ausführungsgeschwindigkeit solcher hochquantisierter Modelle optimiert.

## Modellvergleich

Angesichts der großen Fortschritte bei der Quantisierung von Sprachmodellen stellt sich die Frage, ob es für mobile Geräte mit begrenzten Ressourcen besser ist, ein kleines Modell mit wenigen Parametern und hoher Genauigkeit der Einzelgewichtungen oder ein größeres Modell mit vielen Gewichtungen geringerer Genauigkeit zu verwenden.

Dazu haben wir in der Tabelle die Performance von kleinen Modellen (Qwen 2.5 1.5B und SmolLM V2 1.7B) in verschiedenen Benchmarks mit Modellen verglichen, die mit einer 4-Bit-Quantisierung etwa doppelt so viele Parameter haben, aber deutlich weniger Speicherplatz benötigen (Llama 3.2 3B und Phi 3.5 min 3.8B).

Es wird schnell klar, dass ein Phi 3.5 mini in der quantisierten Version selbst bei einem Genauigkeitsverlust von 10 Prozent gegenüber der FP16-Version immer noch deutlich besser abschneidet als ein SmolLM v2. Und auch das in einigen Bereichen dem SmolLM überlegene Qwen 2.5 1.5B steckt es noch locker in die Tasche, obwohl es weniger Speicher benötigt und bei entsprechender Hardwareunterstützung schneller ist.

Mit der Android-App „H2O AI Assistant“ kann man auf leistungsfähigen Android-Geräten selbst ausprobieren, wie gut das funktioniert. Auf dem Asus Zenfone 9 des Autors lief ein Phi 3.5 mini mit geschätzten 2 Token pro Sekunde. Das ist nichts für Ungeduldige, aber angesichts der einfachen Inbetriebnahme und der Tatsache, dass es komplett lokal auf dem Smartphone läuft, eine beeindruckende Leistung. Die Qualität der Antworten ist oft besser als bei der ersten Version von ChatGPT (v3.5), die vor gerade einmal zwei Jahren veröffentlicht wurde und die Fachwelt in Erstaunen versetzte. Ohne Quantisierung wäre das undenkbar. (hag) 

Studien zu  
Quantisierungsverfahren:  
[ct.de/wcfh](https://ct.de/wcfh)

## Spielend leicht umsteigen!

Einstieg, Gaming und  
praktische Tools auf über 140 Seiten



Diese Top-Themen warten  
u.a. auf Sie:

- Die richtige Distribution, praktische Tools, Troubleshooting
- Blockbuster-Games wie unter Windows spielen
- Installieren und Software auswählen
- Mint oder Fedora statt Windows 11
- Tipps für den Linux-Alltag
- Windows-programme unter Linux



**JETZT  
BESTELLEN!**

Heft für 14,90 € • PDF für 12,99 €  
Heft + PDF 19,90 €

 [shop.heise.de/ct-linuxguide24](https://shop.heise.de/ct-linuxguide24)

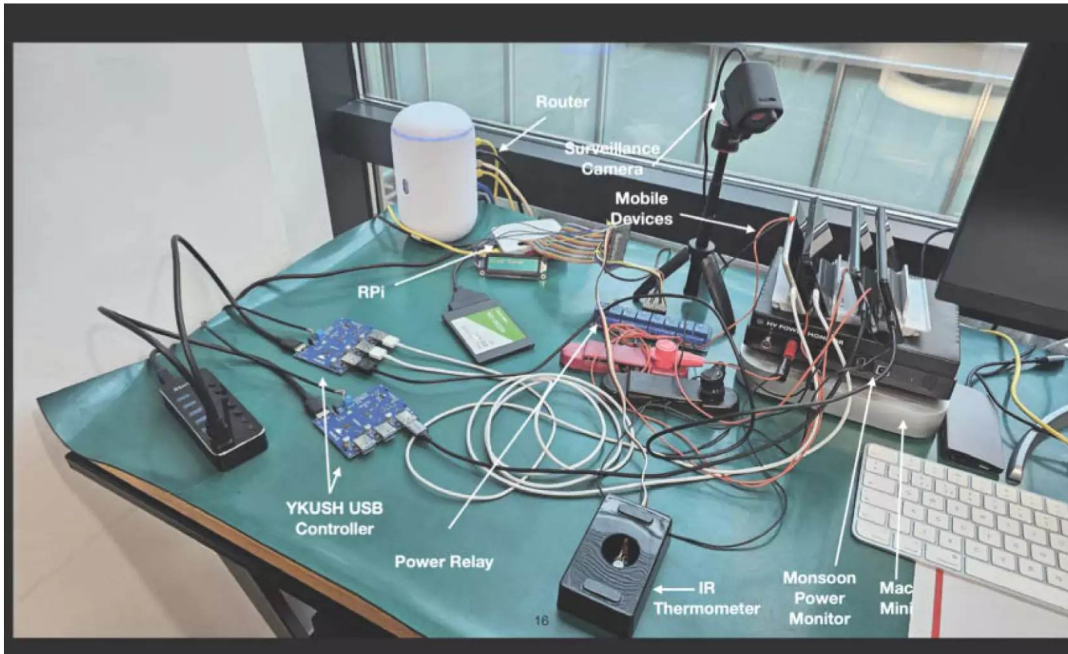


Bild: Brave

# Generative KI ohne Cloud

Ein Forschungsteam hat erstmals systematisch untersucht, wie gut große und kleine KI-Sprachmodelle auf Android- und iOS-Smartphones oder sogar auf hypothetischen, KI-fähigen Routern und Smart-TVs laufen würden.

Von **Monika Ermert** und **Dušan Živadinović**

**E**igentlich ist die Messung der KI-Leistung von Anwendergeräten wie Smartphones nichts Neues: Programme wie Geekbench haben schon mit Erscheinen der ersten, auf Machine Learning spezialisierten Rechenkerne deren Geschwindigkeit gemessen. Und seit August 2024 kann man mit Geekbench AI ermitteln, wie schnell ein Smartphone oder ein PC bestimmte KI-Aufgaben löst. Dabei geht es um Dinge wie die Bild- und Textklas-

sifizierung, Gesichts- und Objekterkennung, Tiefenmessung oder maschinelle Übersetzung.

Doch damit gewinnt man nur einen Eindruck von der Leistung eines im Betriebssystem eingepflanzten kleinen Sprachmodells, das aufwendige Aufgaben dem um Größenordnungen schnelleren großen Bruder in der Cloud zuschiebt.

Es gibt aber gute Gründe, private Daten lokal zu verarbeiten. Oft ist unklar, was der Cloud-Betreiber

mit den erhaltenen Daten anstellt. Außerdem sind Cloud-Dienste generell ein attraktives Ziel für Angreifer.

Aus diesem Dilemma führen mehrere Auswege. Betreiber könnten ihre Clouds härten, etwa so, wie Apple mit Private Cloud Compute vorführt (ct.de/wkzf). Zunehmend rückt auch der lokale Betrieb in den Fokus: Da Smartphones immer mehr leisten, könnten sie große Sprachmodelle vielleicht bald selbst ausführen. Außerdem sind potente, mit KI-Beschleunigern ausgerüstete Heimnetzgeräte wie Heimserver oder NAS-Geräte vorstellbar, an die die Smartphones ihre Prompts weitergeben.

Um die Erfolgsaussichten solcher Szenarien zu erkunden, hat ein Team des Browserherstellers Brave, das der Spezialist für On-Device- und Federated-Learning, Stefanos Laskaridis, leitet, den Benchmark Mobile Evaluation of Language Transformer entwickelt (MELT). Damit, so schreiben die Autoren in einer auf ArXiv erschienenen Facharbeit (ct.de/wkzf), lassen sich verschiedene, „populäre, instruktionsoptimierte Sprachmodelle“ und Frameworks systematisch auf verschiedenen Plattformen analysieren.

## Token, Strom und Temperatur

Anders als etwa Geekbench AI erfasst MELT nicht nur den KI-Durchsatz in Form von Tokens pro Sekunde, sondern auch die Leistungsaufnahme, die Temperaturentwicklung oder die Akkulaufzeiten der Geräte. Die Vielzahl der Messungen ergibt ein umfassendes Bild darüber, ob ein auf Anwendergeräten betriebenes Sprachmodell praxistauglich ist.

Laskaridis sagte im Gespräch mit c't: „Unsere Hauptmotivation war es, große Sprachmodelle Privacy-orientiert aufzusetzen. Die Endgeräte werden immer mächtiger. Wir wollten ausprobieren, ob sich große Sprachmodelle direkt auf dem Endgerät implementieren lassen und das Verhalten sowie die Leistungsgrenzen der Geräte erfassen.“

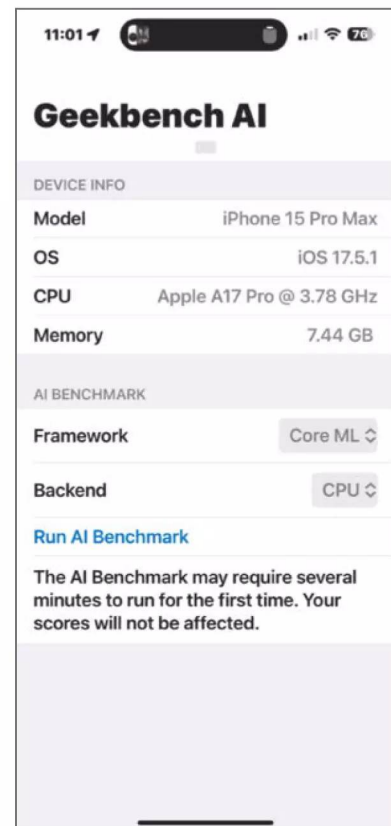
Dafür hat das Team einige Hürden überwunden und die detaillierten Lösungsbeschreibungen machen den Aufsatz zusätzlich interessant. Beispielsweise bringen nur Android-Geräte eine standardisierte Schnittstelle für die Automatisierung mit, die Android Debug Bridge. Für iOS hat das Team einen Service in Python geschrieben, der für Touch- und Wischgesten sowie für Texteingaben eine Bluetooth-Tastatur und -Maus simuliert.

Im Fokus der ersten Tests standen die Apple-Smartphones iPhone 14 Pro und iPhone SE sowie die Android-Modelle Pixel 6a und Galaxy S23. Das

Pixel 6a und das iPhone SE stehen stellvertretend für ältere Mittelklassegeräte. Das Samsung Galaxy S23 und das iPhone 14 Pro gehören zur Oberklasse, die vor zwei bis drei Jahren erschienen ist; heute gibt es schnellere Nachfolger.

## KI-Router am Horizont

Künftige Smartphones könnten KI-Aufgaben per WLAN an Heimserver, Router oder NAS-Geräte mit großen Sprachmodellen auslagern; so bleiben Anwenderdaten im privaten Bereich. Beides zusammen, die Auslagerung und der Privatsphärenschutz, erscheint als tragfähige Grundlage für persönliche Assistenten mit kontextsensitivem Verhalten und Langzeitgedächtnis.



**Schmucklos, aber schnell: Mit Geekbench AI lässt sich nur die Leistung von eingebauten Sprachmodellen der Smartphones ermitteln.**



Um die KI-Leistung solcher hypothetischer Geräte zu simulieren, haben die Forscher die von Nvidia gefertigten und mit Ubuntu Linux betriebenen Entwicklerplattformen Jetson Orin AGX und Nano herangezogen. Mit beiden lassen sich Geräte unterschiedlicher Leistung simulieren, indem man unterschiedliche Energieniveaus konfiguriert; davon hängen die Zahl und der Takt der aktiven Kerne und anderes mehr ab. Der kleine Nano läuft wahlweise mit einer Thermal Design Power von 7 oder 15 Watt, der AGX mit 15, 30 oder sogar 50 Watt.

## Mac kompiliert, Raspi schaltet

Die für den Test erforderlichen Apps hat das Team auf einem Mac Studio kompiliert und per WLAN an die Geräte verteilt. Ein Raspberry Pi 4 koordinierte die Messungen und erfasste die Werte. Für Hardwarebastler und Messtechniker besonders lesenswert erscheinen die Passagen zur Strommessung und -versorgung: Die Akkus wurden mitsamt der Akku-Controller entfernt und die Geräte direkt mit Strom versorgt.

Das ist aus zwei Gründen nützlich: So lässt sich die Leistungsaufnahme direkt erfassen; die Forscher setzen dafür den Monsoon Power Monitor mit einer Abtastrate von 5 kHz ein, was granulare Messwerte liefert. Außerdem liegt so der Zugang zum Stromanschluss der Geräte frei, sodass sie sich bedarfsgemäß ein- und ausschalten lassen. Das erledigte wiederum der Raspi, und zwar über ein an der GPIO-Pfostenleiste angeschlossenes programmierbares Relais.

Auf den Geräten liefen mehrere Sprachmodelle unterschiedlicher Leistungsklassen; im KI-Sprech sind damit Parametermengen gemeint. Das kleinste Modell TinyLlama verarbeitet bis zu 1,1 Milliarden Parameter, das größte – Llama 2 – bis zu 13 Milliarden.

Mit den kleineren Modellen kamen die Testmuster erwartungsgemäß besser zurecht. Ein Maß dafür ist die Menge der berechneten Tokens in verschiedenen Phasen der Verarbeitung. In der anfänglichen Prefill-Phase gab beispielsweise das iPhone 14 Pro mit TinyLlama bis zu 170 Tokens pro Sekunde aus. Das Galaxy lieferte nur 35 Tokens pro Sekunde. Die Mittelklassehandys blieben noch weiter zurück. In der zweiten Phase, der Generierung neuer Tokens, gingen alle Smartphones in die Knie: Selbst das iPhone 14 Pro schaffte nur noch 25 Tokens pro Sekunde.

Das iPhone 14 Pro antwortete fast immer am schnellsten; das Galaxy S23 war aber beispielsweise mit dem 3-Milliarden-Parameter-Modell Sta-

bleLM Zephyr3b schneller. Das Jetson Orin AGX übertraf die Smartphones um Längen, wenn die CPU bis zu 50 Watt verbraten durfte: Mit TinyLlama generierte es 70 neue Tokens pro Sekunde. 50 Watt, das liegt allerdings weit über dem Niveau heutiger Router. Viele nehmen nur 10 bis 20 Watt auf, und zwar für alle Systemkomponenten zusammen.

Überrascht waren die Forscher von der hohen Leistungsaufnahme der Smartphones. Das iPhone 14 Pro zog im Mittel 13,8 Watt und 18 Watt maximal, das Galaxy S im Schnitt 8,5 Watt und 14 Watt maximal. Je nach KI-Modell und Framework erledigten die Geräte mit der Energiemenge, die normalerweise ihr Akku liefert, zwischen 490 (Galaxy) und 590 Prompts (iPhone). Auch fiel die Wärmeentwicklung auf: Das iPhone erwärmte sich teils auf über 47 Grad, was zusammen mit dem hohen Stromverbrauch die Praxistauglichkeit erheblich einschränkt.

## Mehr messen mit MELT

Der neue KI-Benchmark MELT erscheint wertvoll, weil er die Grundlage für Messungen kommender Geräte legt und kleinen Laboren viel Know-how vermittelt. Dazu trägt besonders bei, dass das Wissenschaftler-Team MELT als Open Source freigegeben hat ([ct.de/wkzf](https://ct.de/wkzf)). MELT regt aber auch zu Ideen an. Möglicherweise reizt es Hersteller, ihre Smart-TVs, Router oder NAS-Geräte für KI-Aufgaben zu rüsten.

Anders als die Werbefilmchen von KI-Giganten wie Google oder OpenAI bringt es MELT auch für Anwender nachvollziehbar an den Tag: Die meisten heutigen Smartphones wären gerne groß, sind es aber nicht. „Letztlich ergibt nicht alles, was an lokalisierten großen Sprachmodellen eingesetzt werden kann, wirklich Sinn“, fasst Laskaridis zusammen. „Man kann ein Modell vielleicht auf ein Gerät bringen, aber das heißt nicht zwangsläufig, dass das schon eine Qualität hat, die den Aufwand rechtfertigt.“ Für heutige Smartphones eignen sich lediglich kleinere, angepasste Sprachmodelle.

Doch daran, dass größere KI-Sprachmodelle eines Tages auf Smartphones einziehen werden, hat das Team kaum Zweifel: „Das Ökosystem steckt noch in den Kinderschuhen und sowohl schlauere Algorithmen als auch schnellere Hardware können die Lage schnell verändern. Wir gehen davon aus, dass die Beschleunigung von KI-Kernen und das Co-Design von Hardware und Frameworks den größten Beitrag für einen effizienten Standalone-Betrieb leisten werden.“ (dz) **ct**

**MELT auf GitHub**  
[ct.de/wkzf](https://ct.de/wkzf)

# KI-WORKSHOPS



14. März

## IT-Security: Künstliche Intelligenz für effiziente Verteidigungsstrategien

Lernen Sie, für welche IT-Sicherheitsprozesse der Einsatz von KI sinnvoll ist und wann gängige Methoden Vorrang haben sollten.



14. – 15. April

## Deep Dive in die OpenAI API: KI in eigene Anwendungen integrieren

Lernen Sie, wie Sie die OpenAI API in Ihre Projekte integrieren können, um innovative KI-Lösungen zu entwickeln.



14. – 17. April

## Große Sprachmodelle feintunen und einsetzen

Lernen Sie aktuelle Deep-Learning-Technologien kennen und verstehen Sie, wie Sie richtig einzusetzen.



20. – 22. Mai

## KI für eigene Daten: Retrieval Augmented Generation (RAG) in der Praxis mit LangChain.js

Sie lernen, wie Sie maßgeschneiderte KI-Lösungen mit RAG entwickeln, effizient implementieren und Prompts optimieren.



Bild: KI Midjourney | Collage c't

# Chatbots und der Datenschutz

Die deutschen Datenschutzbehörden bieten eine „Orientierungshilfe“ zum Umgang mit generativer Sprach-KI an. Doch was Unternehmen und Behörden helfen soll, bleibt schwammig und könnte sie an mancher Stelle sogar vor unlösbare Probleme stellen.

Von **Holger Bleich**

**D**en Siegeszug generativer Künstlicher Intelligenz (KI) haben die Macher der Datenschutz-Grundverordnung (DSGVO) 2016 nicht vorhersehen können. Deshalb spielt KI in dem Regelwerk kaum eine Rolle. Der Gesetzgeber ging vom klassischen Input-Output-Schema aus: Personenbezogene

Daten werden erhoben, gespeichert, verarbeitet und ausgegeben. Lediglich Art. 22 DSGVO beschäftigt sich rudimentär mit KI. Er legt fest, dass automatisierte Verarbeitungen und Entscheidungen, die Menschen betreffen, immer auch von Menschen absegnet werden müssen.



Generative KI-Modelle, insbesondere Large Language Models (LLMs) wie GPT, stellen die europäischen Datenschutzbehörden deshalb vor Herausforderungen, die sie kaum in den Griff bekommen können. Als besonders problematisch erweisen sich Modelle, die nicht lokal arbeiten, sondern bei einem Anbieter liegen, der den Zugang über ein Webfrontend und/oder ein API gewährleistet. In aller Regel weiß niemand genau, welche personenbezogenen Daten zum Training dienen, wie das Modell sie genau verarbeitet und welcher Output entstehen könnte. Nicht einmal die Anbieter selbst können da immer Auskunft geben.

Unter solchen Umständen fällt es Unternehmen und Behörden schwer, LLM-Chatbots DSGVO-konform einzusetzen. Was müssen sie beachten, um nicht in Konflikt mit ihren Mitarbeitern, Kunden und den zuständigen Datenschutzbehörden zu geraten? Es fehlen Leitplanken für die sogenannte Compliance, also konkrete Anforderungskataloge, die man abarbeiten kann.

In diese Bresche sind mittlerweile einige EU-Aufsichtsbehörden gesprungen. Die französische Commission Nationale de l'Informatique et des Libertés (CNIL) etwa hat einen „Self-assessment guide for

artificial intelligence (AI) systems“ in Form von konkreten Checklisten veröffentlicht, der französischen Unternehmen helfen soll (siehe [ct.de/wcge](https://ct.de/wcge)). Diese erfahren damit zumindest, wie die nationale Datenschutzaufsicht KI und die DSGVO zusammenbringt und können darauf reagieren.

Die föderale Gliederung der Aufsicht macht dies in Deutschland schwieriger. Einige Landesaufsichtsbehörden haben Checklisten oder ähnliche Papiere erarbeitet, etwa die aus Hamburg, Bayern und Baden-Württemberg [1]. Was bislang fehlt, ist eine gemeinsame Position, die nur von der Datenschutzkonferenz (DSK) als Diskussions- und Beschlussgremium aller Landesdatenschützer und des Bundesdatenschutzbeauftragten geäußert werden kann. Einen Beschluss gibt es zwar nicht, aber immerhin hat die DSK im Mai 2024 eine 15-seitige „Orientierungshilfe KI und Datenschutz“ veröffentlicht (siehe [ct.de/wcge](https://ct.de/wcge)).

## Grobe Auslegung

Die Orientierungshilfe diene „im Sinne einer Checkliste als Leitfaden insbesondere für datenschutzrechtlich Verantwortliche“, betonte die DSK zur Ver-



The screenshot shows the CNIL website. At the top, the CNIL logo is displayed next to the text: "PROTÉGER les données personnelles", "ACCOMPAGNER l'innovation", and "PRÉSERVER les libertés individuelles". Below this is a navigation bar with links: "MY COMPLIANCE TOOLS", "DATA PROTECTION", "TOPICS", "THE CNIL", and a search icon. The main content area has a breadcrumb trail: "Home > Self-assessment guide for artificial intelligence (AI) systems". The title "Self-assessment guide for artificial intelligence (AI) systems" is prominently displayed. Below the title, a paragraph states: "The CNIL offers organisations an analysis grid through which to assess by themselves the maturity of their artificial intelligence systems with regard to the GDPR. It also describes best practices that can be followed."

**Die französische Datenschutzbehörde CNIL hilft Unternehmen auch in englischer Sprache mit konkreten Checklisten zur KI-Nutzung weiter.**

öffentlichung. Doch wer sich eine echte Checkliste zum Abhaken erhoffte, wie sie die CNIL herausgebracht hat, wurde enttäuscht: Die DSK wiederholt in vielen Passagen grobe Auslegungen der DSGVO, wie sie etwa bereits im Hamburger Leitfaden von Ende 2023 zu lesen waren.

So erklärt die DSK, dass eine Organisation allein Verantwortliche im Sinne der DSGVO ist, wenn sie die KI-Anwendung selbst betreibt. Nutze sie dagegen einen externen Anbieter, bestehe „zwischen dem Anbieter und dem Verantwortlichen häufig ein Auftragsverarbeitungsverhältnis gemäß Art. 28 f. DSGVO“. Deshalb solle man daran denken, mit genutzten Anbietern eine Auftragsverarbeitungsvereinbarung abzuschließen. Außerdem solle man prüfen, ob der Anbieter Daten ins außereuropäische Ausland mit einer geltenden Rechtsgrundlage übermittelt.

Bevor eine Organisation LLMs von externen Anbietern einsetzt, muss sie sich der DSK zufolge umfangreich über deren Funktionsweise informieren und vorab das Risiko „für die Rechte und Freiheiten natürlicher Personen“ einschätzen. „Vielfach“ werde dann eine Datenschutzfolgeabschätzung (DSFA) nach Art. 35 DSGVO erforderlich sein. Konkrete Beispiele, wann dies der Fall sein könnte, liefert die Orientierungshilfe nicht und bleibt an dieser Stelle deshalb wenig hilfreich.

Die DSK empfiehlt Unternehmen und Behörden, klare interne Weisungen zur KI-Nutzung zu erteilen. Insbesondere seien die konkreten Zwecke mitzuteilen und untersagte Einsatzszenarien idealerweise anhand konkreter Beispiele zu benennen. Für die berufliche Nutzung sollen Arbeitgeber den Beschäftigten Funktionsaccounts inklusive Funktions-Mail-Adresse des Unternehmens bereitstellen. Diese Accounts sollten demnach nicht auf einzelne Mitarbeiter rückführbar sein.

Bei allen Überlegungen soll der Orientierungshilfe zufolge das Prinzip „data protection by design“ gemäß Art. 25 DSGVO im Vordergrund stehen. So heißt es: „Zum Beispiel können für Accounts, die Beschäftigte nutzen sollen, die Funktionen zur Nutzung von Eingaben für das KI-Training und zur Eingabe-Historie bereits bei Inbetriebnahme eines Accounts so gewählt werden, dass keine Eingabedaten zu Trainingszwecken verarbeitet werden und keine Eingabe-Historie über die Sitzung hinaus gespeichert wird.“

Darüber hinaus empfiehlt die DSK, beim Prompting auf personenbezogene Daten zu verzichten, und gibt dafür Beispiele: „Nenne zwei Romane von Hermann Hesse“ sei ein DSGVO-konformer Prompt,



Bild: Georg Moltreier

### **Der österreichische Datenschutzaktivist Max Schrems hat mit seiner Organisation noyb Beschwerde gegen Open AI eingereicht.**

auch „Schreibe ein Haiku über Datenschutz“ falle darunter. Die Eingabe von „Entwirf ein Arbeitszeugnis im befriedigenden Bereich für einen Kundenberater in Autohaus X“ hingegen „kann Personenbezug aufweisen, wenn erkennbar ist, aus welchem Unternehmen sie zu welchen Zeitpunkt getätigt wurde“. Das mag so sein, aber wem helfen diese Beispiele in seiner eigenen konkreten Situation wirklich weiter?

## **Grundsatz der Datenrichtigkeit**

An zwei wichtigen Punkten äußert sich die DSK viel schärfer, als es bisher die Landesdatenschutzbehörden taten: bei den Anforderungen an die Transparenz sowie der Umsetzung von Betroffenenrechten. So fordert die DSK in dem Papier, dass der Auftragsverarbeiter den Verantwortlichen dabei unterstützen muss, seinen Informationspflichten nachzukommen. Im Fall von ChatGPT etwa würde das heißen, dass Open AI mit den API-Nutzern kooperieren soll. Doch alles, was das US-Unternehmen tut, ist, ein englischsprachiges „Data processing addendum“ (DPA) bereitzustellen, um seinen Verarbeiterpflichten halbwegs nachzukommen.

In Punkt 11.1. ihrer Orientierungshilfe fordert die DSK von Verantwortlichen, dass „betroffene Personen ihre Rechte auf Berichtigung gemäß Art. 16



DSGVO und Löschung gemäß Art. 17 DSGVO ausüben können“ müssen. Die DSK pocht hier auf den Grundsatz der „Datenrichtigkeit“, den der europäische Gesetzgeber in Art. 5 Abs. 1 DSGVO so formuliert hat: „Personenbezogene Daten müssen sachlich richtig und erforderlichenfalls auf dem neuesten Stand sein; es sind alle angemessenen Maßnahmen zu treffen, damit personenbezogene Daten, die im Hinblick auf die Zwecke ihrer Verarbeitung unrichtig sind, unverzüglich gelöscht oder berichtigt werden.“

Doch wie sollen Verantwortliche sicherstellen, dass ein LLM immer korrekte personenbezogene Daten ausgibt? Man befrage dreimal ChatGPT zu seinem eigenen Geburtstag und der Bot wird oft drei verschiedene Daten ausspucken, weil das Modell zu fabulieren beginnt, wenn es nicht mehr weiter weiß.

Eine Berichtigung müsse „in einer KI-Anwendung umsetzbar sein, zum Beispiel durch Korrektur von Daten oder durch ein Nachtraining/Fine Tuning“, lautet darauf die Antwort im Papier. Genauso müsse man falsche Daten löschen lassen können. Beides geht an der technischen Realität vorbei: Ein LLM ist nun einmal keine Datenbank, in der sich Informationen austauschen lassen. Selbst wenn Open AI GPT-4 auf derlei Hinweise hin nachtrainieren würde – eine Garantie, dass das Modell danach korrekter antwortet, gibt es nicht, allenfalls eine höhere Wahrscheinlichkeit.

Datenschutzbehörde ein Verbot gegen ChatGPT verhängt, unter anderem wegen Kritik an der Datenkorrektheit, dann wohl auch aufgrund des öffentlichen Drucks wieder aufgehoben.

Die Untersuchungen der deutschen Datenschutzbehörden unter dem Dach der DSK laufen derzeit noch. Auch sie wollen von Open AI wissen, wie das Unternehmen die Richtigkeit der verwendeten und ausgegebenen Daten sicherstellen will, insbesondere dann, wenn Betroffene Berichtigungen oder Löschungen einfordern. Es ist folglich nicht auszuschließen, dass auch hierzulande Verbote gegen ChatGPT, Gemini und Co. ausgesprochen und vielleicht auch Verantwortliche in die Pflicht genommen werden.

In dieselbe Kerbe schlug auch der österreichische Datenschutzaktivist Max Schrems mit seiner Non-Profit-Organisation noyb: Bereits im April 2024 hat noyb für eine betroffene Person Beschwerde gegen Open AI bei der österreichischen Datenschutzbehörde eingereicht. Obwohl das von ChatGPT angegebene Geburtsdatum des Beschwerdeführers falsch sei, habe OpenAI seinen Antrag auf Berichtigung oder Löschung abgelehnt, lautet die Begründung. Open AI habe angegeben, dass eine Korrektur der Daten nicht möglich sei.

Es sei damit klar, „dass Unternehmen derzeit nicht in der Lage sind, Chatbots wie ChatGPT mit dem EU-Recht in Einklang zu bringen. Wenn ein System keine genauen und transparenten Ergebnisse liefern kann, darf es nicht zur Erstellung von Personendaten verwendet werden. Die Technologie muss den rechtlichen Anforderungen folgen, nicht umgekehrt“, erklärte Maartje de Graaf, Datenschutzjuristin bei noyb. (hob) **ct**

## Literatur

[1] Joerg Heidrich, Datenschutzcheck für KI-Chatbots, Erste behördliche Handreichung zum Umgang mit Sprach-KI, c't 29/2023, S. 160

## Papiere der

## Datenschutzbehörden

[ct.de/wcge](https://www.ct.de/wcge)

## Verbote nicht ausgeschlossen

An diesem Punkt zeigt sich, dass sich die Funktionsweise generativer KI-Modelle mit der DSGVO massiv beißt. Bereits im April 2023 hatte die italienische

Es gibt **10** Arten von Menschen.  
iX-Leser und die anderen.

Jetzt Mini-Abo testen: 3 digitale Ausgaben + Bluetooth-Tastatur nur **19,35 €**  
[www.ix.de/digital-testen](https://www.ix.de/digital-testen)



[www.ix.de/testen](https://www.ix.de/testen)



0511 / 647 22 888



[leserservice@heise.de](mailto:leserservice@heise.de)







# KI-Überwachung im öffentlichen Raum

Kameras überwachen Parkplätze, sammeln Mautgebühren ein und helfen, Kriminelle zu schnappen. Dazu scannen sie Kennzeichen, Gesichter oder das menschliche Verhalten: nicht immer fehlerfrei, nicht immer transparent.

Von **Andrea Trinkwalder**

**A**n immer mehr öffentlichen und privaten Orten tauchen Überwachungskameras auf. Was genau sie analysieren und ob sie dies im Rahmen der geltenden Gesetze tun, sieht man ihnen von außen nicht an. In Parkhäusern und auf Firmenparkplätzen erfassen schrankenlose Systeme

die Nummernschilder ein- und ausfahrender Pkws, um die Parkzeit zu ermitteln und unberechtigtes Parken zu sanktionieren. Mitunter kontrollieren sie auch per Laufwege-Tracking auf Parkplätzen „Nur für Kunden“, ob der Fahrer tatsächlich in den Laden geht oder das Gelände verlässt.

Häufig müssen dann Datenschützer oder Gerichte entscheiden, ob eine spezifische Technik rechtskonform verwendet wird oder werden kann. Dazu müssen sie mitunter tief in die komplexen Verfahren eintauchen, um die dafür relevanten (technischen) Details zu finden und deren Auswirkungen nachzuvollziehen.

Nahezu jede Anwendung birgt ihre eigenen technischen Herausforderungen: Mal gilt es, eher simple Objekte zu erfassen, die sich aber sehr schnell bewegen; ein Beispiel dafür ist die Nummernschilderkennung im fließenden Verkehr. Andere bewegen sich nur gemächlich oder stehen, was für die Objekterkennung ideal ist. Dies trifft in der Regel auf Personen an öffentlichen Plätzen zu. Dafür sind ihre Gesichter wesentlich schwerer zu unterscheiden als Nummernschilder, weil viel mehr und subtilere Merkmale zu berücksichtigen sind als bei den normierten Autokennzeichen.

## Kennzeichen-Knobeleyen

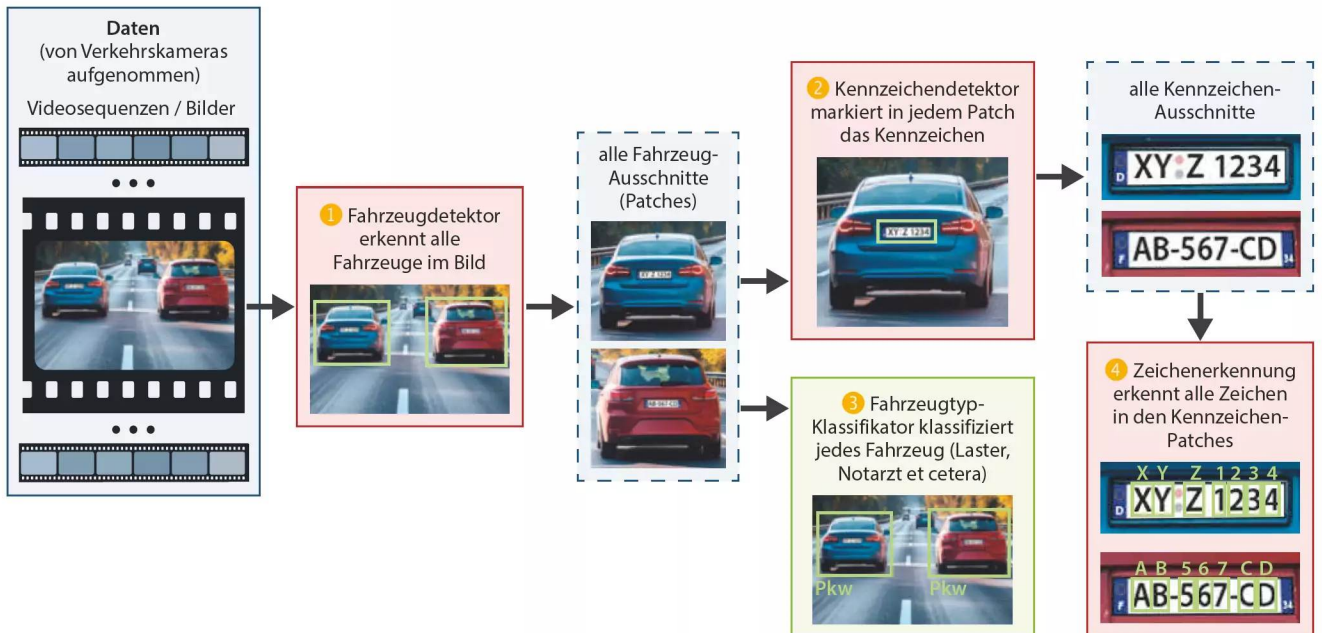
Kfz-Kennzeichen auszulesen, erscheint auf den ersten Blick wie ein längst gelöstes Problem: Die Schilder sind rechteckig, die Schrift kontrastreich, die Glyphen groß, die Anzahl der zu unterscheidenden Klassen gering: 0-9, A-Z. Doch so simpel ist es nicht. Insbesondere Maut- oder Fahndungssysteme müssen die Kennzeichen auch bei hoher Geschwindigkeit und unter widrigen Bedingungen korrekt interpretieren: ungünstige Kamerawinkel, Schnee, Regen, Nebel, Dunkelheit, Gegenlicht, (dadurch bedingte) Reflexionen, Schmutz auf dem Nummernschild et cetera. Wie gut ihnen das gelingt, hängt von der Qualität des aufgezeichneten Bild- und Videomaterials – und damit von der Kamertechnik – sowie von der Robustheit des Algorithmus ab. Nicht zuletzt spielen auch Design, Schriftarten und die verwendeten Ziffern eine Rolle. So

Quelle: Al-Barat et al.

## Funktionsweise eines Echtzeit-Kennzeichenscanners

Maut- und Überwachungssysteme verarbeiten die Videos und Bilder in vier Schritten. Zunächst werden alle Ausschnitte extrahiert, in denen ein Echtzeit-Objektdetektor Fahrzeuge erkannt hat <sup>1</sup>. Anhand dieser

ermitteln weitere Erkennungsalgorithmen die Fahrzeugtypen <sup>3</sup> und extrahieren die Nummernschilder <sup>2</sup>. Eine Zeichenerkennung (OCR) interpretiert die Formen und macht daraus maschinenlesbaren Text <sup>4</sup>.





Quelle: Al-batat et al.

**Damit Kennzeichen-Scanner lernen, mit seltenen und leicht verwechselbaren Glyphen wie „G“ oder „O“ umzugehen, trainiert man sie mit manipulierten Nummernschildern, in denen Originalzeichen zufällig durch solche Glyphen ersetzt wurden.**



Quelle: Al-batat et al.

**Schlechte Witterung, Reflexionen und andere ungünstige Verhältnisse simulieren die Forscher mit Bildbearbeitungsfiltern: Schatten, überstrahlende Farbe und Weichzeichner.**

sind die standardisierten deutschen Schilder mit ihrem starken Schwarz-Weiß-Kontrast für Maschinen grundsätzlich besser lesbar als die mehr oder weniger frei gestaltbaren US-amerikanischen. Bei den Zeichen 0 und O, B und 8, I und 1, S und 5 sowie Z und 2 besteht relativ hohe Verwechslungsgefahr.

Der für deutsche Kfz-Schilder verwendete Font ist die sogenannte fälschungsschwerende Schrift (FE-Schrift). Die charakteristischen Formen der Glyphen sollten Nachbildungen erschweren. Ob sie dies tat-

sächlich gewährleisten, ist umstritten, aber die Zeichen sind darüber hinaus auch auf gute Erkennung und Maschinenlesbarkeit optimiert – ideale Voraussetzungen also für Machine-Learning- und andere Scan-Algorithmen.

Viele Hersteller werben mit fabelhaften Genauigkeiten von 95 Prozent und mehr. Wie gut ihre Systeme im Einsatz tatsächlich funktionieren, ist allerdings nicht transparent und wird in der Regel nicht unabhängig evaluiert. Einige Erkenntnisse kann man immerhin aus wissenschaftlichen Publikationen wie etwa der akribischen Arbeit britischer Forscher unter der Leitung von Reda Al-batat ziehen. Al-batat war Software Engineer an der Universität Westminster und arbeitet mittlerweile beim Überwachungstechnikhersteller Vicom. In ihrem Aufsatz „An End-to-End Automated License Plate Recognition System Using YOLO Based Vehicle and License Plate Detection with Vehicle Classification“ (Link via [ct.de/w1dp](https://ct.de/w1dp)) beschreiben und evaluieren die Wissenschaftler eine selbst entwickelte Echtzeiterkennung, die Fahrzeuge bei hoher Geschwindigkeit klassifiziert und anhand des Nummernschilds identifiziert – auch im Dunkeln und bei schlechter Witterung.

Grundsätzlich besteht die automatische Kennzeichenerkennung (Automatic Number Plate Recognition, ANPR) aus vier Schritten

1. Fahrzeugerkennung
2. Fahrzeugklassifizierung
3. Nummernschilderkennung
4. Optische Zeichenerkennung (OCR)

Ein zentraler Baustein zahlreicher kommerzieller Systeme ist mittlerweile das Open-Source-Framework YOLO (You only look once), das im Jahr 2015 den Durchbruch bei der Echtzeit-Objekterkennung brachte. Dessen Entwickler haben das Klassifizierungsprocedere mit ein paar cleveren Tricks und Kniffen so stark beschleunigt, dass endlich Echtzeitanwendungen realisierbar wurden.

Klassische Algorithmen nutzten ein zweistufiges Verfahren, um alle Objekte im Bild zu lokalisieren und mithilfe von Rahmen (Bounding Boxes) zu markieren: Zunächst trafen sie grobe Vorhersagen darüber, in welchen Bildbereichen relevante Gegenstände oder Personen vorhanden sein könnten. Im zweiten Durchgang wurden diese Vorhersagen verfeinert. YOLO v1 benötigte hingegen nur einen einzigen Durchgang, um potenziell passende Pixelansammlungen aufzuspüren und die Bounding Boxes zu ziehen, was die Berechnungszeit drastisch redu-



# GUTE ARCHITEKTUR SCHAFFT KLARHEIT!

**SPECIAL:  
JETZT BESTELLEN!**

So geht saubere, moderne Softwarearchitektur



**Ein Muss  
für  
Entwickler**



- ▶ Methodisches Vorgehen bei Evolutionärer Architektur
- ▶ Monolithen, Modulithen und Microservices
- ▶ Architektur mit LLMs entwerfen
- ▶ Pragmatische Architekturdokumentation mit Canvas

... und noch viel mehr hilfreiche Artikel

**GLEICH  
REINSCHAUEN!**



**shop.heise.de/ix-softwarearchitektur24**

zierte. Das 2023 veröffentlichte YOLO v8 ist eine sehr universell nutzbare Modellvariante, die nicht nur schnelle, sondern auch sehr akkurate Ergebnisse lieferte. YOLO v10 wiederum wurde auf maximale Geschwindigkeit und Effizienz getrimmt; mit nur knapp halb so vielen Parametern arbeitet es deutlich flotter als seine Vorgänger. Weitere gute Echtzeit-Algorithmen sind Mask-RNN, R-CNN und SSD, mehr dazu siehe [ct.de/w1dp](https://ct.de/w1dp)

Alle Phasen ihrer Echtzeit-Nummernschilderkennung haben die Forscher mit maschinell trainierten Verfahren umgesetzt, drei davon mit einer der YOLO-Versionen: Zunächst spürt YOLO v2 alle Fahrzeuge im Bild auf, anschließend bestimmt ein klassisches Konvolutionsnetzwerk (ResNet50) die Wagentypen. Ein schlankes YOLO v4 tiny extrahiert die zugehörigen Nummernschilder und ein YOLO v4 interpretiert die darauf abgebildeten Zeichen. Doch auch ein YOLO-Detektor ist nur so gut wie die Daten, mit denen er trainiert wurde, wie die Autoren betonen. Während die beiden ersten Phasen kaum Ansprüche an die Trainingsdaten stellen und nahezu fehlerfrei durchlaufen, benötigt man für die beiden letzten sehr vielfältiges, kuratiertes Material, damit das System auch unter realistischen Bedingungen hält, was es im Labor verspricht.

Ein häufiges Problem: Viele Systeme werden mit allzu (länder-)spezifischen Daten trainiert und knifflige Fälle dann über hinzuprogrammierte Regeln und Routinen abgefangen. Leicht verwechselbare Glyphen wie „0“ und „O“ etwa werden in eine gemeinsame Klasse gesteckt und anschließend abhängig von ihrer Position als Buchstabe oder Ziffer interpretiert. Solche länderspezifischen Regeln scheitern spätestens an ausländischen Kennzeichen. Die britischen Wissenschaftler hingegen behandeln jedes Zeichen als eigene Klasse. Damit ihr Mustererkenner lernt, möglichst alle weltweit vorkommenden Schildvarianten korrekt zu interpretieren, nahm die Gruppe rund um Al-batati fünf gängige Datensätze als Basis. Daraus generierten sie außerdem künstliches Trainingsmaterial, indem sie einzelne Glyphen in den Originalschildern durch seltene oder schwer erkennbare Zeichen ersetzten.

Zusätzlich wurde jedes Trainings-Sample dupliziert und mit einer von drei Techniken verfremdet, um natürliche Variationen zu simulieren: Entweder wird ein künstlicher Schatten zufällig über den Ziffern platziert oder eine partielle Färbung angewendet, wie sie auch unter Sonnenlicht auftreten könnte. Als Drittes kommen Unschärfefilter zum Einsatz, um die Effekte von hoher Geschwindigkeit,

ungünstigen Winkeln et cetera nachzuahmen. Insgesamt kommen auf diese Weise knapp 17.000 Samples zusammen.

Das Ergebnis ist eine deutlich robustere Kennzeichenerkennung, die in freier Wildbahn besser durchblickt, auch wenn ihr ein ungewöhnlich gestaltetes Kennzeichen vor die Linse läuft. Im Schnitt soll das System für die gesamte Pipeline eine Genauigkeit (Accuracy) von 90 Prozent erreicht haben. Das bedeutet: 90 von 100 Fahrzeugen hat es korrekt identifiziert. Die zehn Versager enthalten bei diesem Fehlermaß alle Arten von Fehlern: nicht markierte oder falsch klassifizierte Autos, nicht lokalisierte Nummernschilder und eben nicht erkannte oder falsch interpretierte Ziffern. Den größten Anteil haben die drei letzten Fehlertypen.

## Fehler unter Kontrolle

Je besser sich die Umstände kontrollieren lassen, umso besser ist die Erkennungsquote. Ein System mit sehr hochauflösenden Kameras an den überdachten Ein- und Ausfahrten eines Parkhauses kann durchaus mehr als 95 Prozent Genauigkeit erreichen, weil der Betreiber hier eine gute, gleichmäßige Ausleuchtung sicherstellen kann. Elektronischen Mautsystemen (ETC) bescheinigen Experten eine Genauigkeit von bestenfalls 85 bis 90 Prozent, weshalb die Technik hier eher als Ergänzung zu Transpondern gesehen wird, siehe dazu auch den detaillierten, 2007 veröffentlichten Report „Technology Scan for Electronic Toll Collection“, den das Kentucky Transportation Center zusammen mit der University of Kentucky verfasst hat. Alle Studien und Quellen haben wir unter [ct.de/w1dp](https://ct.de/w1dp) verlinkt.

Vor allem für Massenscreenings benötigen die Betreiber deshalb ein effizientes Fehlermanagement, das nicht nur die glasklaren Versager ausleitet und einem manuellen Prüfer vorlegt. Es muss auch in der Lage sein, einen Großteil der unsicheren Kandidaten herauszufischen, bei denen das System möglicherweise Zeichen falsch interpretiert hat. Dazu werden die von einem Maut- oder Parksysteem erfassten Kennzeichen beispielsweise nach Konfidenzlevel sortiert, in dessen Berechnung einfließt, ob der Scanner alle Zeichen klar erkannt oder vielleicht das Zeichen 0 lediglich knapp gegenüber der Ziffer 0 favorisiert hat. Auch externe Faktoren wie etwa die jeweiligen Witterungsbedingungen oder die Bildqualität können in solche Berechnungen einfließen.

Um die manuelle Sichtung zu beschleunigen, muss die Software das aufgenommene Bild sowie



alle erkannten beziehungsweise ermittelten Fahrzeugdetails möglichst klar und übersichtlich aufbereiten. Ein Mensch kann dann rasch vergleichen, ob die ANPR einen Fehler gemacht hat oder nicht, und diesen gegebenenfalls korrigieren, bevor ein falscher Gebührenbescheid herausgeht. Manche Dienstleister agieren hier relativ kulant, indem sie schwer entscheidbare Fälle ad acta legen oder ein funktionierendes Beschwerdemanagement aufbauen. Andere versuchen, durch einen möglichst hohen Automatisierungsgrad ihren Gewinn zu maximieren.

Ob aus einer technisch bedingten Ungenauigkeit ein Problem wird oder nicht, ist also häufig eine Frage des Geschäftsmodells: Viele sogenannte Parkraumbewirtschafter verdienen ihr Geld vor allem über hohe Vertragsstrafen, weil sie den Parkplatzbesitzern ihre Dienste kostenlos anbieten. Mitunter vermieten sie die Parkflächen auch außerhalb der Geschäftszeiten an Anwohner oder verkaufen Datenanalysen für Marketingzwecke.

## Datenschutz-Algorithmen

Das Überwachen von Parkflächen oder eine Fahnung per Kennzeichenscanner ist zwar datenschutzrechtlich heikel, aber nicht per se verboten. Allerdings kommen die Betreiber von Überwachungskameras immer wieder mit dem Gesetz beziehungs-

weise mit den Datenschutzbeauftragten in Konflikt, weil sie zu viele personenbezogene Daten erheben, also gegen das Gebot der Datensparsamkeit verstoßen, oder diese für andere Zwecke verwenden.

Damit solche Begehrlichkeiten gar nicht erst entstehen und um zu verhindern, dass sensible Daten Hackern in die Hände fallen, gibt es eine ganze Reihe von Anonymisierungstechniken, die sich in Echtzeit auf das Videomaterial anwenden lassen. Bei der Parkraumüberwachung beispielsweise müssen Personen und der öffentliche Bereich (angrenzende Straße, Gehweg et cetera) bereits unwiderruflich unkenntlich gemacht werden, bevor die Aufzeichnung gespeichert wird. Dabei kommen ebenfalls Echtzeit-Objekterkennungsalgorithmen zum Einsatz: Sie markieren alle Bildbereiche, die zu Gesichtern gehören, und färben sie schwarz oder verpixeln sie. So landen personenbezogene beziehungsweise biometrische Daten gar nicht erst auf dem Server des Dienstleisters. Auch die Laufwege von Kunden lassen sich so DSGVO-konform verfolgen: Der Algorithmus trackt lediglich die zuvor anonymisierte Person.

Das hessische Parkraumbewirtschaftungsunternehmen Parkvision beispielsweise sah sich dem Vorwurf ausgesetzt, gegen die DSGVO zu verstoßen. Im Raum stand vor allem der Verdacht, dass es zur Laufwegeerkennung auch den öffentlichen Bereich filmt und personenbezogene Daten nicht ausrei-

Sogenannte Parkraumbewirtschafter wie Parkvision bieten Geschäftsleuten oft kostenlos an, ihre Parkplätze zu überwachen. Dabei scannen sie nicht nur Nummernschilder, sondern kontrollieren auch per Laufwegetracking, ob eine Person tatsächlich den Laden betritt.

The screenshot shows the Parkvision website with a navigation bar at the top containing links: Home, Parkraumlösungen, Bewachungen, Vermietung, Datenschutz, Geschäftskunden, and Falschparker. The main content area features a large graphic of a building labeled 'PLAZA' with a red 'Alarm!' speech bubble. Below this, the text reads: 'IHR KOSTENFREIER START MIT DIGITALER PARKRAUMAUF SICHT UND KUNDENIDENTIFIZIERUNG' followed by a button 'JETZT MEHR ERFAHREN'. At the bottom, there are four service tiles: 'DIGITALE AUF SICHT' (with a camera icon), 'DIGITALE SCHRANKEN' (with a barrier icon), 'HYBRID LÖSUNG' (with a building and camera icon), and 'BAUSTELLEN BEWACHUNG' (with a construction site icon).



chend anonymisiert, woraufhin der Hessische Beauftragte für Datenschutz und Informationsfreiheit (HBDI) aktiv wurde. Parkvision spielte auf Zeit und hielt den HBDI so lange mit dürren Informationen hin, bis dieser drohte, den Weiterbetrieb des Systems zu untersagen. Zum Zeitpunkt der dazu stattfindenden Anhörung hatte Parkvision die Technik allerdings nachgebessert. Der HBDI hält die Anlage nun für DSGVO-konform, nannte aber gegenüber c't mit Verweis auf das Geschäfts- und Betriebsgeheimnis keine Details zur eingesetzten Technik oder den Fehlerraten. Von einer Untersagung des Weiterbetriebs sah die Behörde ab, leitete aber ein Ordnungswidrigkeitsverfahren ein, weil Parkvision nicht angemessen kooperiert habe.

Bei Mautsystemen wiederum fallen naturgemäß sehr detaillierte Ortsdaten an, insbesondere wenn der Staat die Gebühren besonders „fair“ oder nach ökologischen Gesichtspunkten erheben möchte: also nicht nur pauschal für einen gewissen Zeitraum oder nach Strecke, sondern zusätzlich nach Schadstoffausstoß, Größe, Tageszeit et cetera. Die vielen einzelnen Datenpunkte verdichten sich dann häufig zu einem umfangreichen Bewegungsprofil auf dem Server des Betreibers, und nicht jeder geht damit verantwortungsvoll um. Im September 2021 musste etwa der norwegische Mautbetreiber Ferde 5 Millionen NOK, also knapp 500.000 Euro Strafe zahlen, weil er Bilder von Zahlungskarten und Nummernschildern zur Verarbeitung an einen chinesischen Dienstleister übertragen hatte.

Ferde betreibt die Mautstationen in den Regionen Vestland, Rogaland und Agder. Ihr System fotografiert die Kennzeichen von Autos, die keinen Prepaid-Sticker an der Windschutzscheibe haben. 2020 hatte der norwegische Fernsehsender NRK herausgefunden, dass die Kennzeichen, die das System nicht automatisiert verarbeiten konnte, an einen chinesischen Dienstleister geschickt wurden, der sie manuell abglich: Zwischen 2017 und 2019 seien 2,5 Millionen Fotos jährlich dorthin übertragen worden. Informationen darüber, wer sich wann wo aufhielt, außer Landes zu geben, stuft die norwegische Datenschutzbehörde als illegal ein.

Ferde kündigte den Vertrag mit dem chinesischen Dienstleister, zahlte die Strafe – und legte die hohen Ausgaben auf die norwegischen Autofahrer um, indem es die Mautgebühren anhob. Somit schließt sich also der Kreis und man landet ganz schnell wieder bei der Frage, welche Tracking-Systeme sich unter welchen Umständen überhaupt wirtschaftlich, sicher und effizient betreiben lassen.

## Zu ungenau, zu unsicher?

Grundsätzliche Bedenken gegenüber der auch in Großbritannien weit verbreiteten ANPR-Technik äußerte sogar der im Oktober 2023 aus dem Amt geschiedene Terrorismus- und Security-Experte Fraser Sampson, der dem Vereinigten Königreich zwei Jahre lang als Commissioner für Biometrie und Überwachungskameras diente. In seinem offenen Abschiedsbrief richtete er einen dringlichen Appell an seinen Nachfolger, die in seinen Augen viel zu ungenauen sowie fehler- und betrugsanfälligen ANPR-Systeme durch eine modernere und verlässlichere Technik zu ersetzen.

Er kritisierte, dass die ANPR allen technischen Fortschritten zum Trotz nach wie vor von einem Stück Plastik abhängt, das leicht zu manipulieren oder gar zu fälschen sei. Einer neueren Schätzung zufolge verwende einer von 15 Fahrern in Großbritannien Anti-ANPR-Technik wie etwa Stealth-Nummernschilder oder reflektierendes Klebeband, um die Schrift unlesbar zu machen. Zugleich werde die ANPR ein immer wichtigeres Instrument für die Strafverfolgung und die strategische Verkehrsplanung und -steuerung.

Die Genauigkeit der aktuellen Systeme hält Sampson in diesen Anwendungsszenarien für nicht akzeptabel. Er rechnete vor, dass „aus der Nutzung dieser Überwachungstechnologie die größte nicht-militärische Datenbank in Großbritannien entstanden ist“, mit ungefähr 15.400 von Kameras überwachten Verkehrsspuren. Diese lieferten derzeit 75 bis 80 Millionen Lesevorgänge pro Tag. „Selbst wenn die Systeme mit einer Genauigkeit von 97 Prozent arbeiten würden, bedeutet dies, dass sie täglich immer noch 2,4 Millionen fehlerhafte Ergebnisse produzieren.“ Mit der Verarbeitung derart ungenauer Daten setze der Staat die unbescholtenen Bürger einem zu hohen Risiko aus und breche möglicherweise auch Datenschutzgesetze, so Sampsons Fazit.

In Deutschland gibt es (noch) kein feingranulares elektronisches Pkw-Mautsystem wie in Norwegen oder Großbritannien, geschweige denn eine mit UK vergleichbare Echtzeit-Videoüberwachung, die ja auch Gesichter und ungewöhnliches beziehungsweise verdächtiges Verhalten zu identifizieren verspricht und mit wissenschaftlich zweifelhaften Verfahren wie der Emotionserkennung experimentiert.

Die hierzulande für Lkws installierten Anlagen erfassen zwar Pkw-Kennzeichen, dürfen aber nicht für die polizeiliche Fahndung zweckentfremdet wer-

## Tarnen und tricksen

Im Internet kursieren allerhand Tipps und Gadgets, wie man sich dem Kennzeichen-Tracking entziehen kann. Dazu gehören reflektierende Klebestreifen oder LED-Leuchstreifen, die strategisch günstig auf dem Schild platziert werden. Angepriesen wird auch Tarnkappenfolie, die das gesamte Schild überdeckt und für Infrarotkameras unsichtbar machen soll. Ob solche Mittel tatsächlich wie erhofft wirken, ist fraglich. Denn bessere Systeme kombinieren mehrere Kameras und Sensoren, etwa eine Bild- mit einer Infrarotkamera. In den meisten Ländern und auch in Deutschland sind Praktiken, die das Ablesen erschweren, außerdem verboten.

Laut § 10 der Fahrzeug-Zulassungsverordnung (FZV) darf am Kennzeichen nichts angebracht oder verändert werden, was zu Verwechslun-

gen führen könnte. Das Kennzeichen muss jederzeit ohne Schwierigkeiten ablesbar sein. Wer dies verhindert, muss mit Konsequenzen rechnen:

- Das Überkleben oder Verändern des Kennzeichens kann als Kennzeichenmissbrauch gewertet werden. Gemäß § 22 Absatz 3 des Straßenverkehrsgesetzes (StVG) droht dafür sogar eine Geld- oder Freiheitsstrafe bis zu einem Jahr.
- Selbst wenn keine Straftat vorliegt, wird das Überkleben des Kennzeichens als Ordnungswidrigkeit geahndet. Es droht ein Bußgeld von bis zu 65 Euro.
- Durch das Überkleben oder Verändern kann das Kennzeichen seine Gültigkeit verlieren, was zu weiteren rechtlichen Konsequenzen führen kann.

den. Dennoch gibt es immer wieder Versuche, solche Systeme zu installieren: den Kennzeichenabgleich „Kesy“ zur Gefahrenabwehr bei schweren und schwersten Straftaten in Brandenburg etwa, „Section Control“ zur Geschwindigkeitsmessung in Niedersachsen oder die 2019 zur Schleierfahndung in Bayern und Baden-Württemberg eingesetzten mobilen und stationären Scanner.

Besonders weit geht der Freistaat Sachsen mit seinem „Personen-Identifikationssystem“ (PerIS), das es eigens von der Firma Opto Precision entwickeln ließ, um an der deutsch-polnischen Grenze Eigentumsdelikte zu verfolgen. Dazu werden Kennzeichen gescannt und mit zur Fahndung ausgeschrieben Fahrzeugen verglichen. Auch Fußgänger erfasste die Polizei per Echtzeit-Gesichtserkennung. Um das System betreiben zu können, hatte die Landesregierung eigens den Paragraph 59 des sächsischen Polizeivollzugsdienstgesetzes geändert. Deshalb darf sie seit 2019 schwerer grenzüberschreitender Kriminalität mit dem „Einsatz technischer Mittel“ begegnen.

Die bis zum 31. Dezember 2023 befristete Regelung sollte nach drei Jahren vom Forschungsinstitut für öffentliche Verwaltung in Speyer evaluiert wer-

den. Anhand dieser Untersuchung zog die Landesregierung das Fazit, dass der „technische und personelle Aufwand sehr groß“ gewesen sei, und verlängerte die Gesetzesnorm nicht weiter. Gleichzeitig wurden Mitte Dezember 2023 weitere stationäre Scanner installiert.

Auch die in Bayern und Baden-Württemberg für die Schleierfahndung eingesetzten und später teilweise verbotenen Kennzeichenscanner produzierten offenbar unverhältnismäßig viele Fehler und hohe Kosten: Selbst unter den Treffern fanden sich nur wenige Fälle, die mit schwerer oder organisierter Kriminalität in Verbindung standen. Dieses Bild zeichnete das Newsportal BuzzFeed anhand mehrerer kleiner Anfragen, die Abgeordnete an die jeweiligen Landesregierungen gestellt hatten.

Eine unabhängige Evaluierung von KI-Systemen findet bislang meist nicht statt, und das wirft die grundsätzliche Frage auf, wie die immer wieder geforderte und beschworene Transparenz solcher Algorithmen denn in der Praxis überhaupt erreicht werden soll. Eine Orientierung, was Betreiber von Überwachungskameras filmen, auswerten und sammeln dürfen und was nicht, gibt unser Rechtsartikel [1].

(atr) **ct**

### Literatur

[1] Joerg Heidrich, Rechtlicher Sichtschutz, Videoüberwachung: Was erlaubt ist und was nicht, c't 17/24, S. 90

### Studien und Quellen

[ct.de/w1dp](https://www.ct.de/w1dp)



# EU-Regeln für KI

Mit dem AI Act will die EU Nutzer von großen Sprachmodellen und Bildgeneratoren vor Risiken schützen. Im Vordergrund stehen Offenlegungs- und Sicherheitspflichten der Anbieter, die Mitte 2025 greifen sollen. Schon heute besteht Rechtsunklarheit wegen geltender EU-Datenschutzvorschriften.

Von **Holger Bleich**

**N**iemand hat die kalifornische Ideologie so gut auf den Punkt gebracht wie Facebook-Gründer Mark Zuckerberg: Sein 2012 zum Börsengang formuliertes Motto „move fast and break things“ beschreibt rückblickend treffend die Geschichte des Weltkonzerns Meta. Im Silicon Valley und in San Francisco behauptet man oft, besser als Regierungen und Gerichte zu wissen, was gut für die Menschen ist; und maximiert nebenbei den eige-

nen Profit. Facebook stieg in den Nullerjahren auch deshalb zum marktbeherrschenden sozialen Netzwerk auf, weil das Unternehmen jede Menge Regeln gebrochen hat, beispielsweise Datenschutzvorschriften. Viele dieser Verstöße blieben ohne negative Folgen für das Unternehmen.

Im Jahr 2022 ereignete sich ein ähnlicher kalifornischer Moment wie einstmals die Vorstellung des iPhones 2007: Das Unternehmen OpenAI in San

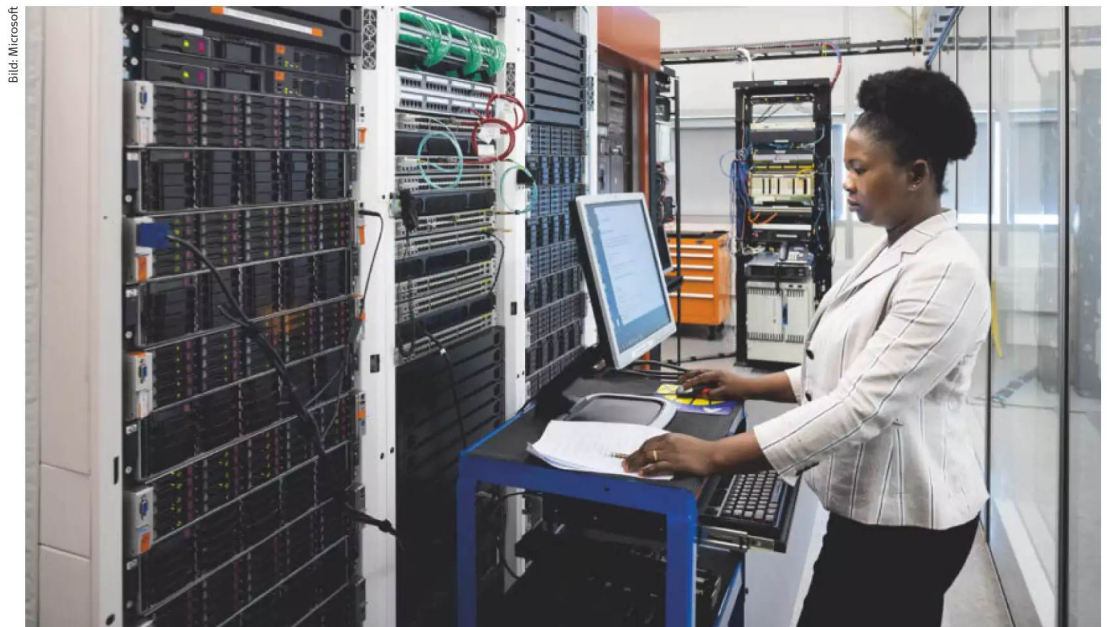


Francisco gab die Webversion seines KI-Bots ChatGPT frei, und die Welt staunte ob der Fähigkeiten des zugrundeliegenden Sprachmodells GPT. Das Wettrennen um die leistungsfähigsten Modelle war eröffnet. Google und Meta sprangen umgehend auf den KI-Hypetrain auf, während Skeptiker erste Fragen nach Ethik und Datenschutz stellten.

Diesmal, so viel scheint klar, will die europäische Politik beim Tempo der potenziellen Kaputtmacher mithalten und ihnen Grenzen setzen, bevor allzu viel zu Bruch geht. Doch wie soll das funktionieren, ohne Innovationen abzuwürgen? Eine erste Antwort auf diese Frage glaubt die EU mit dem AI Act (manchmal auch KI-Verordnung oder einfach KI-VO genannt) gefunden zu haben, dem weltweit ersten Gesetz zur Regulierung von Entwicklern und Anbietern künstlicher Intelligenz. Man stelle mit der Verordnung sicher, „dass KI-Systeme die Grundrechte, die Sicherheit und die ethischen Grundsätze achten“, behindere aber nicht „die Investitionen und die Innovation im Bereich KI in der gesamten EU“. So stellt es die EU-Kommission dar.

Das klingt nach einem durchdachten Konzept. Dabei hatte es sich 2022 einfach nur gut getroffen, dass die EU-Kommission ohnehin längst eine Verordnung für risikobasierte Regulierung von Entscheidungsalgorithmen in der Mache hatte, die sie etwas großspurig KI-Gesetz betitelte. Weil mit den großen Sprachmodellen und den KI-Bildgeneratoren aber etwas gänzlich Neues auftauchte, das sich keinen konkreten Einsatzzwecken zuordnen ließ, musste dieses KI-Gesetz überarbeitet und ergänzt werden.

Herausgekommen ist ein eigenes Kapitel in der KI-Verordnung mit dem Namen „KI-Modelle mit allgemeinem Verwendungszweck“ (Kapitel V, ab Art. 51 KI-VO), auf Englisch bezeichnet als „General Purpose AI“ kurz GPAI. In diesem Kapitel definiert die EU, was sie unter GPAI versteht und welche Pflichten Entwickler, Anbieter und Anwender von GPAI, insbesondere von KI-Sprachmodellen, bald erfüllen müssen. Diese Regeln gelten ein Jahr, nachdem die KI-Verordnung in Kraft getreten ist. Weil dies am 1. August des vergangenen Jahres geschah, greifen sie ab dem 2. August 2025. Allzu viel Zeit bleibt also nicht.



**Microsoft stellt unter anderem in einem Rechenzentrum bei Iowa enorm viel Azure-Cloud-Infrastruktur bereit, auf der OpenAI seine GPT-Sprachmodelle trainieren kann. Die Kapazität dürfte EU-Schwellwerte für systemische Risiken deutlich überschreiten.**

## Konkrete Risikofaktoren

Die gesamte KI-Verordnung zu erläutern würde den Rahmen dieses Artikels sprengen. Eine Einordnung lesen Sie in c't 11/2024 [1]. Nur so viel: Diese Verordnung gilt seit August 2024 unmittelbar als Gesetz in jedem EU-Mitgliedsstaat (den Volltext finden Sie unter [ct.de/wr3a](https://eur-lex.europa.eu/eli/reg/2024/4181/oj)). Sie unterteilt KI-Systeme nach vielen Kriterien in vier Risikoklassen, abgestuft von „verboten“ bis „unbedenklich“, weshalb die EU von einem „risikobasierten Ansatz“ spricht. Je grundrechtsgefährdender die Anwendung eines KI-Systems ist, desto heftigere Pflichten treffen den Anbieter. KI, die die EU-Kommission als verboten deklariert, musste schon im Februar 2025 vom Markt verschwinden.

Die Einordnung in Risikoklassen hat der EU-Gesetzgeber auch bei den zuletzt angeflanschten Artikeln zu generativer KI mit allgemeinem Verwendungszweck übernommen, also auch für große Sprachmodelle. Allerdings beschränkte er sich hier auf zwei Stufen: Die Standardmodelle sowie jene GPAI, von der ein „systemisches Risiko“ ausgehen könnte. Was meint der Gesetzgeber damit? In Art. 51, Abs. 1a heißt es, GPAI mit systemischem Risiko verfüge „über Fähigkeiten mit hohem Wirkungsgrad, die mithilfe geeigneter technischer Instrumente und Methoden, einschließlich Indikatoren und Benchmarks, bewertet werden“.

Sehr konkret wird es dann in Art. 51 Abs. 2: „Bei einem KI-Modell mit allgemeinem Verwendungszweck wird angenommen, dass es über Fähigkeiten mit hohem Wirkungsgrad gemäß Absatz 1a verfügt, wenn die kumulierte Menge der für sein Training verwendeten Berechnungen, gemessen in Gleitkommaoperationen, mehr als  $10^{25}$  beträgt.“ Der Gesetzgeber zieht also die Grenze bei einer konkreten FLOP-Anzahl der fürs Training benötigten Rechenleistung, bis ein KI-Modell in den produktiven Einsatz geht.

An dieser Schwelle wurde bereits viel Kritik geübt, weil sie willkürlich scheint. Derzeit überschreiten wohl nur die großen US-Anbieter wie OpenAI und Google diese Linie, und damit bleiben nachzeitigem Stand Anbieter aus der EU vorerst in der sanfteren Reglementierung. Verbirgt sich dahinter ein europäischer Protektionismus? Immerhin behält sich die EU in Art. 51 Abs. 3 der KI-Verordnung vor, die Grenze mit einem sogenannten delegierten Rechtsakt jederzeit nach oben verschieben zu können, beispielsweise, wenn sich die „Hardwareeffizienz erhöht“.

Doch ohnehin gilt: Keine Regel ohne Ausnahme. Die EU-Anbieter sollten sich nicht zu sicher vor harten Pflichten fühlen, denn Art. 51, Abs. 1b definiert den Vorbehalt, dass die EU-Kommission „von Amts wegen oder aufgrund einer qualifizierten Warnung des wissenschaftlichen Gremiums“ auch außerhalb der in Abs. 2 genannten quantitativen Kriterien einem generativen KI-Modell systemische Risiken attestieren kann. Im Anhang XIII der KI-Verordnung sind Kriterien genannt, die dabei eine Rolle spielen könnten, etwa die Anzahl der gewerblichen Nutzer eines Modells.

## Transparenzpflichten

In den auf Art. 51 folgenden Artikeln führt die KI-Verordnung all die Pflichten auf, die Anbieter von allgemeinen KI-Modellen ab August 2025 erfüllen müssen. In erster Linie geht es der EU augenscheinlich um Information und Transparenz. Anbieter sollen qua Gesetz gezwungen werden, sich mit der Entstehung und Wirkung ihrer KI-Modelle auseinanderzusetzen. Gewerbliche und private Anwender möchte man in die Lage versetzen, Informationen zu eben jenen Faktoren zu erhalten und damit eine informierte Entscheidung zu treffen, ob sie das Modell nutzen möchten.

An dieser Stelle sei erwähnt, wie diffizil, aber wichtig die Begrifflichkeiten in der KI-Verordnung sind. Sie unterscheidet beispielsweise zwischen KI-Modellen und KI-Systemen, in denen die Modelle zur Anwendung kommen. Man findet diese Begrifflichkeiten nicht im Gesetzestext, sondern in den nicht rechtsverbindlichen Erwägungsgründen zur Verordnung. Dort liest man in Grund 101 beispielsweise den zum Verständnis der Verordnung wichtigen, länglichen Satz: „Anbieter von KI-Modellen mit allgemeinem Verwendungszweck nehmen entlang der KI-Wertschöpfungskette eine besondere Rolle und Verantwortung wahr, da die von ihnen bereitgestellten Modelle die Grundlage für eine Reihe nachgelagerter Systeme bilden können, die häufig von nachgelagerten Anbietern bereitgestellt werden [...]“.

Daraus folgt wohl: „Anbieter“ von Modellen stehen am Beginn der Verwertungskette. Anbieter können aber auch Wiederverkäufer der Modelle sein, also beispielsweise Mailprovider, die ein Sprachmodell via API in ihr Frontend einbauen und den Kunden zur Verfügung stellen. Andererseits spricht die KI-VO von „Anwendern“, die auch gewerblich sein können. Das trifft etwa Unternehmen, die einen Chatbot via API ihren Mitarbeitern bereitstellen.



Spricht man mit EU-Beamten, beschleicht einen bisweilen der Eindruck, dass auch bei Diskussionen ums Gesetz nicht jeder immer wusste, was der andere gerade meinte.

Art. 53 KI-VO definiert die technischen Dokumentationspflichten für alle GPAI-Anbieter und verweist für Konkretes auf Anhang XI der Verordnung. Demnach müssen sie eine allgemeine Beschreibung ihres KI-Modells bereithalten, beispielsweise vorgesehene Einsatzzwecke, „Architektur und Anzahl der Parameter“ und „Modalität“ (etwa Sprachmodell oder Bildgenerator). Außerdem soll eine „ausführliche Beschreibung der Elemente“ enthalten sein. Hier geht es beispielsweise um das Trainingsverfahren, die Herkunft der dazu verwendeten Daten und deren Aufbereitung sowie um den „bekannten oder geschätzten Energieverbrauch“. Art. 53 Abs. 1c fordert außerdem von GPAI-Anbietern eine „Strategie zur Einhaltung des Urheberrechts der Union und damit zusammenhängender Rechte“.

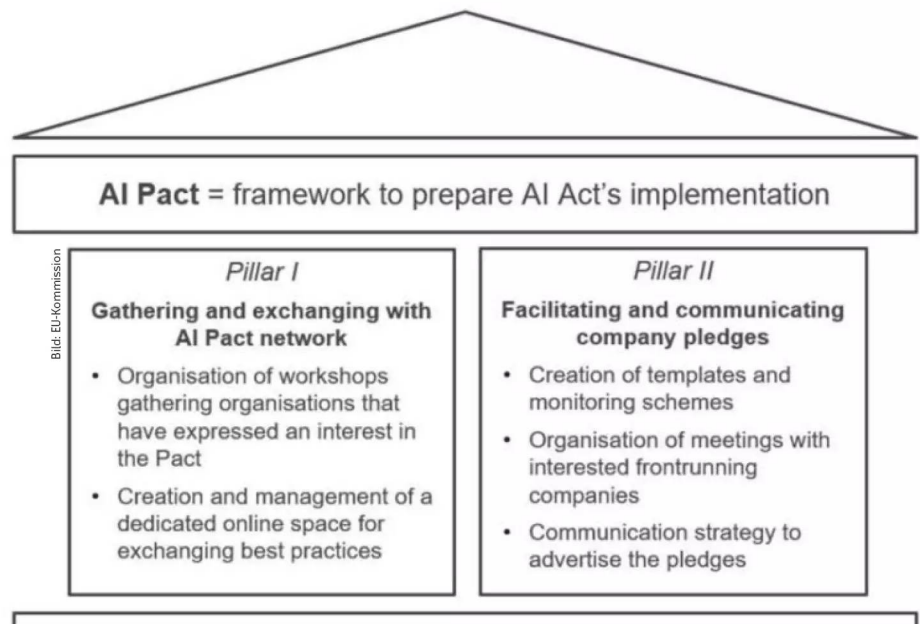
All diese Transparenzpflichten gelten gemäß Art. 53 Abs. 2 nicht für Anbieter von Open-Source-Modellen. Ob beispielsweise die quelloffenen Sprachmodelle von Meta und dem französischen Anbieter

Mistral in diese Kategorie fallen, muss sich noch zeigen. Ohnehin gilt die Ausnahme nicht, wenn das Modell systemische Risiken birgt, was nach der KI-VO-Definition zumindest bei den großen freien Modellen durchaus der Fall sein könnte.

Für alle GPAI-Anbieter gilt laut KI-VO: Sie müssen für die zuständige Aufsichtsbehörde jederzeit erreichbar sein. Nach Art. 53 sollen Anbieter aus Nicht-EU-Staaten deshalb einen Bevollmächtigten in der EU benennen, der mit weitreichenden Kompetenzen ausgestattet ist. Diese Person soll beispielsweise jederzeit dazu bereitstehen, angeforderte Dokumentationen zu übergeben.

## Systemische Risiken

Ob sein KI-Modell die Schwelle zum systemischen Risiko nach den oben genannten Kriterien überschreitet, soll ein Anbieter permanent selbst bewerten. Ist dies der Fall, muss er es Art. 53 Abs. 1 KI-VO zufolge innerhalb von zwei Wochen der EU-Kommission melden. Die Kommission darf aber auch selbst Prüfungen anstoßen und KI-Modelle als risikant deklarieren.



**Ein Pakt zwischen GPAI-Anbietern und EU-Kommission soll auf zwei Säulen ruhen und gewährleisten, dass die Anbieter wissen, wie sie ihre neuen Pflichten erfüllen können.**



Gemäß Art. 53 KI-VO gelten für KI-Modelle mit systemischen Risiken zusätzliche, verschärfte Transparenzpflichten, die man im zweiten Abschnitt des Anhang XI nachlesen kann. Darin geht es um die Prüfung des Modells auf Fehler und Risiken sowie um den Nachweis der Sicherheit, etwa mit Angriffstests. Außerdem verlangt die KI-VO eine ausführliche Beschreibung der Systemarchitektur, einschließlich des Zusammenspiels der verwendeten Software-Komponenten. Die Anbieter müssen künftig also viele Formulare ausfüllen und Dinge preisgeben, die sie vor der Konkurrenz oder den Rechteinhabern ihrer Trainingsinhalte gerne geheim halten würden.

Spannend wird es insbesondere in Art. 55 der Verordnung. Hier verlangt der Gesetzgeber von Anbietern zu GPAI-Modellen mit systemischen Risiken, diese zu benennen und Vorgehenspläne vorzulegen, die sie beseitigen können. „Einschlägige Informationen über schwerwiegende Vorfälle und mögliche Abhilfemaßnahmen“ sind zu erfassen, zu dokumentieren und unverzüglich der zuständigen Aufsichtsbehörde zu melden. An dieser Stelle lässt die Datenschutz-Grundverordnung (DSGVO) schön grüßen.

## Code of Practice

Summa summarum verlangt die KI-Verordnung den Anbietern von Sprachmodellen ab, sich Gedanken über deren Wirkung auf die Gesellschaft und die Grundrechte zu machen, eine Ethik zu entwickeln und diese zu dokumentieren. Außerdem sollen sich nachgelagerte Anbieter ein Bild vom Modell machen können und dies wiederum dokumentieren (siehe KI-VO Anhang XII). All die Transparenzpflichten lassen Kritiker der KI-Verordnung von einem Compliance-Monstrum sprechen. Der Gesetzgeber hält dagegen, er habe in der Verordnung Erleichterungen angelegt.

Dies führt zur Frage, wer die Aufsicht führt. Neben anderen Gremien steht vor allem das bei der EU-Kommission angesiedelte KI-Amt (AI Office) im Fokus. Dort soll die Kompetenz versammelt sein, und dort sollen die Meldungen der Modellanbieter einlaufen. Außerdem sollen noch zu benennende nationale Aufsichtsbehörden mit dem KI-Amt kooperieren und ihm zuarbeiten. Nach derzeitigem Stand wird dies in Deutschland die Bundesnetzagentur übernehmen.

Damit die Modellanbieter Einfluss auf die technische Seite der neuen Compliance-Standards nach KI-VO nehmen können, hat der ehemalige EU-Binnenmarktkommissar Thierry Breton einen „KI-Pakt“ initiiert. In diesem Rahmen sollen die Unternehmen

bis Mai 2025 selbst „Best Practices“ erarbeiten. Der Pakt wurde Ende September 2024 in Abwesenheit von Breton von über hundert erstunterzeichnenden Unternehmen besiegelt, darunter etwa OpenAI, Microsoft, Amazon, Palantir, Aleph Alpha und die Deutsche Telekom. Nicht dabei waren zunächst jedoch die genannten Firmen Meta und Mistral.

Parallel zu dieser Aktion hat die EU-Kommission dazu aufgerufen, Vorschläge und Wünsche für einen „Code of Practice“ für das GPAI-Kapitel einzureichen, der die Transparenzanforderungen konkretisieren soll. Dieses Verfahren ist genau so in der KI-Verordnung vorgesehen. Gemäß Art. 56 Abs. 9 KI-VO muss dieser Code of Practice bis zum 2. Mai 2025 veröffentlicht sein. Die Kommission hat angekündigt, dass sie diesen Rahmen auszuschöpfen gedenkt. Nach Veröffentlichung bleiben den GPAI-Anbietern folglich drei Monate Zeit, um die Anleitungen bis zum 2. August 2025 umzusetzen.

Tun sie es unzureichend oder gar nicht, drohen drakonische Strafen gemäß Art. 71 KI-VO. Bei Mängeln in Transparenz, Risikomanagement oder Sicherheit können nationale Aufsichtsbehörden oder das KI-Büro demnach Geldbußen von bis zu 20 Millionen Euro oder bis zu 4 Prozent des weltweiten Vorjahresumsatzes verhängen. Bei unvollständigen Angaben gegenüber dem KI-Büro könnten immerhin bis zu 10 Millionen Euro Strafe fällig werden. Bußgelder dürfen die Behörden aber nicht nur gegen die Anbieter selbst aussprechen, sondern auch gegen nachgelagerte Anbieter, Händler oder gewerbliche Nutzer.

All diese Adressaten tun also gut daran, sich mit der neuen EU-Regulierung von KI-Systemen zur allgemeinen Verwendung zu befassen. Das gilt insbesondere für alle, die mit den großen Sprachmodellen hantieren, die wir im vorherigen Artikel vorgestellt haben. Zusätzlich zur KI-Verordnung, die erst kommt, stehen derzeit aber auch ungelöste Konflikte mit dem Datenschutzrecht im Blickpunkt.

## Risikofaktor Datenschutz

Zur Rechtsunsicherheit tragen sicherlich die gegenläufigen Äußerungen hiesiger Datenschutzbehörden bei. Die Deutsche Datenschutzkonferenz (DSK) als gemeinsames Gremium der Behörden etwa hat im Mai 2024 eine (wenig konkrete) Orientierungshilfe dazu veröffentlicht, wie Unternehmen Sprachmodelle datenschutzkonform einsetzen könnten [2].

Zwei Monate später, im Juli 2024, widersprach aber bereits der als pragmatisch bekannte Hambur-

**Steile These: Laut dem Hamburgischen Beauftragten für Datenschutz und Informationsfreiheit Thomas Fuchs ist die Speicherung von personenbezogenen Daten in Sprachmodellen kein Fall für den Datenschutz.**

Bild: Daniel Bockwoldt/dpa



## Literatur

[1] Holger Bleich, Aufgeweichtes Regelwerk, Was die neue KI-Verordnung der EU für KI-Anbieter und Anwender ändert, c't 11/24, S. 140

[2] Holger Bleich, KI-Chatbots unter Druck, Warum sich Datenschutzbehörden mit ChatGPT & Co. so schwertun, c't 13/24, S. 154

## KI-Verordnung und weitere Infos

[ct.de/wr3a](https://ct.de/wr3a)

gische Beauftragte für Datenschutz und Informationsfreiheit Thomas Fuchs indirekt den Prämissen der DSK. Unter dem Titel „Diskussionspapier: Large Language Models und personenbezogene Daten“ veröffentlichte er seine These, dass die DSGVO auf die Speicherung von Daten in Sprachmodellen gar nicht anwendbar sei. Deshalb wirke sich ein „gegebenenfalls datenschutzwidriges Training nicht auf die Rechtmäßigkeit des Einsatzes eines solchen Modells in einem KI-System aus“ (siehe [ct.de/wr3a](https://ct.de/wr3a)).

Momentan ist völlig unklar, welche Prozesse beim Aufbau und Betrieb datenschutzrechtlich zu behandeln sind. Nichtsdestoweniger laufen in Europa erste Verfahren gegen die Anbieter großer Sprachmodelle. Während es OpenAI mit der italienischen Datenschutzbehörde zu tun bekommen hat, griff die irische Datenschutzaufsichtsbehörde DPC Google an. Sie geht in einem Mitte September 2024 eröffneten Verfahren dem Verdacht nach, dass Google vor dem Training seines Pathways Language Model 2 (PaLM 2) keine oder nur eine unzureichende Datenschutz-Folgeabschätzung vorgenommen hat. Google hat PaLM 2 im Mai 2023 offiziell vorgestellt, heute verwendet der Alphabet-Konzern das Modell in verschiedenen Anwendungen, etwa in Gmail.

Spannend wird nun zu beobachten sein, ob die DPC weitere Verfahren folgen lässt. Die Behörde ist

für jene US-Unternehmen federführend zuständig, die ihre europäische Hauptniederlassung in Irland haben. Das betrifft einen Großteil der Techunternehmen wie Meta, Apple, X und eben auch Alphabets Google. Sowohl Meta mit seinem Llama-Modell als auch X, das das von X-AI betriebene Grok-LLM zugeführt bekommt, stoppten nach Intervention der DPC vorerst das Training ihrer Modelle mit DSGVO-geschützten Daten.

Einer repräsentativen Erhebung des Branchenverbands Bitkom im Spätsommer 2024 zufolge befürchten 57 Prozent von 605 befragten Unternehmen, der Datenschutz Sorge dafür, dass die Anwendung von KI in der EU eingeschränkt werde. 52 Prozent gehen sogar davon aus, der Datenschutz vertreibe Unternehmen aus der EU, die KI entwickeln. Ein Grund dafür (laut 50 Prozent der Unternehmen): Der Datenschutz erschwere es, KI-Modelle mit genügend Daten trainieren zu können.

In Hinblick auf den Einsatz von KI ist für den Bitkom die Lage prekär. Er verbindet seine Bestandsaufnahme mit einer deutlichen Forderung: „Wir brauchen beim Datenschutz Einfachheit und Klarheit. Der Datenschutz hat tiefgreifende Auswirkungen auf Unternehmen ebenso wie auf die Gesellschaft, deshalb muss er verständlich und praxistauglich gemacht werden“, betont Susanne Dehmel, Mitglied der Bitkom-Geschäftsleitung. (hob) **ct**

# Interview: Wie der AI-Act entstand

**Gabriele Mazzini hat als Teamleiter die europäische KI-Verordnung maßgeblich gestaltet und erwarb sich mit seinem umfangreichen juristischen und technischen Sachverstand auch den Respekt von Regulierungskritikern. Im Gespräch mit c't (geführt auf der TED-AI-Konferenz in Wien) erläutert er, wie akribisch geforscht, wie bedacht abgewogen werden muss – und wie er mit Lobbyismus umgeht.**

Von **Andrea Trinkwalder**

**c't: Sieben Jahre lang haben Sie am AI Act gearbeitet: eine fachlich anspruchsvolle Thematik. Haben Sie einen technischen Hintergrund?**

**Gabriele Mazzini:** Nein, ich bin Jurist. Ich war immer als Anwalt tätig, aber nicht in einer Anwaltskanzlei. Nach dem Jurastudium habe ich promoviert und dann für die EU gearbeitet. Zunächst war ich Rechtslinguist und habe im Wesentlichen Dokumente für den Gerichtshof übersetzt. Dann habe ich als Anwalt für das Europäische Parlament gearbeitet. Anschließend bin ich in die USA gegangen und habe als Unternehmensjurist für eine NGO gearbeitet, die sich mit internationaler Entwicklung beschäftigt. Außerdem war ich für zwei Start-ups tätig. Als ich zurück nach Brüssel ging, begann ich mit der politischen Arbeit zu KI. Mein Hintergrund war also immer eine Mischung aus Recht und Politik, sowohl innerhalb als auch außerhalb der Institutionen. Ich glaube, dass dies in gewisser Hinsicht auch der Knackpunkt war, denn ich denke nicht wie ein Bürokrat. Was ich zu sagen versuche: Mein Hintergrund ist ein wenig gemischt.

**c't: Und wie kam es, dass Ihnen die KI-Gesetzgebung übertragen wurde?**

**Mazzini:** Ich war, wie gesagt, sehr an Technologie interessiert. So begann ich zunächst in einer Kommission zu arbeiten, die sich mit den haftungsrechtlichen Auswirkungen von KI beschäftigte. Wie Sie

vielleicht schon gehört haben, gibt es jetzt eine Richtlinie zur KI-Haftung, die noch aussteht. Als mir aber klar wurde, dass ich nicht hinter dieser Initiative stehe, verließ ich sie. Ich dachte: Okay, mein Leben mit KI ist vorbei, und war kurz davor, in die EU-Delegation nach Kabul zu wechseln.

Aber dann hatte ich das Gefühl, dass ich vielleicht noch einen Moment warten sollte, denn es war eine Zeit des Übergangs von der Juncker-Kommission zur Präsidentschaft von Ursula von der Leyen. Und als sie 2019 ins Amt kam, hatte sie in ihren programmatischen Leitlinien tatsächlich einen Rechtsrahmen für KI. Damit hatten die Kommissionsmitglieder nicht gerechnet.

Sie fragten mich, ob ich Interesse an diesem Projekt hätte, weil ich einer der wenigen, wenn nicht sogar der Einzige in der Kommission war, der sich bereits intensiv mit der Schnittstelle zwischen KI und neuem Recht beschäftigt hatte. Ich hatte nämlich schon während der ersten zwei Jahre, in denen ich in der Kommission war, über die KI-Thematik nachgedacht, ohne für eine Regulierung zu plädieren. Ich war der Ansicht, dass wir zuerst intensiv studieren und lernen müssen, denn KI-Governance ist komplex.

Aber irgendwie hat von der Leyens Ansinnen, eine Gesetzgebung vorzuschlagen, diesen Prozess abgekürzt. Sie sagte sogar: innerhalb von hundert Tagen, wenn ich im Amt bin. Das war unrealistisch, aber es gab bei mir schließlich den Ausschlag. In gewisser Weise war ich bereit, und im Grunde gaben sie mir einen Freibrief.



**c't: KI vernünftig zu regulieren, ist in der Tat ein sehr komplexes Thema. Wie ist es Ihnen gelungen, herauszufinden, was für die Gesetzgebung wirklich wichtig ist?**

**Mazzini:** Zuerst habe ich versucht, zu verstehen, ob es eine Lücke gibt. Und so habe ich versucht, herauszufinden: Okay, wir haben möglicherweise dieses Problem, aber ist dieses Problem bereits durch einen anderen rechtlichen Rahmen gelöst? Daraufhin begann ich, Wege zu finden, wie wir den AI Act so gestalten könnten, dass wir einerseits dieses Mandat einhalten, uns aber andererseits nur auf die Dinge konzentrieren, die nicht reguliert sind. Und das war die Idee, die letztendlich zu diesem Ansatz der Produktgesetzgebung führte.

**c't: Es gibt sehr viele unterschiedliche Interessengruppen und Akteure. Vor allem als die Idee aufkam, die Grundlagenmodelle gesondert zu regeln, schien der Lobbyismus zuzunehmen. Wie groß ist der Druck? Wie schafft man es, sich nicht manipulieren zu lassen?**

**Mazzini:** Lobbyarbeit ist normal. Ich habe das Gefühl, dass sie manchmal einen negativen Beigeschmack hat. Die Leute reden über Lobbyisten als etwas Schlechtes. Aber es gibt Unternehmen, die Lobbyarbeit betreiben, es gibt zivilgesellschaftliche Organisationen, die Lobbyarbeit betreiben, es gibt Regierungen, die Lobbyarbeit betreiben, jeder betreibt Lobbyarbeit. Denn jeder hat seine eigenen Interes-

sen. In gewisser Weise ist Lobbyarbeit meiner Meinung nach also in Ordnung, sie ist Teil des Spiels. Das macht den Job eigentlich erst interessant, weil man mit so vielen Leuten spricht und jeder eine andere Meinung hat.

Die Kunst besteht darin, dass man als politischer Entscheidungsträger in der Lage ist, zu verstehen und zu filtern, was einem gesagt wird. Denn wenn ich, sagen wir mal, mit Google spreche und sie schicken mir einen Ingenieur, der mir erklärt, wie ein LLM funktioniert: Dann sollte ich besser auch einen Ingenieur von der Kommission mitbringen, der überprüfen kann, ob die Informationen, die von Google kommen, tatsächlich die vollständige Darstellung sind oder vielleicht nur eine Möglichkeit, die Thematik zu betrachten. Denn ich weiß, wenn ich mit Google spreche, werden sie mir nur sagen, was in ihrem Interesse ist.

Hier kommt es also auf die Reife des politischen Entscheidungsträgers an. Und man braucht Kompetenz. Man muss in der Lage sein, zu verstehen, woher die Informationen kommen, ob sie verlässlich sind, und man muss in der Lage sein, sie zu filtern. In diesem Sinne sind die guten Lobbyisten diejenigen, die tatsächlich mit Informationen an den Tisch kommen. Nicht diejenigen, die versuchen, Ihnen zu sagen, Sie sollten dies tun, Sie sollten das tun. Aber um noch einmal auf Ihre Frage zurückzukommen, ob es Druck von Interessengruppen gab. Ich habe das ganze KI-Gesetz entwickelt, ohne dass mir jemand gesagt hat, wie ich es machen soll. Das ist die Realität.

**c't: Wie gehen Sie vor, um zwischen ernsthaften und interessengeleiteten Bedenken zu entscheiden, etwa bei der häufig geäußerten Befürchtung: "Dies oder jenes wird Europa daran hindern, KI zu entwickeln."? Wie können Sie einschätzen, ob es eine wirklich realistische Sorge ist oder nur der Versuch, möglichst viel für die eigene Branche herauszuschlagen?**

**Mazzini:** Zuerst studiert man, macht also seine eigene Due-Diligence-Prüfung, dann lässt man sich beraten. Man hat also eine Idee und fängt an, sich mit Experten zu beraten. Aus diesem Grund arbeitet die Kommission regelmäßig mit externen Experten zusammen, aber selbst die Experten haben manchmal unterschiedliche Ansichten. Ich denke also, dass es auch keine exakte Wissenschaft ist, auf jeden Fall ist es immer eine Meinung, oder? Werden wir mehr oder weniger wachsen, wenn wir dies oder jenes tun?

Man muss aber zumindest ein ordentliches Verfahren haben, um eine ordentliche Due-Diligence-

**Der promovierte Jurist Gabriele Mazzini ist federführender Autor des AI Act, der ab Februar regelt, wie künstliche Intelligenz in der EU eingesetzt werden darf. Die nachträglich ergänzten Passagen zur Regulierung von Grundlagenmodellen – also Sprach- und anderen generativen KIs – empfindet er als Fremdkörper.**

Bild: Robert Leslie / TED AI Vienna



Prüfung, ordentliche Entwürfe und dann eine ordentliche Konsultation durchzuführen, damit man irgendwann sagen kann: Es gibt jetzt einen ausreichenden Konsens und ein ausreichendes Verständnis für die Fragen, und dann wollen wir uns auf diese konzentrieren und nicht auf andere. Es ist also eine Kombination aus einerseits technischem Fachwissen und Kenntnissen und andererseits dem, was aus der Außenwelt kommt: von den Menschen, die tatsächlich mit KI arbeiten, KI anwenden oder von KI betroffen sein könnten. Das ist der Prozess. Und wie gesagt, nicht perfekt, aber es ist ein Prozess, bei dem man, wenn man ihn gut macht und zuhört, wahrscheinlich die meisten Dinge richtig macht.

**c't: Der AI Act war fast fertig, als mit dem Erfolg von ChatGPT die Diskussion über die Grundlagenmodelle und speziell die großen Sprachmodelle aufflammte. In puncto Regulierung schien alles wieder von vorn loszugehen. Gab es einen Punkt, an dem Sie dachten, dass der AI Act nie zu einem Ende kommen würde?**

**Mazzini:** Nein, ich hatte zu dem Zeitpunkt vielmehr das Gefühl, dass genau diese Diskussion mehr Zeit braucht, und hatte eher Bedenken, zu früh fertig zu werden. Denn diese Debatte hat gezeigt, wie komplex die Sache ist. Und wenn eine Sache komplex ist, dann muss man das meiner Meinung nach auch anerkennen. Es ist schließlich keine Atombombe, die zu explodieren droht.

ChatGPT hat viele Menschen überrascht und sicherlich einige Probleme verursacht. In Italien wurde ChatGPT beispielsweise von der Datenschutzbehörde zwei Wochen lang gesperrt. Aber sagen wir mal so: Das alles ist kein Grund, ein Gesetz zu erlassen, bevor man überhaupt verstanden hat, welche Auswirkungen die LLMs beziehungsweise Grundlagen- oder multimodale Modelle auf die Wirtschaft haben würden. In meinen Augen war daher vor allem problematisch, dass in gewisser Weise dieses Gefühl herrschte: Wir müssen es rasch angehen, gemäß dem politischen Zeitplan, weil wir sonst scheitern werden. Ich war der gegenteiligen Meinung, nämlich dass es wichtiger ist, sich zuerst gut aufzustellen, weil wir eindeutig noch nicht so weit waren.

**c't: Welche wichtigen Prozesse fehlten in Bezug auf die Regulierung von Grundlagenmodellen, beziehungsweise wurden abgekürzt?**

**Mazzini:** Es gab noch keine Folgenabschätzung. Das ist eine Studie, die von der Kommission durchgeführt wird, bevor sie einen Vorschlag veröffentlicht. Im Wesentlichen handelt es sich dabei um eine Analyse der potenziellen Auswirkungen der Verordnung, die Sie entwerfen wollen. In der Regel müssen wir dann eine Prüfung durch einen Ausschuss für Regulierungskontrolle bestehen, der diese Studie, diese Folgenabschätzung, analysiert und beurteilt, ob sie gut genug ist oder nicht. Dieses Instrument ist zwar nicht perfekt, aber es ist ein gutes Instrument.

Es liefert also eine Begründung, warum man [die Kommission, Anm. d. Red.] sich für dieses und nicht für jenes entschieden hat. Daher wird in der Regel jeder Vorschlag der Kommission von den beiden Mitgesetzgebern irgendwie unterstützt. Während des rechtlichen Verfahrens oder des Gesetzgebungsverfahrens ändern die beiden Mitgesetzgeber den Vorschlag also normalerweise nicht vollständig. Sie passen ihn nur an, verfeinern ihn. Der endgültige Text folgt im Allgemeinen der Struktur, den grundlegenden Ideen des Kommissionsvorschlags – mit Ausnahme dessen, was bei den LLMs geschah.

**c't: Könnten Sie diesen Unterschied bitte näher erläutern?**

**Mazzini:** Der endgültige Text des AI Act enthält ein neues Kapitel: ein völlig neues Kapitel über General Purpose AI, das im Kommissionsvorschlag nicht enthalten war. Und dieses Kapitel weicht auch vom ursprünglichen Ansatz des Kommissionsvorschlags ab, der darauf abzielt, KI in Abhängigkeit von Anwendungsfällen zu regulieren. Dieses neue Kapitel hingegen reguliert Basismodelle: Ja, das ist auch eine Änderung der Logik.

Dieses neue Kapitel beeinflusste auch das Governance-System [also den festzulegenden Ordnungsrahmen mit seinen Regeln und Prozessen, Anm. d. Red.]. Dafür hatten wir aber noch keine Evidenz, außer einige Gespräche und einige Papiere mit einigen Interessenvertretern, aber keine Studie, keine interne Analyse. Das war es, was mir ein wenig Unbehagen bereitete.

**c't: Es wäre in Ihren Augen also besser gewesen, die ursprünglichen Texte und Kapitel des AI Act so zu überarbeiten, dass automatisch auch die generativen KIs adäquat reguliert werden, anstatt diese Spielart gesondert zu behandeln?**



Wie gut Unternehmen die Anforderungen bewältigen können, hängt laut Mazzini von den Leitlinien ab, erste sind bereits erarbeitet. Eine große Bedeutung kommt etwa der KI-Kompetenz von Mitarbeitern zu. Das der Dualen Hochschule Baden-Württemberg zugehörige Projekt AI Comp hat schon ein Modell mit den verschiedenen für KI erforderlichen Kompetenzfeldern erstellt.

Bild: AI Comp/DHBW

## Ethische Kompetenz

Ethisches Bewusstsein für KI-bezogene Fragen

### Lernzielbeschreibungen

#### Stufe 1-2: Wissen & Verstehen – Anwenden & Analysieren:

Die Lernenden sollen in der Lage sein, ethische Dilemmata, die mit KI-Anwendungen verbunden sind, zu erkennen und zu beschreiben.

#### Stufe 3-4: Anwenden & Analysieren – Bewerten & Entwerfen:

Die Lernenden sollen fähig sein, die ethischen Auswirkungen verschiedener KI-Anwendungen zu analysieren und zu bewerten, um verantwortungsbewusste Entscheidungen zu treffen.

#### Stufe 5-6: Bewerten & Entwerfen – Erschaffen:

Die Lernenden sollen in der Lage sein, eigenständige ethische Richtlinien für den Einsatz von KI zu entwickeln und zu implementieren, die sowohl persönliche als auch organisationale Verantwortung reflektieren.



Ethische Kompetenz bezeichnet die Fähigkeit, ethisch mit KI-relevanten Themen umzugehen und verantwortungsvoll zu handeln.

**Mazzini:** Nun, das ist einer der Einwände, die vorgebracht wurden, dass dieses neue Kapitel im Wesentlichen eine Technologie an sich reguliert. Es regelt also im Wesentlichen die Transformer, die Grundlagenmodelle, während das KI-Gesetz als Konzept neutral ist und künstliche Intelligenz unabhängig von einer bestimmten Architektur regelt. Es spricht von KI-Systemen, und dann muss man natürlich definieren, was ein KI-System ist. Dies schließt im Grunde alles maschinelle Lernen ein, richtig? Dazu gehören auch Expertensysteme, im Grunde ist das Konzept des maschinellen Lernens sehr umfassend.

Im Gegensatz dazu wird es jetzt spezifische Bestimmungen für eine bestimmte Art des maschinellen Lernens geben, nämlich die Transformer. Und das ist eine Frage, die sich meiner Meinung nach stellt: Warum unterscheidet sich diese spezielle Technologie von anderen? Eigentlich glaube ich nicht, dass es der richtige Ansatz ist, eine bestimmte Technologie zu regulieren. Denn schon morgen wird es eine andere geben.

**c't:** Welche kritischen Regulierungslücken wurden denn im ursprünglichen Entwurf identifiziert, die es erforderlich machten, Grundlagenmodelle gesondert zu behandeln?

**Mazzini:** Es gibt eine Frage, die diese neuen Werkzeuge, also generative KI, gegenüber dem ursprünglichen Kommissionsvorschlag aufwerfen. Der Kommissionsvorschlag konzentrierte sich auf Anwendungsfälle: also auf einen engen Anwendungsbereich und damit typischerweise auf KI-Systeme, die auf diesen engen Bereich trainiert wurden.

Eine generative KI wie etwa einen Chatbot hingegen können Sie bitten, eine Reihe von Aktivitäten durchzuführen; in mehreren, einschließlich Hochrisikobereichen. Die eigentliche Frage, die ich mir damals gestellt habe und die durch das KI-Gesetz nicht wirklich gelöst wird: Wie können wir sicherstellen, dass der risikobasierte Ansatz greift? Dass wir regulieren, sobald der Einsatz von KI ein hohes Risiko darstellt, zum Beispiel wenn Firmen eine KI beim Recruiting von Menschen einsetzen. Dann wollen wir sicherstellen, dass das System überprüft und kontrolliert wird. Wie erreichen wir also, dass dieses Szenario tatsächlich als hohes Risiko angesehen wird, wenn es als solches im Kommissionsvorschlag definiert ist? Und wie stellen wir sicher, dass auch generative KI abgedeckt ist, wenn ich sie benutze, um dasselbe Ziel zu erreichen?

Das ist die Komplexität, denn generative KI kann viele, viele verschiedene Aufgaben außer dem Recruiting erledigen und wird auf eine andere Weise



trainiert als das traditionelle KI-System. Wenn man aber die [im Kommissionsvorschlag definierten] Anforderungen für ein hohes Risiko auf generative KI anwendet, dann funktionieren sie als solche nicht wirklich. Das ist meiner Meinung nach ein Problem, das immer noch nicht gelöst ist.

**c't:** Für die praktische Umsetzung wurden und werden Leitlinien und Verhaltenskodizes erarbeitet, um es Unternehmen zu erleichtern, sich an die Vorgaben des KI-Gesetzes zu halten. Wie stark beeinflussen die verbliebenen Unklarheiten diesen für die Praxis wichtigen Prozess?

**Mazzini:** Diese Aufregung um generative KI und ChatGPT führte dazu, dass die Erzählung über KI in Richtung von Risiken verschoben wurde, die sehr undefiniert, systemisch und sogar existenziell sind. Außerdem wurde zusätzliche Komplexität in einen rechtlichen Rahmen gebracht, der ohnehin bereits sehr komplex war.

Das Hinzufügen dieses neuen Kapitels bedeutete, dass die Kommission jetzt Leute einstellt, die nur an der Umsetzung dieses völlig neuen Kapitels arbeiten. Sie sollen versuchen, den Sinn dieser Bestimmungen zu verstehen, die im Moment ohnehin sehr vage sind.

Viele Details werden also gemeinsam mit der Industrie ausgearbeitet. Das lenkt einen großen Teil der Aufmerksamkeit auf dieses Thema. Dadurch werden viele Ressourcen von der tatsächlichen Umsetzung anderer Teile des KI-Gesetzes abgelenkt, die für uns relevanter sind, nämlich das hohe Risiko. Im Wesentlichen haben wir also Regeln geschaffen, die wir meiner Meinung nach nicht brauchen. Oder zumindest wussten wir nicht genau, wie wir sie umsetzen sollten, weil wir keine Analyse hatten. Und jetzt erzwingen diese Regeln eine Art von Aktivität, die die Umsetzung anderer Teile des Gesetzes komplexer und schwieriger machen wird. In gewisser Weise ist das nicht gut für die europäische Industrie.

**c't:** Ein häufiger Kritikpunkt ist, dass das Ausmaß an Nachweisen bezüglich Transparenz, Trainingsdaten et cetera vor allem für kleinere Firmen nur schwer zu bewältigen sein dürfte. Worauf kommt es jetzt bei der Entwicklung der Leitlinien an?

**Mazzini:** Das alles muss in der Tat erst definiert werden. Jetzt geht es ans Eingemachte, denn im Mo-

ment ist das KI-Gesetz noch sehr substanzlos. Wir brauchen also die Standards, wir brauchen die Rechtsakte, die Durchführungsrechtsakte. Ohne diese wird es für die Unternehmen schwierig sein, die Vorschriften einzuhalten. Aber zumindest sollten sie damit in der Lage sein, zu wissen, was sie tun sollen.

Vor allem kleinere und mittlere Unternehmen (KMU) sowie Start-ups haben weniger Ressourcen. Ein Mangel an Klarheit darüber, wie die Vorschriften einzuhalten sind, wird sich vor allem auf die kleineren Akteure auswirken. Im Wesentlichen wird der Maßstab für die Einhaltung durch den Stand der Technik gegeben sein. Wenn Sie ein KMU sind, wird es für Sie viel schwieriger sein, Nachforschungen anzustellen, die Literatur zu sichten, um zu verstehen, was der Stand der Technik für eine technische Lösung ist. Die großen Unternehmen werden weniger Probleme haben.

**c't:** Sie hatten viel Gestaltungsfreiheit, aber immer unter dem Mandat, ein KI-Gesetz auf den Weg zu bringen. Wie würden Sie KI regulieren, wenn Sie vollkommen frei entscheiden könnten?

**Mazzini:** Wenn ich heute das Gleiche mit voller Freiheit tun könnte, würde ich kein KI-Gesetz machen, sondern auf meine ursprüngliche Idee zurückkommen. Wir bräuchten jemanden im Büro des Kommissionspräsidenten, der sich mit KI auskennt und der den Abteilungen sagen kann: Okay, wir haben mit der KI vielleicht ein paar Probleme bei den Finanzdienstleistungen. Dann verfassen Sie dafür einige Anleitungen. Und wenn Sie Probleme mit der Regulierung medizinischer Geräte haben, können Sie dort etwas tun. Aber diese Leute müssten koordiniert werden. Es müsste also jemanden geben, der einen 360-Grad-Blick hat und diese Hebel in Bewegung setzen kann.

Aber diese Möglichkeit gibt es nicht, wir haben jetzt das KI-Gesetz. Ich habe also versucht, das KI-Gesetz so zu gestalten, dass wir zumindest diese Spannungen mit bestehenden Gesetzen vermeiden. Aber das ist mir letztendlich nicht gelungen. Mein ursprünglicher Entwurf war nicht ideal und verbesserungsbedürftig, aber ich denke, er war relativ in Ordnung. Der endgültige Text hingegen macht in mancher Hinsicht diese Spannungen zwischen anderen Rechtsrahmen meiner Meinung nach potenziell komplizierter. (atr) **ct**

**KI-Verordnung, Zeitplan  
und Informationen:**

[ct.de/w1zq](https://ct.de/w1zq)

# IMPRESSUM

## Redaktion

Postfach 61 04 07, 30604 Hannover  
Karl-Wiechert-Allee 10, 30625 Hannover  
Telefon: 05 11/53 52-300  
Telefax: 05 11/53 52-417  
Internet: [www.heise.de](http://www.heise.de)

**Leserbriefe und Fragen zum Heft:**  
[sonderhefte@ct.de](mailto:sonderhefte@ct.de)

Die E-Mail-Adressen der Redakteure haben die Form [xx@heise.de](mailto:xx@heise.de) oder [xxx@heise.de](mailto:xxx@heise.de). Setzen Sie statt „xx“ oder „xxx“ bitte das Redakteurs-Kürzel ein. Die Kürzel finden Sie am Ende der Artikel und hier im Impressum.

**Chefredakteur:** Torsten Beeck (tbe, verantwortlich für den Textteil), Dr. Volker Zota (vza)

**Konzeption:** Hartmut Gieselmann (hag)

**Koordination:** Jobst Kehrnhahn (keh, Leitung), Pia Groß (pia), Tom Leon Zacharek (tlz)

**Redaktion:** Jo Bager (jo), Holger Bleich (hob), Ronald Eikenberg (rei), Greta Friedrich (gref), Hartmut Gieselmann (hag), Arne Grävmeyer (agr), Nico Jurrán (nij), André Kramer (akr), Andrea Trinkwalder (atr), Dr. Sabrina Patsch (spa), Dorothee Wiegand (dwi), Christof Windeck (cwi), Christian Wölbert (cwo), Dušan Živadinović (dz)

**Mitarbeiter dieser Ausgabe:** Thomas Brandstetter, Isabel Grünewald, Joerg Heidrich, Falk Steiner

**Assistenz:** Susanne Cölle (suc), Tim Rittmeier (tir), Martin Triadan (mat)

**DTP-Produktion:** Vanessa Bahr, Anne-Marie Berndt, Dörte Bluhm, Lara Bögner, Beatrix Dedek, Madlen Grunert, Cathrin Kapell, Steffi Martens, Leonie Preuß, Lisa Reich, Marei Stade, Matthias Timm, Christiane Tümmeler, Nicole Wesche

**Digitale Produktion:** Christine Kreye (Leitung), Thomas Kaltschmidt, Martin Kreft, Pascal Wissner

**Fotografie:** Andreas Wodrich, Melissa Ramson

**Illustrationen, Aufmacherfotos:** Rudolf A. Blaha, Albert Hulm, Thorsten Hübner, Sebastian Steppuhn

**Titel:** Steffi Martens, [www.freeipik.com](http://www.freeipik.com)

## Verlag

Heise Medien GmbH & Co. KG  
Postfach 61 04 07, 30604 Hannover  
Karl-Wiechert-Allee 10, 30625 Hannover  
Telefon: 05 11/53 52-0  
Telefax: 05 11/53 52-129  
Internet: [www.heise.de](http://www.heise.de)

**Herausgeber:** Christian Heise, Ansgar Heise, Christian Persson

**Geschäftsführer:** Ansgar Heise, Beate Gerold

**Mitglieder der Geschäftsleitung:** Jörg Mühle, Falko Ossmann

**Anzeigenleitung:** Michael Hanke (-167)  
(verantwortlich für den Anzeigenteil),  
[www.heise.de/mediadaten/ct](http://www.heise.de/mediadaten/ct)

**Anzeigenverkauf:** Verlagsbüro ID GmbH & Co. KG,  
Tel.: 05 11/61 65 95-0, [www.verlagsbuero-id.de](http://www.verlagsbuero-id.de)

**Leiter Vertrieb und Marketing:** André Lux (-299)

**Service Sonderdrucke:** Julia Conrades (-156)

**Druck:** Firmengruppe APPL Druck GmbH & Co. KG,  
Senefelder Str. 3-11, 86650 Wemding

**Vertrieb Einzelverkauf:**  
DMV DER MEDIENVERTRIEB GmbH & Co. KG  
Meßberg 1  
20086 Hamburg  
Tel.: 040/3019 1800, Fax: 040/3019 145 1815  
E-Mail: [info@dermedienvertrieb.de](mailto:info@dermedienvertrieb.de)  
Internet: [dermedienvertrieb.de](http://dermedienvertrieb.de)

**Einzelpreis:** € 14,90; Schweiz CHF 27,90;  
Österreich € 16,40; Luxemburg € 17,10

**Erstverkaufstag:** 28.02.2025

Eine Haftung für die Richtigkeit der Veröffentlichungen kann trotz sorgfältiger Prüfung durch die Redaktion vom Herausgeber nicht übernommen werden. Kein Teil dieser Publikation darf ohne ausdrückliche schriftliche Genehmigung des Verlages in irgendeiner Form reproduziert oder unter Verwendung elektronischer Systeme verarbeitet, vervielfältigt oder verbreitet werden. Die Nutzung der Programme, Schaltpläne und gedruckten Schaltungen ist nur zum Zweck der Fortbildung und zum persönlichen Gebrauch des Lesers gestattet.

Für unverlangt eingesandte Manuskripte kann keine Haftung übernommen werden. Mit Übergabe der Manuskripte und Bilder an die Redaktion erteilt der Verfasser dem Verlag das Exklusivrecht zur Veröffentlichung. Honorierte Arbeiten gehen in das Verfügungsrecht des Verlages über. Sämtliche Veröffentlichungen erfolgen ohne Berücksichtigung eines eventuellen Patentschutzes.

Warennamen werden ohne Gewährleistung einer freien Verwendung benutzt.

Hergestellt und produziert mit Xpublisher:  
[www.xpublisher.com](http://www.xpublisher.com)

Printed in Germany.

Alle Rechte vorbehalten.

© Copyright 2025 by  
Heise Medien GmbH & Co. KG

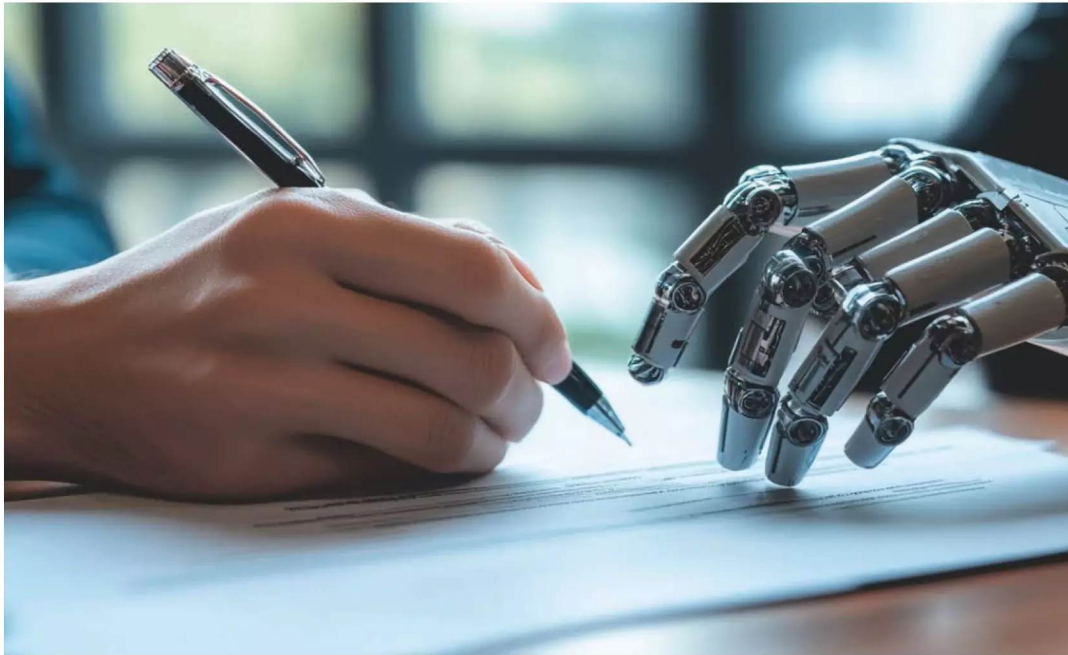


Bild: Midjourney / Joerg Heidrich

# Neue Verträge für KI-generierte Werke

KI-generierte Inhalte im kreativen Bereich erfordern angepasste Nutzungsverträge. Sie führen zu der absurden Situation, dass ein unbeteiligter Dritter solche Inhalte faktisch einfacher nutzen kann als ein zahlender Erwerber.

Von **Joerg Heidrich**

**E**gal ob Stockfotos, Texte oder Videos: Wer digitale Inhalte erwirbt oder verkauft, schließt über diesen Vorgang einen Vertrag ab. Dieser regelt im Kern die Erlaubnis, die urheberrechtlich geschützten Inhalte für einen bestimmten Zweck zu nutzen. Weil im europäischen Recht das Urheberrecht unveräußerlich ist, werden nur die Nutzungsrechte übertragen und Lizenzen erteilt. Diese können unterschiedliche Schwerpunkte haben.

So wechseln bei einem Autorenvertrag üblicherweise viele Rechte zum Vertragsnehmer. Neben dem eigentlichen „Copyright“, also dem Nutzungsrecht, vereinbaren die Vertragspartner meist auch Rechte zu Bearbeitung, Übersetzung oder Weitergabe. Zu Stockfotos regelt der Vertrag, in welchem Bereich die Bilder oder Grafiken eingesetzt werden dürfen – etwa nur online, für Werbezwecke oder auch auf dem Cover eines Printmagazins.



Bei KI-generierten Inhalten, seien es Bilder, Texte, Grafiken, Musik oder Filme, stellt sich die rechtliche Situation anders dar. Denn das klassische Urheberrecht fußt auf der Vorstellung eines menschlichen Schöpfers, der eine originell-kreative Leistung erbringt. Ziel dieses Rechts ist es, die persönliche geistige Leistung umfassend zu schützen und dem Schöpfer eine wirtschaftliche Verwertung zu ermöglichen.

Worin liegt nun die geistige Leistung desjenigen, der ChatGPT oder Midjourney anwendet? Er denkt sich einen mehr oder weniger originellen Prompt aus. Danach übernimmt die Maschine und erzeugt automatisch das Werk. Das Ergebnis lässt sich niemals genau voraussehen. Bei generativer KI kann man die menschliche Leistung eher mit der eines Mäzens vergleichen, der einem Maler die Grundlagen eines Bilds vorgibt. Die eigentliche „kreative“ Leistung entspringt zweifellos dem KI-Modell.

## Kein Schutz ohne Vertrag

Urheberrechtsexperten haben sich daher schnell geeinigt: Die Ergebnisse generativer KI sind in aller Regel nicht urheberrechtlich geschützt. Es mangelt an einer persönlichen geistigen Schöpfung. Uneingeschränkt gilt dies allerdings nur für Inhalte, die direkt aus der KI gefallen sind und nicht erheblich umgearbeitet wurden.

Bearbeitet man etwa ein KI-generiertes Bild aufwendig nach oder schreibt einen Text nahezu komplett um, kann dies durchaus einen eigenen urheberrechtlichen Schutz begründen. Ob das tatsächlich der Fall ist, hängt wie so oft vom Einzelfall ab und dürfte zunehmend die Gerichte beschäftigen.

Sind die Ergebnisse generativer KI nicht urheberrechtlich geschützt, können sie auch nicht über die gängigen Verträge weitergegeben werden, wie sie bislang im Rechtshandel üblich sind. Diese Verträge funktionieren in aller Regel rechtlich nicht bei KI-generierten Inhalten. Erstaunlicherweise haben bislang nur wenige Rechtsverwalter Anpassungen vorgenommen. So finden sich auch bei vielen größeren Stockfoto-Agenturen zu dieser Problematik allenfalls sehr oberflächliche Hinweise.

Die Anwendung traditioneller Vertragsformen kann sogar unangenehme rechtliche Konsequenzen nach sich ziehen. Im Kern übertragen sie Nutzungsrechte an Inhalten gegen Entgelt. Bestehen diese Rechte aber gar nicht, so wird der Vertrag durch die Übergabe nicht geschützter Inhalte nicht erfüllt. Deshalb kann sich der KI-Nutzer unangenehmen

Regressforderungen seiner Vertragspartner gegenübersehen. Im schlimmsten Fall könnte ein Gericht diesen Fall sogar als Betrug im Sinne des Paragraphen 263 Strafgesetzbuch (StGB) werten, wenn nämlich der KI-Nutzer wider besseres Wissen seine Vertragspartner darüber täuscht, dass er ihnen tatsächlich nicht vorhandene Nutzungsrechte überträgt.

## Vertragsanforderungen

Trotz dieser problematischen Situation lassen sich Ergebnisse von Claude, Midjourney oder ChatGPT durchaus rechtskonform und kostenpflichtig übertragen. Alles andere wäre beispielsweise für Inhalteanbieter wie Foto- oder Werbeagenturen, die für ihre Kunden exklusive Inhalte produzieren, auch katastrophal. Denn dort ist es mittlerweile üblich, Auftragsarbeiten auch mithilfe von KI zu erstellen oder zumindest anzupassen. Sowohl derlei Unternehmen als auch kreativ schaffende Einzelkämpfer sollten ihre bestehenden und neuen Verträge kritisch überprüfen und anpassen. Dabei müssen sie das Verhältnis der beiden Vertragsparteien und die übertragenen Leistungen neu denken.

Wie können also Nutzungsverträge aussehen, die KI-generierte Inhalte einbeziehen? Eine wichtige Rolle spielt die Präambel, also das Vorwort eines Vertrags. Die Lizenzgeber sollten dort erläutern, dass die übertragenen Werke einer generativen KI entstammen und an ihnen keine Urheberrechte bestehen. Mit dieser eindeutigen Festlegung können sie Unsicherheiten der Vertragsparteien im Hinblick auf bestehende Rechte von Anfang an ausräumen.

Der Nutzungsvertrag muss regeln, dass Abnehmer die KI-Inhalte übertragen bekommen und zu vereinbarten Zwecken nutzen dürfen. Er sollte auch definieren, wie die Inhalte bereitstehen. Der Lizenzgeber sollte die Inhalte im Vertrag möglichst präzise spezifizieren. Bei KI-generierten Bildern kann er dies beispielsweise mit einem Bildtitel, einer detaillierten Beschreibung des Motivs sowie den Metadaten umsetzen. Falls ein Abnehmer den KI-generierten Output eigens beauftragt hat, kann man im Vertrag vereinbaren, dass auch der zugrundeliegende Prompt zum lizenzierten Inhalt gehört.

Neben der Vergütung muss der Vertrag festlegen, ob die Werke exklusiv überlassen werden. Vielleicht möchte der Ersteller sie selbst auch verwenden oder anderen zu anderen Bedingungen geben? Ist vorgesehen, dass die Inhalte nur auf einem bestimmten Kanal, etwa Online oder Print, laufen dürfen? Ist die Anzahl von Veröffentlichungen beschränkt? In diesen

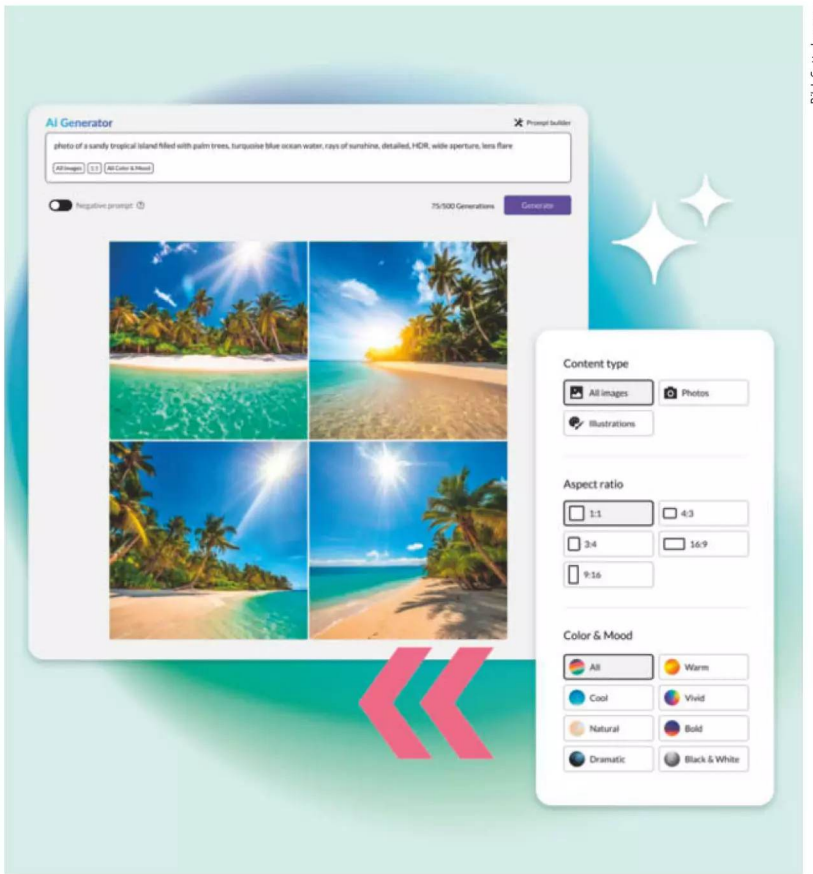


Bild: Getty Images

**Getty Images füttert eine bildgebende KI mit dem eigenen Fotobestand und verkauft den Zugang dazu, damit die Kunden selbst urheberrechtsfreie Werke auf eigener Datenbasis generieren können.**

Fällen sollte man etwaige Exklusivitätsklauseln tunlichst meiden oder im Einzelfall streichen.

Darüber hinaus sollten die Verträge eindeutige Haftungsklauseln enthalten. Sie regeln, wer bei möglichen Rechtsverletzungen durch KI-Inhalte verantwortlich ist. Insbesondere: Wie geht man mit Plagiatvorwürfen oder unbeabsichtigten Ähnlichkeiten zu bestehenden Werken um? Hier steht die Rechtsprechung noch am Anfang. Deshalb ist besondere Vorsicht angebracht, um spätere Konflikte zu vermeiden.

## Paradoxe Rechtslage

Sie könnten sich jetzt fragen: Warum sollte ich einen womöglich teuren Vertrag über die Nutzung von

Werken abschließen, an denen gar keine Urheberrechte bestehen? In der Tat führt die beschriebene Rechtslage zu dieser reichlich paradoxen Situation. Ein Vertragspartner, der KI-generierte Inhalte erwirbt, wird durch den Vertrag mit dem Ersteller faktisch weniger Rechte erhalten als jeder Dritte, der die Inhalte einfach ohne Nachfrage übernimmt.

Einschränkungen auf einsetzbare Medien oder die Nutzungsdauer etwa gelten nur zwischen den Beteiligten, aber nicht für jeden anderen. Mangels urheberrechtlichem Schutz ist es jedem Dritten erlaubt, die erworbenen Inhalte zu übernehmen und sie frei und ohne Einschränkungen zu nutzen. Er darf sie sogar bearbeiten, übersetzen und weitergeben, wenn er sie denn in die Finger bekommt.



Dennoch ist es für Kreative meist sinnvoll, die KI-generierten Werke mit Nutzungsverträgen de facto über die Rechtslage hinweg zu schützen. Denn es geht mitunter auch um teuer erworbene Erfahrung, um Exklusivität, um Zeitvergütung und eben um hohe Qualität bei der Nachbearbeitung. Immer öfter dürften kommerziell generierte KI-Inhalte künftig daher mit Wasserzeichen versehen sein beziehungsweise nur mit Zugangshürden in hoher Qualität zu beziehen sein.

Große Stockfoto-Anbieter finden zudem gerade ganz neue Wege, ihren riesigen Datenbestand mit generativer KI zu verschränken und neu auszuwerten. Getty Images etwa bietet selbst gar keine KI-generierten Bilder an. Stattdessen verkauft der riesige Bildermakler den Zugang zu einer bildgebenden KI, die mit Fotos von Getty Images gefüttert wurde und deshalb stattdessen „potenziell rechtskonforme“ Ergebnisse liefert – „unbeschränkte Haftungsfreistellung in jedem Download enthalten“.

## Fazit

Einigkeit besteht unter Juristen darüber, dass die Ergebnisse von generativer KI nicht urheberrechtlich geschützt sind. Die sich daraus ergebenden Konsequenzen müssen mit der zunehmenden Verbreitung dieser Technik noch Eingang in das rechtliche Alltagsleben finden. Dazu gehört auch die Frage, wo genau die Grenze zwischen rechtlich geschützten und frei nutzbaren Ergebnissen bei „Mischwerken“ zu ziehen ist.

Will man Rechte an KI-generierten Bildern, Texten oder Videos übertragen, gilt es, bestehende Nutzungsverträge entweder anzupassen oder sie gänzlich neu zu formulieren. Auf jeden Fall sollten Inhalte produzierende Unternehmen genauso wie Kreativschaffende hier besondere Sorgfalt walten lassen, um rechtliche Risiken zu minimieren und faire Bedingungen für die Nutzung von KI-Inhalten zu schaffen. (hob) **ct**

# Ihr Windows-Ratgeber

- Hochsicherheits-Windows
- Notfallsystem bauen, Viren suchen, Probleme lösen
- Unerwünschte Treiber- und Firmwareupdates verhindern
- Home- in Pro-Edition umwandeln
- Windows-Umzug mit c't-WIMage
- Whitelisting: Möglichkeiten, SRP/Restrictor, AppLocker

Heft für 14,90 € • PDF für 12,99 €  
Heft + PDF 19,90 €



[shop.heise.de/ct-windows24](https://shop.heise.de/ct-windows24)

**JETZT  
BESTELLEN!**

**Auch als  
Heft+PDF  
mit 28%  
Rabatt**





Bild: Albert Hulm

# Erste Lösungsansätze im KI-Urheberstreit

Seit dem Boom der generativen KI-Systeme ist der Nachweis darüber, dass ein Kunstwerk unrechtmäßig verwendet wurde, viel schwerer zu führen als bisher. Es geht nicht mehr um simple Kopien, sondern darum, dass generative KI Teile vieler Werke verwendet, um daraus etwas Neues zu generieren. In der Diskussion darum melden sich nun Verwertungsgesellschaften zu Wort.

Von **Falk Steiner**

**W**er die Werke anderer nicht nur zur Inspiration nutzt, sondern in Teilen kopiert, benötigt dafür eine Einwilligung vom Urheber oder dessen Erben. So schreibt es das europäische Recht vor, es sei denn, der kopierte Künstler ist bereits mehr als 70 Jahre tot (Regelschutzfrist). In der Frühzeit der Digitalisierung ging es bei Verletzungen des Urheberrechts um Fotokopien oder Musikmitschnitte. Heute ist die Situation ungleich

komplizierter. Den Kampf mit großen KI-Unternehmen führen nun vor allem Juristen.

Die kontinentaleuropäische Auffassung unterscheidet sich beim Urheberrecht stark von der angelsächsischen Betrachtung. Das wird bereits in der Wortwahl deutlich: Während das Urheberrecht auf der geistigen Schöpfung eines Urhebers fußt, geht es im Copyright, dem Vervielfältigungsrecht, um die Frage, wer ein Werk weiterverbreiten darf. Historisch

konnten in den USA nur registrierte Urheber über die Verbreitung ihrer Werke bestimmen, während dieses Recht in Europa bereits aufgrund der Schöpfung durch den Kreativen galt.

Generative KI eröffnet völlig neue Möglichkeiten, aus einem Original etwas Neues zu kreieren. Beispiele für Werkzeuge, die auf Basis vorhandener Werke ähnliche, aber maschinengenerierte Inhalte fabrizieren, gibt es mittlerweile viele. Unzählige KI-Textprogramme schreiben auf ein Prompt hin Texte aller Art, Plattformen wie Udio oder Suno generieren Musik und Bild-KIs wie Dall-E und Midjourney spucken auf Zuruf Bilder in beliebigem Stil aus. Mittlerweile gibt es die ersten KI-Werkzeuge, die ganze Videos produzieren. Bevor die Tools solches Material liefern können, müssen sie viel lernen. Tera- bis petabyteweise wurden die KI-Algorithmen mit sogenannten Token gefüttert, kleingehacktem Material aus vorhandenen Text-, Bild- und Audio-Dateien. Diese Bruchstücke dienen als Grundlage für das KI-generierte Material.

Wenn Nutzer einer KI ihren Prompt passend formulieren, kann das Ergebnis realen Kunstwerken recht ähnlich sehen. Ein Beispiel: Lautet die Nutzereingabe an ein KI-Musikprogramm „Song über die enttäuschte Liebe einer Frau, die ihre Loslösung und

Rachefantasien thematisiert, gesungen im Stil einer jungen amerikanischen Popsängerin mit leichtem Countryeinfluss“, dann könnte das Ergebnis verdächtig nach Taylor Swifts „Look what you made me do“ klingen. Die Bild-KI Midjourney ist berüchtigt dafür, dass sie sogar die Wasserzeichen von Vorschaubildern abkuppert, denen sie beim Abgrasen des Internets begegnet ist.

Solche KI-Produkte werfen die Frage auf, ob überhaupt zulässig ist, was die generative KI möglich macht. Juristisch existieren zur Zeit Einschränkungen bei den Rechten der Urheber. Weil nämlich die EU und ihre Mitgliedstaaten die Entwicklung von KI in Europa fördern wollen, haben sie KI-Unternehmen im Urheberrecht mit der Richtlinie zum Digitalen Binnenmarkt von 2018 privilegiert. Nach deutscher Rechtslage ist seitdem die Tokenisierung urheberrechtlich geschützten Materials für kommerzielles Training möglich, solange dagegen kein maschinenlesbares Verbot, ein sogenannter Nutzungsvorbehalt, besteht und Unternehmen beim Training ihrer KIs keine dauerhaften Kopien anfertigen. Eine zweite Einschränkung erlaubt die Verwendung geschützten Materials zu Forschungszwecken.

## Erstes Gerichtsurteil

Der Berliner Fotograf Robert Kneschke wollte prüfen lassen, wie weit diese Einschränkungen gehen. Vor dem Landgericht Hamburg verklagte er den Verein Laion, der unter anderem Fotos von Kneschke ausgewertet haben soll, um aus Bildern und Bildbeschreibungen eine Datenbank für das KI-Training aufzubauen. Im Oktober 2024 verlor Kneschke gegen Laion in erster Instanz. Doch was das LG Hamburg in diesem Fall urteilte, ist nur der Auftakt zu langen und intensiven Auseinandersetzungen.

Die Hamburger Richter sahen Laion als Sonderfall an, weil der Verein wissenschaftliche Forschung betreibt. Da Laion sich auf das Wissenschaftsprivileg berufen kann, muss der Verein sogenannte Nutzungsvorbehalte nicht in gleicher Weise berücksichtigen wie kommerzielle KI-Anbieter. Allerdings ist diese Auslegung erst einmal nur die Ansicht eines Gerichts. Weitere Verfahren vor anderen Richtern werden sicherlich folgen.

Noch verzwickter macht die Lage, dass Künstler nicht nur Token-Lieferanten, sondern längst auch selbst Nutzer von KI-Systemen sind, in der Postproduktion von Filmen, bei Werbung, intensiv auch bei der Produktion von Musik. „KI wurde schon lange genutzt“, erklärt Florian Drücke vom Bundes-



Home » Für Mitglieder » Rechte kennen

### Rechte kennen

Das Fundament unserer Arbeit als Illustratoren:in ist das **Urheberrecht**. Es schützt unsere geistigen und künstlerischen Leistungen und definiert mit den darin verankerten **Nutzungsrechten** die Rahmenbedingungen ihrer Verwertung. Eine weitere wichtige Grundlage für Auftragsverhandlungen sind die **Allgemeinen Geschäftsbedingungen**, die die IO für den individuellen Einsatz zur Verfügung stellt. Mit der Abgabe des ersten Angebots spielt auch das **Umsatz- und Steuerrecht** für Kreativschaffende eine wichtige Rolle.

**Berufsverbände wie die Illustratoren Organisation (IO) informieren ihre Mitglieder ausführlich über deren Rechte, insbesondere im Zusammenhang mit dem Urheberrecht.**



verband Musikindustrie. „Ein bekanntes Beispiel ist der neue Song ‚Now and Then‘ von den Beatles gewesen, wo KI eine große Rolle gespielt hat. Aber eben nicht generative KI, sondern eher als Werkzeug zur Unterstützung.“

Für Drücke ist klar, dass generative KI eine ganz andere Dimension hat, weil sie auf menschengemachten Werken beruht – und darum gehe die Auseinandersetzung nun. Komponisten, Künstler, Bands, Produzenten und Hersteller von Tonträgern haben seiner Ansicht nach eigene Rechte an den KI-Werken, was eine Einigung durchaus komplex mache.

Noch sind Klagen von Künstlern, Fotografinnen oder Musikern gegen KI-Unternehmen selten. Es wäre ein ungleicher Kampf. Immerhin sind Firmen wie OpenAI, Alphabet, Meta und Mistral finanziell problemlos in der Lage, jahrelang mit den besten Anwälten zu prozessieren. Gegen die großen KI-Firmen könnten einzelne Künstler wenig ausrichten, fürchten viele Betroffene.

## Angemessene Vergütung gefordert

„Wird generative KI eingesetzt, müssen die Kreativschaffenden und ihre am Kreativprozess beteiligten Partner angemessen vergütet werden“, sagt Juliane Fiedler von der Gesellschaft zur Verwertung von Leistungsrechten (GVL). „Es muss sichergestellt werden, dass beim Einsatz von KI – wie für jede urheberrechtliche Nutzung – die entsprechenden Rechte zu erwerben sind.“

Noch deutlicher wird die Verwertungsgesellschaft Wort. Sie meint, „dass die gesetzliche Erlaubnis für kommerzielles Text- und Datamining bei der Nutzung von geschützten Werken für KI-Training nicht anwendbar ist“, teilt Sprecherin Anette Frankenberger auf c't-Anfrage mit. Zugleich allerdings geht die VG Wort mit einer Änderung am sogenannten Wahrnehmungsvertrag einen Schritt hin zur kollektiven Rechteeräumung für KI-Nutzung. Diese Änderung soll es der VG Wort ermöglichen, Unternehmen zur internen Verwendung bestimmte Rechte an den Werken von VG-Wort-Mitgliedern zu gewähren. Das soll sowohl klassische KI-Anwendungen als auch generative KI betreffen. Allerdings sieht der Vertragsentwurf vor, dass die Lizenznehmer ihre Ergebnisse lediglich einem „bestimmt abgegrenzten Kreis von Teilnehmern“ zugänglich machen. Für die VG Wort ist diese Einschränkung zentral: „Es geht demnach nicht um Nutzungen durch große KI-Ent-

wickler wie OpenAI und andere“, betont VG-Wort-Sprecherin Frankenberger.

Im November 2024 hat die Gesellschaft für musikalische Aufführungsrechte (GEMA) als erste Verwertungsgesellschaft weltweit eine Klage wegen unlizenzierter Nutzung von geschützten Musikwerken gegen einen KI-Anbieter erhoben. Konkret geht es um das US-amerikanische Unternehmen OpenAI, den Anbieter von ChatGPT. Der Vorwurf: OpenAI habe die Texte der Lieder von GEMA-Mitgliedern für das Training von ChatGPT widerrechtlich verwendet. „Die Songs unserer Mitglieder sind nicht der kostenlose Rohstoff für die Geschäftsmodelle der Anbieter generativer KI-Systeme“, sagt GEMA-Geschäftsführer Tobias Holzmüller. OpenAI müsse Lizenzen erwerben und die Urheber fair vergüten. Die GEMA glaubt, dass sie OpenAI bei der Frage der Songtexte einen Verstoß gegen das Urheberrecht gut nachweisen kann. Nicht nur die Verarbeitung, sondern auch unerlaubte Abwandlungen seien durch ChatGPT entstanden – aus Sicht der GEMA Verstöße gegen das Urheberpersönlichkeitsrecht.

Auch die GEMA hat ein erstes KI-Lizenzmodell entwickelt. Darin fordert sie 30 Prozent aller Netto-Einnahmen, die ein generatives Modell erwirtschaftet, falls im Vorfeld GEMA-Repertoire für das Training genutzt wurde. Selbst wenn der Anbieter eines KI-Modells keinen Gewinn erwirtschaftet, soll nach den Wünschen der GEMA eine Mindestvergütung fällig werden. Zudem soll ein zweiter Teil der Lizenz die anschließende Nutzung ebenfalls abgabepflichtig machen: Dann würde – Training mit GEMA-Material vorausgesetzt – jeder KI-generierte Song, der im Supermarkt dudelt, den Künstlern ein paar Cent in die Kasse spülen.

Bei der VG Bild-Kunst ist eine Lizenz ebenfalls in Planung. Fotograf Robert Kneschke warnt allerdings davor, solch einer Lizenzerteilung zuzustimmen, weil sich Urheber damit seiner Meinung nach des Rechtes berauben könnten, individuell vor Gericht gegen KI-Firmen vorzugehen.

So weit wie VG Wort, VG Bild-Kunst und GEMA sind andere Verwertungsgesellschaften noch nicht. Bei der Verwertungsgesellschaft der Film- und Fernsehproduzenten (VFF) etwa, die für die Auftragsproduzenten von Fernsehsendern Leistungsschutzrechte vertritt, denkt man noch darüber nach, wie eine eigene KI-Lizenz aussehen könnte. Eine „Blaupause für eine Vergütung für die massenhafte Nutzung von Rechten durch KI-Anbieter könnte das gut eingespielte System der Vergütungsansprüche für die sogenannte Leermedienabgabe sein“, sagt Geschäftsfüh-



Aktuelles > Dossier: Künstliche Intelligenz > KI-Klage



## Die GEMA klagt für eine faire Vergütung

Die GEMA hat als erste Verwertungsgesellschaft weltweit eine Klage wegen unlizenzierter Nutzung von geschützten Musikwerken gegen einen Anbieter von Systemen generativer Künstlicher Intelligenz (KI) erhoben.

Konkret geht es um das US-amerikanische Unternehmen OpenAI, den Betreiber autogenerativer Chatbot-Systeme. Die GEMA wirft OpenAI vor, geschützte Songtexte von deutschen Urheberinnen und Urhebern wiederzugeben, ohne dafür Lizenzen erworben beziehungsweise die Urheberinnen und Urheber der genutzten Werke vergütet zu haben.

→ Zur Pressemitteilung vom 13.11.24

Mit einer Musterklage will die GEMA die Vergütungspflicht von KI-Anbietern in Europa klären.

# VG WORT

Die VG WORT   Teilnahme   Auszahlungen   Einnahmen-Tarife   Dokumente

[Startseite](#) / [Veröffentlichungen](#) / [Fragen / Antworten zur Änderung](#)

Kollektive Lizenzen mit erweiterter Wirkung

International

Termine

**Veröffentlichungen**

Änderung der Wahrnehmungsbedingungen

**Fragen / Antworten zur Änderung**

Geschäftsberichte

Transparenzberichte

Repräsentationsvereinbarungen

Reden / Vorträge

Stellungnahmen

Broschüren

Wortreport

## Fragen und Antworten zur neuen KI-Lizenz der VG WORT

1) Was ist der Gegenstand der neuen KI-Lizenz? +

2) Welche Unternehmen hat die VG Wort als potentielle Lizenznehmer im Blick? +

3) Welche Werke können seitens der Unternehmen genutzt werden? +

4) Was ist darunter zu verstehen, dass die Werke zuvor rechtmäßig erworben sein müssen? +

5) Was genau ist unter einer unternehmensinternen Nutzung zu verstehen? +

6) Geht es um die Nutzung der geschützten Werke als „Input“ für das Training der KI oder wird auch der „Output“ lizenziert? +

Die VG Wort hat ein KI-Lizenzmodell entwickelt, das aber nur für KI mit abgegrenztem Nutzerkreis verwendet werden dürfte.

c't KI-Wissen 2025

EU-Regulierung & Datenschutz 175

rer Albrecht Bischoffshausen. Auch er sieht allerdings klarere gesetzliche Regelungen für nötig an.

## Sorge bei KI-Unternehmen

Bei der Computer & Communications Industry Association CCIA Europe, der europäischen Interessenvertretung der IT-Industrie, sei man gesprächsbereit, erklärt Senior Policy Manager Boniface de Champris. Allerdings warne der Industrieverband „nachdrücklich vor der Einführung verpflichtender Lizenzmodelle, da diese die Wettbewerbsfähigkeit europäischer Unternehmen beeinträchtigen und Nutzer, Forscher, Entwickler und Unternehmen von der Nutzung modernster KI-Technologie abhalten würden.“

Ein großer Vorteil für die Techbranche liegt in dem System der Ausschüttungen über die Verwertungsgesellschaften: Würden sich die Anbieter selbst individuell Rechte einräumen lassen, müssten die KI-Modelle etwa auch berücksichtigen, wann die

jeweiligen Rechte erlöschen. Das wiederum können die Verwertungsgesellschaften deutlich einfacher handhaben.

Gespräche zwischen der CCIA Europe und den Verwertungsgesellschaften haben bisher aber wohl noch nicht stattgefunden. „Die Gesprächsbereitschaft hält sich in Grenzen“, berichtet zumindest die GEMA. „Bislang ist kein KI-Anbieter proaktiv auf die GEMA zugegangen, um eine Lizenz zu erwerben.“

## KI-VO ändert Rechtslage erneut

Die europäische KI-Verordnung verlangt in Artikel 53 von KI-Herstellern und -Anbietern, dass diese mittels „modernster Technologien“ die Wahrung von Urheberrechten sicherstellen. Wie genau das technisch ablaufen soll, schreibt die KI-VO nicht vor. Praktisch dürfte die Formulierung so zu verstehen sein: Die Richtlinie fordert von Urhebern, die ihre Werke nicht zum KI-Training verwertet sehen wollen, dass sie

## Gemeinsam stark?

1965 führte das „Gesetz über die Wahrnehmung von Urheberrechten und verwandten Schutzrechten“ in der Bundesrepublik die sogenannten Verwertungsgesellschaften ein. Die Idee: Wahrnehmungsberechtigte Mitglieder räumen einer Gesellschaft bestimmte Rechte ein. Bei der Verwertungsgesellschaft Wort (VG Wort) ist es beispielsweise das Recht, mit Verbänden von Kopierer- und Speicherherstellern über eine Urheberrechtsabgabe zu verhandeln; Verhandlungsbasis sind Schätzungen dazu, wie viele Kopien mithilfe dieser Hardware im Schnitt angefertigt werden.

Die Rechte von Kreativen der Musikbranche nimmt unter anderem die Gesellschaft für musikalische Aufführungs- und mechanische Vervielfältigungsrechte (GEMA) wahr. Sie sammelt treuhänderisch für ihre Mitglieder Abgaben ein, die anfallen, sobald Musik aus dem GEMA-Repertoire etwa auf einer öffentlichen Veranstaltung gespielt wird. Auch Pauschalabgaben für USB-Sticks, Festplatten oder Vereinbarungen zur Nutzung urheberrechtlich geschützter Musik auf Internetplattformen werden von der GEMA erhoben. Die Aufsicht über alle in Deutschland ansässigen Verwertungsgesellschaften übt das Deutsche Patent- und Markenamt (DPMA) aus.

### Liste der in Deutschland tätigen Verwertungsgesellschaften

Kurzname	vollständiger Name	Website
AGICOA	AGICOA Urheberrechtsschutz Gesellschaft mbH	<a href="http://www.agicoa.de">www.agicoa.de</a>
Corint Media	Corint Media GmbH	<a href="http://www.corint-media.com">www.corint-media.com</a>
GEMA	Gesellschaft für musikalische Aufführungs- und mechanische Vervielfältigungsrechte	<a href="http://www.gema.de">www.gema.de</a>
GÜFA	Gesellschaft zur Übernahme und Wahrnehmung von Filmaufführungsrechten mbH	<a href="http://www.guefa.de">www.guefa.de</a>
GVL	Gesellschaft zur Verwertung von Leistungsschutzrechten mbH	<a href="http://www.gvl.de">www.gvl.de</a>
GWFF	Gesellschaft zur Wahrnehmung von Film- und Fernsehrechten mbH	<a href="http://www.gwff.de">www.gwff.de</a>
GWVR	Gesellschaft zur Wahrnehmung von Veranstalterrechten mbH	<a href="http://www.gwvr.de">www.gwvr.de</a>
TWF	Treuhandgesellschaft Werbefilm mbH	<a href="http://www.twf-gmbh.de">www.twf-gmbh.de</a>
VFF	Verwertungsgesellschaft der Film- und Fernsehproduzenten mbH	<a href="http://www.vff.org">www.vff.org</a>
VG Bild-Kunst	Verwertungsgesellschaft Bild-Kunst	<a href="http://www.bildkunst.de">www.bildkunst.de</a>
VG Musikedition	Verwertungsgesellschaft Musikedition	<a href="http://www.vg-musikedition.de">www.vg-musikedition.de</a>
VG Wort	Verwertungsgesellschaft Wort	<a href="http://www.vgwort.de">www.vgwort.de</a>
VGF	Verwertungsgesellschaft für Nutzungsrechte an Filmwerken mbH	<a href="http://www.vgf.de">www.vgf.de</a>

maschinenlesbar einer Nutzung widersprechen. 2018 hätte man das etwa als Hinterlegung in der auf Webservern üblichen Datei robots.txt interpretieren können. Doch dass allein das Rechtswirkung entfalten soll, dürfte nach der KI-VO-Formulierung kaum haltbar sein: Die Fähigkeit von Crawlern, Nutzungsvorbehalt-Textstücke wie „Alle Rechte vorbehalten“ zu erfassen, dürfte in jedem Fall unter diese Formulierung fallen. Damit steigt der Druck auf die KI-Unternehmen weiter, sich um zweifelsfrei legale Wege zur KI-Fütterung zu bemühen. Der relevante Teil der KI-Verordnung tritt Anfang August 2025 in Kraft.

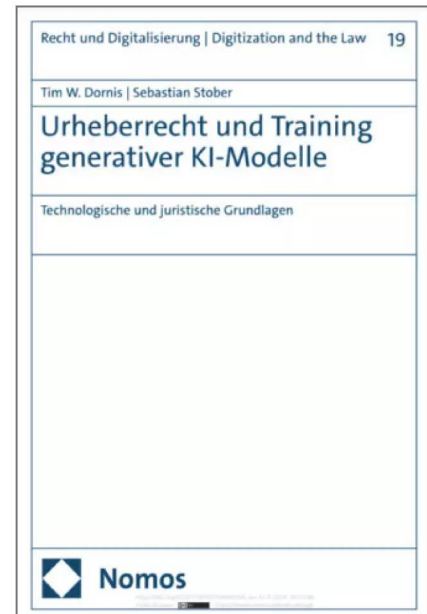
Für die GEMA ist das ein wichtiger Punkt. „Wenn Lizenzlösungen weiter auf sich warten lassen, muss vor allem auf EU-Ebene über weitere Maßnahmen nachgedacht werden, wie der Anteil der Urheberinnen und Urheber an der Wertschöpfung, die letztlich auf Grundlage ihrer Werke generiert wird, gesichert werden kann“, erklärt sie auf c’t-Anfrage. Die KI-Verordnung durchzusetzen wäre Aufgabe der zuständigen Behörden, in Deutschland also der Bundesnetzagentur.

Bei der CCIA Europe ist man sich des Problems mit der KI-Verordnung durchaus bewusst. Es gehe um eine schwierige Balance zwischen KI-Förderung und Urheberrecht, die in ihrer jetzigen Form unbedingt erhalten werden müsse, erklärt de Champris. „Diese umzukehren könnte dramatische Folgen für Nutzung und Entwicklung von KI haben“, warnt er die Politik vor Aktionismus. Um ein rechtssicheres Opt-out zu erreichen, müssten „flexible, skalierende, anpassbare, erkennbare und durchsetzbare automatische Lösungen“ her. Man sei bereit, auch darüber mit den Beteiligten zu sprechen.

## Klare Regeln verzweifelt gesucht

Verbände wie der Deutsche Journalistenverband empfehlen Autoren, zumindest für eine Nutzung in abgeschlossenen Gruppen, beispielsweise unternehmensintern, dem Lizenzvorhaben der VG Wort zuzustimmen. Generell stellt sich für Künstler die Frage: Lieber den kollektivrechtlichen Spatz auf dem Konto oder die Taube auf dem Dach der Gerichte?

Doch nicht nur Urheber und KI-Anbieter wollen profitieren: Die ersten Verlage fordern, dass die Autoren ihnen umfassende Rechte für ihr internes KI-Training einräumen. Der Berliner Wissenschaftsverlag De Gruyter will sich beispielsweise das Recht einräumen lassen, jedes Werk „in eigenen oder fremden Datenbanken zu speichern, um mittels einer automatischen Analyse Informationen über dessen



**In dieser umfangreichen Studie stuften die Autoren das Training generativer KI-Modelle als „Vervielfältigung urheberrechtlich geschützter Werke im Sinne des § 16 Abs. 1 UrhG“ ein.**

Muster und Merkmale zu gewinnen und dadurch maschinelles Lernen zu ermöglichen; dies umfasst das Training von Modellen, die aufgrund der gelernten Muster und Merkmale neue Daten generieren können (generative künstliche Intelligenz).“

Andere Verlage haben sich kurzerhand die Rechte für jedwede denkbare, künftig mögliche Nutzungsart einräumen lassen. Ob das tatsächlich auch die Nutzung durch fremde KI-Modelle umfasst, ist unter Juristen durchaus umstritten. Noch fehlen Regelungsvorbilder, die die KI-Entwicklung nicht verhindern, aber zugleich Urheberrechte wahren. Bisher hat kein Land der Welt abschließend klare Regelungen zum Verhältnis von Urheberrecht und KI beschlossen, sagen Experten. Und so könnte es sein, dass am Ende doch die EU-Kommission noch einmal prüfen muss, ob nicht klarere Regeln erforderlich sind. Ein konkretes Vorhaben dazu steht derzeit nicht auf der Agenda. Aber EU-Kommissionspräsidentin Ursula von der Leyen hat eine „KI-Strategie für die Kultur- und Kreativwirtschaft“ für die laufende EU-Legislaturperiode angekündigt. (dwi) **ct**



## Vorschau: c't Besser Wissen

Ab dem 28. März im Handel und auf [ct.de](https://ct.de)

### Digitale Helfer für mehr Ordnung im Kopf

Das neue c't Sonderheft bietet einen umfassenden Überblick über Methoden und Tools, das eigene Know-how zu organisieren, zu managen und zu ordnen. „c't Besser Wissen“ erklärt, wie man Wissen effizient sammelt, Gedanken strukturiert und die gesammelten Informationen übersichtlich aufbereitet – auch jenseits von Schule und Studium. Dafür haben Autoren unter an-

derem die wichtigsten Programme rund um das Thema Wissensmanagement getestet. Dabei spielt natürlich auch KI eine wichtige Rolle, sowohl beim Sammeln als auch beim Aufbereiten von Wissen. Außerdem präsentiert das Heft Methoden, Gedanken besser zu sortieren.

Weitere c't Sonderhefte: [heise.de/s/00MxL](https://heise.de/s/00MxL)

## Themenschwerpunkte

### Wissen sammeln

- Wikipedia richtig nutzen
- Überblick und Beratung zu RSS-Readern
- Wissenschaftliche Literatur verwalten und mit KI durchforsten

### Gedanken strukturieren

- Notizprogramme im Vergleich
- Einführungen in Obsidian, Notion und ProNotes
- Backups und Datenexporte aus Notizprogrammen
- Notizen organisieren mit der Zettelkasten-Methode

### Ergebnisse darstellen und zusammenfassen

- Mindmap-Tools im Vergleich
- Einstieg in Datawrapper
- KI fürs wissenschaftliche Schreiben und Publizieren
- Eigene E-Books erstellen

 heise academy

Für erfolgreiche  
IT-Teams  
von morgen



Interesse geweckt? Hier mehr erfahren:  
[heise-academy.de/academy-pass](https://heise-academy.de/academy-pass)

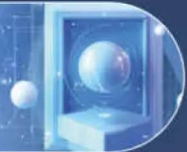
# KI-Webinare 2025



## KI-Modellkunde

Fundierte Entscheidungshilfe für KI-Modelle im Unternehmen

11.03.2025 | 21.10.2025



## KI-Governance

So bringen Sie generative KI wirksam in Arbeitsprozesse

18.03.2025



## EU AI Act

Die neue KI-Verordnung verstehen und umsetzen

15.04.2025



## KI-Change-Management

Mensch und Maschine erfolgreich im Unternehmen zusammenbringen

20.05.2025



## Das bringt generative KI wirklich

Cases & Ergebnisse aus dem Alltag der heise group

17.06.2025



## Generative KI im Einsatz

Fallbeispiele aus deutschen Unternehmen

15.07.2025



## KI-Agenten

Grundlagen, Anwendungsgebiete und Zukunftsperspektiven

19.08.2025



**Jetzt Ticket sichern:**

[heise-academy.de/formate/webinare](https://heise-academy.de/formate/webinare)

# FREITAG IST c't-TAG!\*

Jetzt 5x c't lesen

für 24,00 €  
statt 31,75 €\*\*

\*\* im Vergleich zum Standard-Abo

30%  
Rabatt!



\*Endlich Wochenende! Endlich genug Zeit, um in der c't zu stöbern. Entdecken Sie bei uns die neuesten Technik-Innovationen, finden Sie passende Hard- und Software und erweitern Sie Ihr nerdiges Fachwissen. **Testen Sie doch mal unser Angebot: Lesen Sie 5 Ausgaben c't mit 30 % Rabatt – als Heft, digital in der App, im Browser oder als PDF. On top gibt's noch ein Geschenk Ihrer Wahl.**

Jetzt bestellen:

[ct.de/meintag](https://ct.de/meintag)

