

Rechtssichere Hochschulprüfungen mit und trotz generativer KI

– Plädoyer für eine große Studien- und Prüfungsreform –

Ass. jur. Sarah Rachut

Geschäftsführerin

TUM Center for Digital Public Services

School of Social Sciences and Technology

Technische Universität München

Dilemma ChatGPT...und nun?



5 Mythen zum Hype um einen Gerichtsbeschluss

... und was wir daraus lernen können

VG München, Beschl. v. 28.11.23 – 3 M 23.4371

SEINE BEWERBUNG WAR**ZU GUT, UM WAHR ZU SEIN**

(ALSO WORTWÖRTLICH)

Ein Bewerber wurde an der TU München **für einen Master-Studiengang abgelehnt**.

DAS IST PASSIERT:

Den Prüfenden fiel auf, dass Sprache und Inhalt des Bewerbungs-Essays **viel besser waren als das, was sie sonst gewohnt sind**. Der Bewerber habe sich ein Jahr vorher schon mal mit einem wohl schlechteren Text beworben.

Sie folgerten: **Der Inhalt muss von einer KI erstellt worden sein**. Das verstößt gegen die wissenschaftliche Sorgfalt.

Der Fall landete vor Gericht. Urteil: Die Uni hat recht. Die Indizien reichten, um ihn abzulehnen. Es ist laut FAZ **das erste Urteil zur Verwendung von Künstlicher Intelligenz an einer Hochschule**.

Quelle: Bayerische Staatskanzlei, FAZ

Bewerber darf nicht zur Uni: So gut wird's nur mit ChatGPT



Im ersten Anlauf schaffte es ein Bewerber nicht in einen Masterstudiengang an der TU München, weil sein Essay nicht gut genug war. Beim zweiten Versuch war er dann so gut, dass die Uni ihn ausschloss, weil man von einem KI-Text ausging. Das VG München hält die Vermutung für valide.

Mythos 1

Der Fall ist entschieden, mit diesem Urteil haben wir Rechtssicherheit!

Mythos 2

Es ist doch gerecht, wenn einer der täuscht, durchfällt!

Mythos 3

Die Nutzung von KI in Prüfungen ist verboten, solange sie nicht als Hilfsmittel ausdrücklich erlaubt ist!

Mythos 4

Man kann eine solche Nutzung generativer KI durch KI-Detektoren beweisen!

Mythos 5

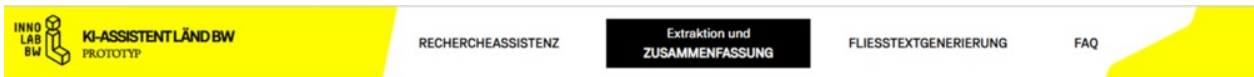
Die Erkenntnis KI-typischer Merkmale in Texten begründet einen
Anscheinsbeweis – es obliegt dem Prüfling, dies zu entkräften!

Generative KI ist gekommen, um zu bleiben
und wird (so oder so) genutzt.

ChatGPT

☀ Examples	⚡ Capabilities	⚠ Limitations
"Explain quantum computing in simple terms" →	Remembers what user said earlier in the conversation	May occasionally generate incorrect information
"Got any creative ideas for a 10 year old's birthday?" →	Allows user to provide follow-up corrections	May occasionally produce harmful instructions or biased content
"How do I make an HTTP request in Javascript?" →	Trained to decline inappropriate requests	Limited knowledge of world and events after 2021

Research Preview: ChatGPT is optimized for dialogue. Our goal is to make AI systems more natural to interact with, and your feedback will help us improve our systems and make



⚠ ACHTUNG! Prototyp: Aus Datenschutz-Gründen keine personenbezogenen Daten oder Verschlusssachen eingeben oder hochladen!

Sie als digitales Assistenzsystem sinnvoll zu nutzen, wird am modernen Arbeitsplatz erwartet. Das Studium sollte darauf vorbereiten.

The interface is titled 'STANDARD 1 KV-VERMERK 1'. It contains a main instruction: 'Laden Sie ein Dokument (.doc/.docx/.pdf, max. 20 MB) hoch, oder geben Sie einen Text ein, zu dem Sie eine Zusammenfassung erstellen möchten.' Below this are two input options: 'Upload 1' with a dashed box containing 'Keine Datei ausgewählt' and a 'DATEI AUSWÄHLEN' button, and 'Texteingabe 1' with a text area 'Text zum Zusammenfassen einfügen'. A '-oder-' separator is between them. To the right is the 'EINSTELLUNGEN' panel with 'Länge 1' (options: kurz, mittel, lang) and 'Art 1' (options: von oben nach unten, nach Themen sortiert). A 'ZUSAMMENFASSEN' button is at the bottom right. A footer bar contains an information icon and the text: 'Die Zusammenfassung wird nach der Erstellung unten angezeigt'.

Ihr Einsatz ist ohne explizite Regelung nicht „*ohnehin verboten*“. Verbote müssen hinreichend bestimmt sein, weil die Enttäuschung berechtigter Erwartungen zugleich das Grundrecht auf Chancengleichheit verletzt.



ORDNUNG DER WISSENSCHAFT

Heft 2 / 2024

Aufsätze

Dirk Heckmann und Sarah Rachut Rechtssichere Hochschulprüfungen mit und trotz generativer KI **85-100**

Selbst, wenn ein Verbot in Einzelfällen vertretbar erscheint: Es wäre unter rechtsstaatlichen Maßstäben nicht überprüfbar, durchsetzbar und damit als staatliche Steuerungsmaßnahme ungeeignet.

Umgekehrt verstößt eine undifferenzierte Erlaubnis oder das einfache Dulden der skizzierten Einsatzszenarien ebenso das Grundrecht auf Chancengleichheit, weil es eine angemessene Bestenauslese torpediert.



Rechtssichere Hochschulprüfungen mit und trotz generativer KI

- Plädoyer für eine große Studien- und Prüfungsreform -

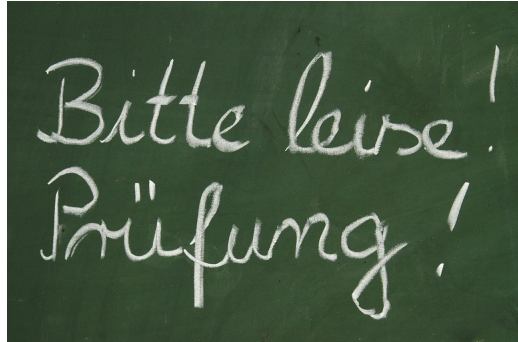
Plädoyer für eine große Studien- und Prüfungsreform

1. Schritt: Modulkataloge & Modulbeschreibungen überarbeiten



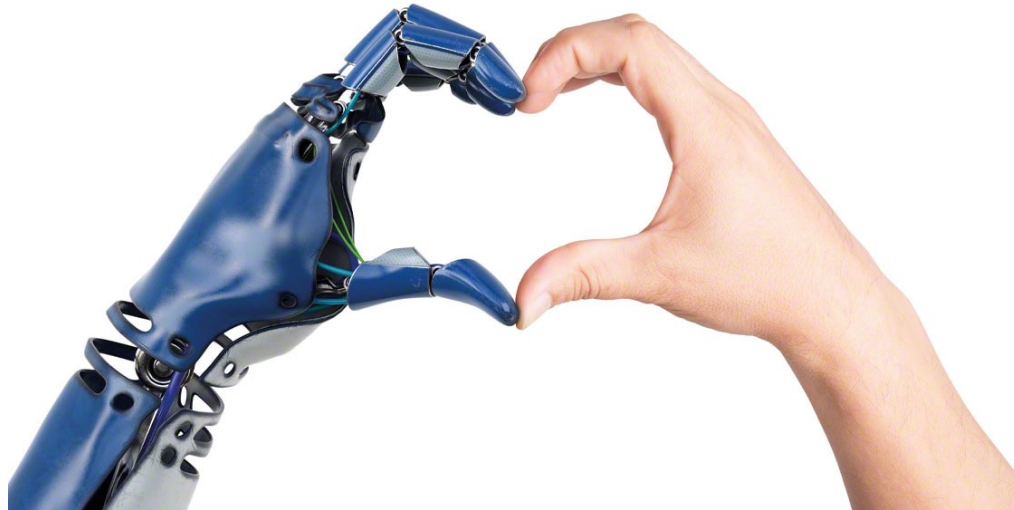
Plädoyer für eine große Studien- und Prüfungsreform

2. Schritt: Prüfungsformate überdenken



Plädoyer für eine große Studien- und Prüfungsreform

3. Schritt: Mut zu Co-Creation!



Reaktionen



Rachut, NJW 2024, 1052
Braegelmann, RdI 2024, 188

Jedenfalls kann man aus der Entscheidung nicht im Ernst ableiten, dass eine enorme Verbesserung in der schriftlichen Leistung im Vergleich zu früheren schlechteren Leistungen demnächst immer *prima facie* für einen KI-basierten Betragsversuch sprechen müsse. So viel Glauben an den Menschen und seine eigene Verbesserungsfähigkeit sollte man sich bewahren.

Merkwürdigerweise hat die Universität selber versucht, einen ähnlichen Essay mit ChatGPT zu generieren. Unklar bleibt, mit welcher Version (mit GPT 3.5, schlecht und kostenlos, oder GPT 4.0, sehr gut, nur mit Ab?) Das Gericht erwähnt dies, als sei es relevant. Ist es aber nicht. Selbstverständlich können Chatbots ähnliche Texte wie Menschen generieren, das ist ja ihre Aufgabe! Die Frage ist aber: Ist das hier geschehen? Das kann man nicht dadurch belegen, indem man späterhin noch einen ähnlichen Text per KI generiert. Das Gericht hat also nicht entschieden, dass die Tatsache, dass eine KI einen ähnlichen Text generieren kann, allein schon dafür spricht, dass auch der ursprünglich vorgelegte Text von einer KI generiert wurde. Das wäre auch eine kuriose Regel. Nur weil ChatGPT et al. sehr gut ähnliche Texte verfassen, macht das andere menschengeschriebene nicht nachträglich zu KI-generierten Texten. Vorzüglich wurde - als Gegenprobe - für diese Anmerkung die Entscheidung des VG München in ChatGPT hochgeladen (Urheberrechtskonform, vgl. § 5 UrhG) und der Bot gebeten, ein gegenteiliges Urteil zu schreiben: Das ging in 20 Sekunden, tadellos. Das bedeutet eben auch nicht, dass das VG München seine Entscheidung (partiell) mit ChatGPT geschrieben hätte.

Auf die KI-Überprüfungssoftware kann es im Ergebnis auch nicht an, denn diese hat auch große Probleme, wie es das Gericht selbst sagt: „Aus der Tatsache, dass Textpassagen nicht als KI-generiert gekennzeichnet sind, kann daher auch auf Grundlage der Softwareüberprüfung nicht gefolgert werden, dass sie nicht von künstlicher Intelligenz verfasst sind.“ – wenn man das ernst nimmt, müsste es ja auch andersherum gelten, dass Textpassagen, die als KI-generiert von der Software gekennzeichnet worden sind, doch nicht von KI generiert worden sind! Zukünftig sollte man, ähnlich wie im Bußgeld-Sachen bei Tempo- und Rodlichverstößen darauf drängen, als Betroffener, so eine Überprüfungssoftware von einem Forensiker untersuchen zu lassen (Vorlage des Quellcodes und des Algorithmus!). Im Moment scheint so eine KI-Überprüfungssoftware eine Blackbox zu sein, was zumindest im öffentlich-rechtlichen Bereich nicht ausreichend für ein nachvollziehbares Verwaltungshandeln der Verwaltung. Es gibt zwar noch Hoffnungen, dass Software-KI-Erkennung im Großen und Ganzen funktioniert, aber (so eingehend aus der Perspektive der universitären Praxis und mit Prüfung verschiedener KI-Tools Kesper, KI-Text in Prüfungsarbeiten erkennen

– Probleme und Lösungen: S. 7, abrufbar unter: <https://beck-link.de/da76>):

„Sich alleine auf das [KI-Texterkennungswerkzeug] zu verlassen, ist indes nicht angezeigt. So gilt es zu bedenken, dass KI-erzeugter Text soweit nachbearbeitet worden sein kann, dass eine automatische Detektion zu falschen Ergebnissen führt. (...) [E]s ist darauf hinzuweisen, dass Urteile durch ein AI-Detection-Tool prinzipiell fehleranfällig sind, denn sie beruhen auf Wahrscheinlichkeitsberechnungen, die immer auch Irrtümern unterliegen können. Bisher fehlten uns verlässliche Daten zu anzunehmenden Irrtumswahrscheinlichkeiten, die mutmaßlich auch textartenspezifisch erhoben werden müssten. (...) Bei gewichtigen Entscheidungen, etwa einer auszusprechenden Exmatrikulation wegen eines Verstoßes gegen die Eigenständigkeitsklärung, bleibt eine ausschließliche Argumentation mit einem AI-Detection-Ergebnis problematisch. (...) Wenn Prüfer/-innen den Verdacht äußern, dass eine Prüfungsarbeit mehrheitlich mit KI-Unterstützung verfasst wurde, werden sie dafür mit Sicherheit keine statistische Analyse im Kopf vorgenommen haben.“

Je besser die Large Language Models der gängigen Chatbots, desto weniger wird man wohl KI-generierten Text automatisch (mit der Betonung auf „automatisch“, also per Software, nicht nach eingehender Prüfung) ausfindig machen können. Es bleibt beim naheliegenden Fazit zu KI-Detektoren des Chatbot-Experten *Ethan Mollick* in seinem FAQ zu automatischen, also software-basierten AI Detectors (Mollick, abrufbar unter: <https://beck-link.de/6f3kf>, Übersetzung ins Deutsche):

„KI-Detektoren funktionieren nicht. Sofern sie überhaupt funktionieren, können sie durch geringfügige Änderungen am Text überwinden werden. Und was noch schlimmer ist: Sie haben eine hohe Falsch-Positiv-Rate (false-positive rate) und neigen dazu, Menschen der Verwendung von KI zu beschuldigen, obwohl diese keine KI verwendeten, insbesondere Schüler/Studenten, für die [die jeweilige Unterrichtsprache] nicht die Muttersprache ist. Die fälschlicherweise Beschuldigten haben keine Möglichkeit, sich zu wehren, weil sie nicht beweisen können, dass sie keine KI verwendet haben. Man kann KI-Text nicht automatisch erkennen. (...) [Menschen] denken vielleicht, dass sie gut darin sind, KI-Texte zu erkennen, aber sie sind nur gut darin, schlechte KI-Texte zu erkennen, und sie kombinieren das mit ihren eigenen Vorurteilen und Heuristiken darüber, wer KI verwenden könnte. Nach ein paar Prompts wird KI-generiertes Schreiben nicht wie generisches KI-Schreiben. (...)

Happy to discuss!



Sarah-Rachut



sarah.rachut@tum.de

www.TUM-CDPS.de