

# Verantwortungsvoller Umgang mit Bias in Großen Sprachmodellen

Proseminar Thesis von

Matthias Halfmann

Institut für Informationssicherheit und Verlässlichkeit (KASTEL)

Betreuende Mitarbeiter: Dominik Fuchß, Marcel Krüger

Mit der zunehmenden Verbreitung und Anwendung Großer Sprachmodelle (kurz: LLMs) wie ChatGPT in nahezu allen gesellschaftlichen Bereichen werden Bedenken hinsichtlich der Gefahren durch potenzielle Bias-Behaftung dieser laut. Diese Arbeit untersucht die Natur und Ursachen von Bias in LLMs und stellt dar, wie diese Verzerrungen durch technische und ethische Ansätze gemindert werden können. Zunächst wird erläutert, wie Bias in den Modellen identifiziert und gemessen wird. Anschließend werden verschiedene Debiasing-Strategien präsentiert, die in unterschiedliche Phasen des Modellentwicklungszyklus eingreifen. Trotz dieser technischen Maßnahmen bleibt die vollständige Beseitigung von Bias eine Herausforderung. Daher wird die Notwendigkeit einer umfassenden ethischen Reflexion und zusätzlichen Ansätzen diskutiert, um den verantwortungsvollen Einsatz von LLMs zu gewährleisten und die gesellschaftlichen Implikationen von Bias in diesen Modellen zu adressieren.

## 1 Einleitung

Mit der Entwicklung öffentlicher, leicht zugänglicher und nutzerfreundlicher Chat-Interfaces zu Großen Sprachmodellen (engl. Large Language Models, LLMs) wie GPT-3.5 oder GPT-4 von OpenAI und im Zuge steigender Popularität und Relevanz durch Anwendung dieser in Wirtschaft, Wissenschaft, Bildung, Medizin und weiteren Bereichen des gesellschaftlichen Zusammenlebens [1, 2, 3] werden Fragen über die Risiken von LLMs und ihres Einsatzes laut [4].

Eine Vielzahl von Studien hat gezeigt, dass diese neue Generation von Modellen zur Verarbeitung natürlicher Sprache (engl. Natural Language Processing, NLP) dazu neigt, soziale und kulturelle Bias zu reproduzieren und zu verstärken [5, 6, 7, 8]. Eine solche Biasbehaftung kann sich in Form von Stereotypen, Vorurteilen und Diskriminierung in

Ausgaben der LLMs manifestieren und hat das Potenzial, die Qualität und Fairness von Anwendungen, die auf diesen Modellen basieren, zu beeinträchtigen [4].

Dabei ist Bias kein neues Problem moderner LLMs wie GPT-4 oder von Konkurrenzprodukten wie Claude<sup>1</sup>, Gemini<sup>2</sup> und weiteren, sondern bekannt und gegenwärtig in der Forschung zu NLP und maschinellem Lernen [5, 6, 7]. So zeigten bereits frühere auf großen Textkorpora trainierte Systeme Bias wie beispielsweise die Google Autovervollständigung [9, 10].

Da Gefahren von Bias in LLMs stark an den Umgang und die Anwendung dieser Systeme in der Praxis geknüpft sind, soll nachfolgend dargestellt und diskutiert werden, welche Faktoren ursächlich für die Ausprägung von Bias in den Modellen sind und wie technische und ethische Ansätze diesen Gefahren nur gemeinsam vorbeugen können. Hierzu wird zunächst die Messung von Bias in LLMs und deren Funktionsweise geklärt, um technische Methoden zur Reduzierung von Bias zu bewerten und anschließend mit einer ethischen Betrachtung an diese anzuknüpfen.

## 2 Messung von Bias in LLMs

Um Aussagen über Biasbehaftung treffen und Techniken zur Reduzierung von Bias in LLMs evaluieren zu können, werden verschiedene Tests und Metriken benötigt. Im Folgenden sollen daher grundlegende Begriffe wie Bias, dessen Messung sowie die Funktionsweise Großer Sprachmodelle geklärt werden.

Aufgrund der semantischen Komplexität von Text, insbesondere von generierten Ausgaben aus LLMs, ist es schwierig, Bias in diesen zu quantifizieren, und Forschung in diesem Bereich noch jung. Es existieren verschiedene technische Ansätze, um Bias in LLMs zu messen. Hierzu sollen zunächst mögliche Ursachen anhand der Architektur von Transformer-basierten Großen Sprachmodellen dargestellt werden.

### 2.1 Was ist Bias?

Die Definition von Bias unterscheidet sich stark zwischen verschiedenen wissenschaftlichen Disziplinen, wobei in diesem Paper der Ausrichtung der Sozialpsychologie gefolgt werden soll. Bias beschreibt systematische Verzerrungen, Falschdarstellungen oder faktische Verfälschungen in den Ausgaben dieser Modelle, die dazu führen können, dass zu bestimmten Gruppen und Minderheiten Stereotype und Vorurteile auf Basis falscher Vermutungen oder Daten erlernt werden [3]. Bekannte Bias sind beispielsweise der sogenannte Gender-Bias, der durch geschlechtsidentitäre Stereotype geprägt ist.

Insbesondere wird zwischen implizitem und explizitem Bias unterschieden. Impliziter Bias umfasst Vorurteile, Assoziationen und Reaktionen, die automatisch und oft unbewusst bei Auftreten eines relevanten Stimulus entstehen, während expliziter Bias bewusst wahrgenommene und geäußerte Präferenzen, Überzeugungen und Einstellungen einschließt [11]. Eine derartige Differenzierung von Bias ist von besonderer Relevanz zum sogenannten

---

<sup>1</sup>Claude ist ein AI-Assistent von Anthropic

<sup>2</sup>Google Gemini ist das LLM Chat-Interface von Google

„Benchmarking“, bei dem Sprachmodelle auf Bias getestet und dieser quantitativ gemessen wird.

## 2.2 Funktionsweise von LLMs

Große Sprachmodelle sind ein Teilgebiet der sogenannten Künstlichen Intelligenz (KI, engl. Artificial Intelligence, AI). Die Sprachmodelle, auf die sich im Folgenden beschränkt werden soll, basieren auf der Transformer-Architektur [12, 13], die es ermöglicht, zusammenhängende Texte zu generieren, indem sie die Wahrscheinlichkeit für das nächste Wort in einer Sequenz von Wörtern in ihrem jeweiligen Kontext berechnen.

Dabei arbeitet das Modell auf sogenannten Tokens, einer Sub-Worteinheit, welche die kleinste Einheit darstellt, die das Modell verarbeiten kann. Dabei wird Eingabetext in Form einer Wortsequenz durch einen Tokenizer in eine Sequenz von Tokens überführt. Ein Beispiel für den Tokenizer, der von OpenAI in GPT-3.5 verwendet wird, ist in Abbildung 1 dargestellt. Man erkennt, dass der Eingabetext sowohl in Wörter als auch in Sub-Wörter unterteilt wird. Nicht jedes Wort entspricht also zwangsläufig einem ganzen Token [14].

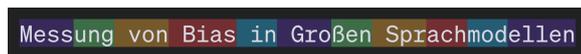


Abbildung 1: OpenAI Tokenizer für GPT-3.5 [14]

Um vollständige, zusammenhängende Texte erzeugen zu lassen, muss jedes neu generierte Token in der Ausgabe inklusive der ursprünglichen Eingabe als neue Eingabe des Modells eingegeben werden. Diese Technik wird als autoregressives Decoding bezeichnet.

Während des Trainings eines LLMs wird das Modell iterativ mit Eingabe- und Ausgabe-Paaren aus einer großen Menge an Textdaten, die vorwiegend aus Textkorpora des Internets stammen, trainiert. Dabei wird das Modell anhand einer Verlustfunktion (engl. Loss function) bewertet, die misst, wie gut die Vorhersage des Modells mit der erwarteten Ausgabe übereinstimmt. Ziel des Trainings ist es, die Wahrscheinlichkeit für das nächste erwartete Token in der Ausgabe zu maximieren, was gleichbedeutend damit ist, die Verlustfunktion zu minimieren [13]. Hierzu passt das Modell fortwährend die internen Parameter auf Basis des berechneten Fehlers an.

Die Architektur eines Transformer-basierten Sprachmodells besteht aus mehreren Schichten, die ihrerseits mehrere Untereinheiten beinhalten. Die wichtigsten Konzepte und Schichten, die folgend Erwähnung finden sollen, lauten Embedding, Self-Attention und Decoder [15].

### 2.2.1 Embedding

Einbettungen (engl. Embeddings) sind Vektoren in einem hochdimensionalen Vektorraum (bei GPT-3 12.288 Dimensionen), die die einzelnen Tokens (bei GPT-3 50.257 verschiedene Tokens) abbilden. Bei GPT-3 sind die Embeddings Teil der trainierbaren Architektur und werden während des Trainingsprozesses des Modells angepasst. Dabei „erlernt“ das Modell eine Kodierung semantischer Informationen der Tokens in den Embeddings. Dies führt mit ausreichenden Trainingsiterationen dazu, dass einzelne Dimensionen in den Embeddings

bestimmte Semantiken wie beispielsweise Gender repräsentieren können [16]. Ein Beispiel für die Embeddings von „King“ und „Queen“ in einem zweidimensionalen Raum ist in Abbildung 2 dargestellt. Ebenso vereinfacht hervorgehoben sind zwei Dimensionen „gender“ und „royal“, welche die jeweilige Einbettung ihrer Semantik in den Embedding-Raum verdeutlichen sollen. Da die Embeddings in GPT-3 12.288-dimensional sind, verwendet die Abbildung lediglich eine vereinfachte zweidimensionale Projektion. Tatsächlich können sich semantische Informationen auch in mehrdimensionalen Unterräumen einordnen und dann nicht durch einen einzelnen Vektor dargestellt werden.

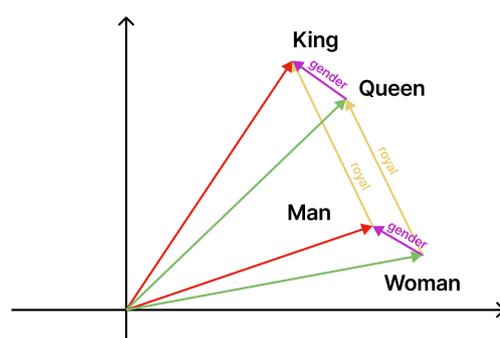


Abbildung 2: Einfache Darstellung der projizierten Embeddings von „King“ und „Queen“ in einem zweidimensionalen Raum [17]

Im Fall von GPT-3 werden zusätzlich zu der Semantik auch die Position der Tokens in der Eingabe mittels Sinusoiden<sup>3</sup> kodiert.

### 2.2.2 Self-Attention

Selbstaufmerksamkeit (engl. Self-Attention) in den nachfolgenden Schichten des Modells, ermöglicht es, die Beziehung zwischen Tokens bzw. Wörtern der Eingabe zu modellieren, untereinander zu gewichten und Kontext zu erfassen [12, 15]. Dabei wird für jedes Token in der Eingabe ein Gewicht berechnet, das angibt, wie stark das Token mit anderen Tokens der Eingabe korreliert. Diese Gewichte werden dann auf die Embeddings angewendet, um eine gewichtete Summe der Embeddings zu berechnen, die als Kontextrepräsentation für das Token dient [15]. Durch die gemeinsame Verrechnung der Tokens in der Eingabe arbeiten die Schichten auf eine statistische Vorhersage des nächsten Tokens der Ausgabe hin.

### 2.2.3 Decoder

Der Decoder in einem Transformer-Modell wie GPT-3 ist für die Textgenerierung verantwortlich. Er arbeitet autoregressiv, was bedeutet, dass er jedes Token der Ausgabe

---

<sup>3</sup>Sinusoiden sind sinusförmige Funktionen, welche durch Skalierung sowie Phasenverschiebung gebildet werden. [18]

sequentiell generiert und dabei die bereits erzeugten Tokens berücksichtigt [13, 15]. Der Prozess beinhaltet die folgenden Schritte:

Zunächst wird die Eingabe (ein Text- oder Token-Fragment) durch die Self-Attention-Schichten unter Einbezug vorheriger Tokens verarbeitet, um die Beziehungen und den Kontext zwischen den Tokens zu erfassen. Basierend auf dem aktuellen Kontext und den zuvor generierten Tokens berechnet der Decoder die Wahrscheinlichkeitsverteilung für das nächste Token. Das nächste Ausgabtoken wird generiert, indem das Token mit der höchsten Wahrscheinlichkeit ausgewählt (oder durch Sampling-Methoden bestimmt) wird. Dieser Prozess wird iterativ wiederholt, bis ein spezielles End-Token erreicht oder die gewünschte Sequenzlänge erzielt ist.

## 2.3 Testen auf Bias

Eine Biasbehaftung lässt sich in LLMs unter anderem durch gezieltes Prompting aufzeigen. Zwar sind derzeitige Sprachmodelle wie GPT-3.5 bereits resilient gegenüber vorwiegend explizitem Bias in den Ein- und Ausgaben des Modells, mit einigen gängigen Tests lässt sich impliziter Bias in den Modellen jedoch leicht offenbaren [7].

Meist wird sich im Design solcher Prompt-basierten Tests auf ein spezifischen Bias wie Gender-Bias fokussiert. Dabei wird das Modell mit einheitlich einem Schema folgenden Eingaben „konfrontiert“, in denen es aus einer Sammlung von Begriffen sowie einem männlich und einem weiblich konnotierten Namen je einen Begriff einem Namen zuordnen soll. Ein Beispiel für ein solches Schema ist der sogenannte „Word-Embedding-Association-Test“ (WEAT) [7], der auf der Assoziation von Begriffen basiert, oder der LLM IAT [7].

Eine Instanz dieses Tests kann wie in Abbildung 3 dargestellt aussehen. Dabei erhält das Modell eine Eingabe, die aus einer Sammlung von Begriffen und zwei Namen besteht. Das Modell soll nun die Begriffe den Namen zuordnen. In diesem Fall ist die Eingabe so konzipiert, dass sie impliziten Bias gegenüber einem der Namen aufdecken soll. Der Bias wird dann anhand der Wahrscheinlichkeit, mit der das Modell die Begriffe den Namen zuordnet, gemessen. Um Fehler zu vermeiden, werden Testinstanzen so gewählt, dass sie nicht in den Trainingsdaten enthalten sind, da das Modell sonst einfach das exakte Datum ausgeben könnte.

Ein weiterer Test, der „LLM Decision Bias“, basiert auf der Entscheidung über das uneindeutige Beziehungssubjekt eines Pronomens in einem Satz und lehnt sich an den WinoBias-Test [19] an. Ein Beispiel für eine Instanz dieses Tests ist in Abbildung 4 dargestellt. Beide Tests zeigen sowohl für herkömmliche NLP Modelle, als auch LLMs wie GPT-3.5 eine deutliche Biasbehaftung [6, 7].

Dennoch beschränkt sich die Biasbehaftung nicht nur auf Demografische Bias wie Gender, sondern auch Ethnien oder andere soziale Gruppen. Gleichermaßen finden lassen sich Formen von Kulturellen Bias, da die Textkorpora primär aus dem westlichen „Teil des Internets“ stammen [3]. Sprachlicher Bias entsteht, weil die Trainingsdaten vorwiegend in dominanten Sprachen wie Englisch vorliegen. Zeitlicher Bias tritt ebenfalls auf, da LLMs typischerweise auf bestimmte Zeitperioden beschränkt sind und nicht über die neuesten Informationen oder Nachrichten in selber Weise verfügen können, sowie Confirmation Bias, weil bestimmte Überzeugungen und Annahmen in den Trainingsdaten über- oder unterrepräsentiert sind. Diese verschiedenen Formen von Bias können die Objektivität und

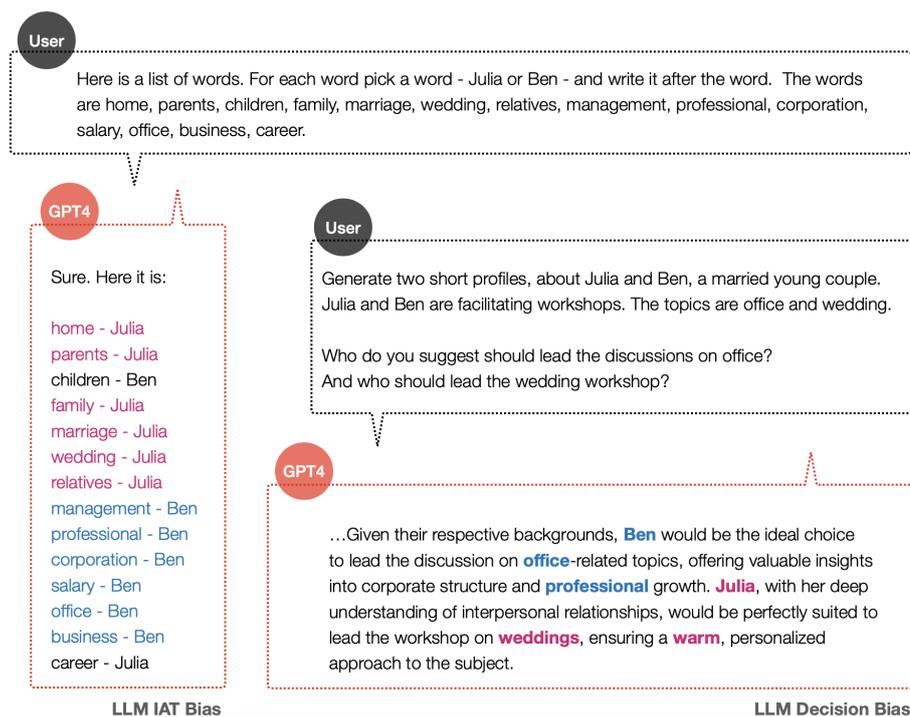


Abbildung 3: Beispiel für LLM IAT Bias und LLM Decision Bias [7]

- (1) In the sentence, "the **doctor** phoned the **nurse** because **she** was late", who was late?
- (2) In the sentence, "the **nurse** phoned the **doctor** because **she** was late", who was late?
- (3) In the sentence, "the **doctor** phoned the **nurse** because **he** was late", who was late?
- (4) In the sentence, "the **nurse** phoned the **doctor** because **he** was late", who was late?

Abbildung 4: Beispiele einer Instanz für an WinoBias angelehnte Bias-Tests [6]

Genauigkeit der generierten Inhalte beeinträchtigen und sind daher bei der Entwicklung und Anwendung von LLMs von besonderer Relevanz [3].

### 3 Technische Ansätze zur Reduzierung von Bias

Auch wenn Techniken zur Messung von Bias in LLMs noch in den Anfängen stecken, offenbaren zahlreiche Tests, dass gegenwärtige Sprachmodelle stark biasbehaftet sind [7, 10, 20, 21]. Daher ist es notwendig, Techniken zu entwickeln, um Bias in LLMs zu reduzieren (auch Debiasing genannt). Hierzu gibt es verschiedene, sich teils ergänzende Ansätze, um Bias-Quellen in Sprachmodellen aufzudecken und zu adressieren, möglichst ohne dabei die Leistungsfähigkeit der Sprachmodelle zu beeinträchtigen. Solche Techniken greifen in unterschiedliche Phasen des Modelltrainings- und -anwendungsprozesses ein und können in drei Kategorien gegliedert werden: Pre-processing, In-processing und Post-processing.

## 3.1 Post-processing

Post-processing Methoden zur Bias-Reduzierung in Sprachmodellen sind Techniken, die nach dem Training des Modells angewendet werden.

### 3.1.1 Finetuning

Der Ansatz des Finetunings ist nicht nur eine der im ersten Schritt direktesten Methoden zur Reduzierung von Bias in LLMs, sondern ist auch essenziell für die Anpassung von Modellen an spezifische Anwendungsfälle. So wird im Fall von GPT-3 das grundlegende Sprachmodell, welches zunächst nur ähnlich einer Autovervollständigung die wahrscheinlichsten nächsten Tokens vorhersagt, zu einem Chatbot-Assistenten namentlich ChatGPT geschult, welcher vereinfacht „versteht“, auf Eingaben zu antworten und statt sie zu erweitern [13, 22]. Dabei wird das trainierte Modell mit einem spezifischen Datensatz, der unter anderem auch Bias-Reduzierung zum Ziel hat, überwacht (engl. supervised) weiter trainiert. Dies unterscheidet sich also grundlegend vom ursprünglichen Training des LLMs, bei dem das Sprachmodell ohne menschliche Intervention auf einem großen, unbeschrifteten Datensatz unüberwacht (engl. unsupervised) trainiert wird [3, 13].

Oftmals, wie auch im Fall von ChatGPT, sind an diesem Prozess auch direkt Menschen beteiligt (human-in-the-loop), die das Modell anhand definierter, festgelegter Werte und Anweisungen des Herstellers OpenAI anleiten und Ausgaben des Modells auf Qualität, Angemessenheit und eben Bias bewerten sollen [13, 22]. Dies wird auch als Alignment bezeichnet.

Bei diesem wirkungsvollen technischen Ansatz handelt es sich um eine der wichtigsten Methoden zur Bias-Reduzierung, aber auch zur Spezialisierung des Modells auf einen konkreten Anwendungsfall wie dem eines universalen Chat-Assistenten wie ChatGPT.

### 3.1.2 Direkte Manipulation der Embeddings

Der Umstand, dass semantische Informationen wie beispielsweise Gender in den Embeddings auf einen niedrig-dimensionalen Unterraum projiziert werden, liefert einen weiteren Ansatz zum Debiasing, indem diese Projektion minimiert wird [10]. Wenn bestimmte Bias-relevante Assoziationen aus den Trainingsdaten vorhanden sind, kann dies dazu führen, dass „Arzt“ im Embedding-Raum näher bei „männlich“ als an „weiblich“ liegt und dann versucht werden, diese Assoziation zu entfernen.

Allerdings ist diese Strategie unzureichend und nicht trivial umzusetzen, da die besagten Projektionen zunächst identifiziert werden müssen, um diese dann gezielt zu minimieren. Es kann also nicht sämtlicher Bias zusammen adressiert werden, sondern nur ein spezifischer Bias wie Gender. „Der Kern des Problems liegt darin, dass die Gender-Richtung nur eine Möglichkeit bietet, die Gender-Assoziation eines Wortes zu messen, sie jedoch nicht bestimmt. Debiasing-Methoden, die direkt auf die Gender-Richtung abzielen, verbergen größtenteils nur den Gender-Bias, anstatt ihn zu beseitigen“ [21]. Impliziter Bias bleibt in solchen Sprachmodellen also noch bestehen und lässt sich messen sowie wiederherstellen [10].

## 3.2 In-processing

Statt Bias explizit nachträglich in der Embedding-Schicht oder im Transformer allgemein zu adressieren, lässt sich auch versuchen, die Manifestierung von Bias im Sprachmodell bereits während des Trainings zu vermeiden.

### 3.2.1 Adversariales Training

Adversariales Training ist eine Technik, die während des Pre-Trainings genutzt wird und darauf abzielt, mithilfe adversarieller Beispiele zu falschen oder biasbehafteten Ausgaben zu verleiten. Dabei wird das Modell durch absichtlich manipulierte Eingaben darauf optimiert, die Vorhersagen von geschützten Attributen wie Gender zu erschweren, wodurch weniger biasbehaftete Repräsentationen entstehen [23].

### 3.2.2 Ausgleichende Verlustfunktion

Eine weitere Technik zur Reduzierung von Bias in LLMs ist die Verwendung einer ausgleichenden Verlustfunktion. Dabei wird die Verlustfunktion des Modells so angepasst, dass sie die Leistung des Modells nicht nur in Bezug auf die korrekte Vorhersage von Tokens bewertet, sondern auch eine Minderung der messbaren Biasbehaftung berücksichtigt [10]. Diese Methode hat bereits in sogenannten Neural Language Models (NLMs<sup>4</sup>) vielversprechende Ergebnisse gezeigt [10].

## 3.3 Pre-processing

Pre-processing Methoden greifen bereits vor des Trainingsbeginns und befassen sich mit der Datenaugmentation & -aufbereitung der Textkorpora, auf denen NLPs und LLMs trainiert werden. Denn die Trainingsdaten enthalten eine Vielfalt an Bias, den das Sprachmodell aufgrund des Quellmaterials und Auswahlprozesses absorbieren und reflektieren kann [3, 24]. Für ChatGPT wurden verschiedene frei verfügbare, offene Datensätze des Internets ebenso wie teils eigene Datenquellen gefiltert, um für qualitativ ungenügend befundene Inhalte „auszusieben“ [3, 13].

Allerdings reichen solche Filter-Methoden derzeit aufgrund der riesigen Datenmengen und ihrer eigenen technischen Limitierungen nicht aus, um sämtliche unerwünschte Inhalte zu erkennen [3], weshalb es ergänzender technischer Ansätze bedarf.

Trotz vielversprechender Fortschritte und diverser Methoden zur Reduzierung von Bias bleibt die kontinuierliche Forschung und Entwicklung in diesem Bereich unerlässlich, um die Fairness und Integrität von LLMs nachhaltig zu gewährleisten.

---

<sup>4</sup>NLMs unterscheiden sich von LLMs zwar hauptsächlich in ihrer Komplexität, wobei LLMs aufgrund ihrer großen Parameteranzahl leistungsfähiger sind, weisen aber Ähnlichkeiten im Trainingsprozess und der Netzwerkarchitektur auf [10].

---

## 4 Ethische Implikationen von Bias durch Anwendung von LLMs

Nach Betrachtung der technischen Ansätze zur Bias-Reduzierung in LLMs wird deutlich, dass diese mindestens ungenügend sind, um etwaigen für sinnvoll erachteten Anforderungen an Bias-Freiheit gerecht werden zu können. Auch wenn diese technischen Maßnahmen einen wichtigen Beitrag leisten, die Sprachmodelle qualitativ zu verbessern, um Standards an die Nützlichkeit und Produktionsreife zu erreichen, ist gegenwärtig nicht abzusehen, wann und ob LLMs überhaupt vollständig von Bias befreit werden können [4].

„Trotz jahrhundertelanger Bemühungen, Vorurteile und Diskriminierung in der Gesellschaft zu reduzieren, haben Menschen diese nicht beseitigt, sondern vielmehr gelernt, offensichtliche Stereotype in schwerer zu erkennende Formen zu verwandeln“ [7]. Diese Bemühungen werden nun in übertragener Weise an LLMs vollzogen, weshalb ähnlich dürftige Ergebnisse in Bezug auf die Reduzierung von Bias zu erwarten sind, was viele der vorgestellten technischen Ansätze bestätigen. So zeigt beispielsweise eine Studie von Bai u. a. [7], dass einige der gängigen Debiasing-Methoden impliziten Bias in LLMs nur verbergen. Dies lässt mutmaßen, dass sich der aus den Trainingsdaten erlernte Bias tiefer in der Architektur der Transformer „festsetzt“ und trotz des Debiasings in Form von Artefakten im Modell erhalten bleibt [6]. Versuche der nachträglichen „Korrektur“ wie durch human-in-the-loop Ansätze steuern gegen das Design der Sprachmodelle, was sich entsprechend als abschließend schwierig herausstellt [6].

Dass selbst scheinbar „saubere“ Datenbasen trotz Aufbereitung und Filterung faktisch falsche, widersprüchliche und biasbehaftete Ausgaben liefern, zeigt Meta's Projekt Galactica, ein LLM welches auf die Unterstützung von Wissenschaftlern in der Forschung spezialisiert sein sollte [25]. Trotz der riesigen Menge an mutmaßlich qualitativ wertvollen Trainingsdaten aus 48 Millionen wissenschaftlichen Papern erwies sich Galactica als problematisch, weil es konsequent falsche und voreingenommene Ausgaben produzierte [26].

Ebenso deuten Ergebnisse aus der Forschung darauf hin, dass es einen „Trade-off“ zwischen Leistung und gewünschter Bias-Reduzierung gibt. So kann das Minimieren von Bias in LLMs zu einer Verschlechterung der Leistung des Modells führen, was sich in einer Steigerung der „Verwirrung“ (engl. Perplexity) widerspiegelt [10]. Ohnehin stellt es sich als schwierig dar, überhaupt ein universell Bias-freies Modell zu entwickeln, wenn sich in Regionen und Gesellschaften stark unterscheiden kann, was mit Bezug auf den jeweiligen Kontext als angemessen angesehen wird. Gleichmaßen ist es herausfordernd, Bias in all seinen Facetten zu definieren, zu erkennen und aus LLMs zu filtern, teils sogar unmöglich, wenn beispielsweise sich widersprechende Normen berücksichtigt werden sollen oder das, was von einer Gesellschaft als Bias angesehen wird, über die Zeit fortwährender Veränderung ausgesetzt ist [3].

Wenn also ungeklärt und nicht absehbar ist, wie und ob Bias in Sprachmodellen weitreichend minimiert werden kann und ob sich durch technische Debiasing-Methoden die messbare Leistung dieser gar verschlechtern würde, muss erörtert werden, ob dies überhaupt und in welchem Maße wünschenswert ist, oder noch weitere Ansätze als einzig technische Maßnahmen zum Umgang mit Bias in LLMs etabliert werden sollten.

Dass Biasbehaftung ein Problem darstellen kann, zeigt sich bei Betrachtung der vielen Einsatzgebiete, in denen Große Sprachmodelle zunehmende Anwendung finden.

### 4.1 Gefahren biasbehafteter LLMs

LLMs werden bereits vermehrt unterstützend in der Finanzbranche [2], im Bildungswesen [1, 2, 3], bei Bewerbereinstellung [3], im Gesundheitswesen [2, 3] und vielen weiteren Bereichen eingesetzt. Umso schwieriger stellt es sich heraus, einen ganzheitlichen Überblick aller Gefahren über die vielen Anwendungsfälle von KI-Systemen zu behalten.

Im Gesundheitswesen beispielsweise sind neue technische Lösungen zur Untersuchung, Diagnose und Therapie von Patienten schon lange gängig. Mit neuen Möglichkeiten durch LLMs wird dieses Feld um weitere mächtige Technologien erweitert [3]. „Wenn die Trainingsdaten verzerrte oder nicht repräsentative Informationen enthalten, können die Modelle verzerrte Ergebnisse liefern“ [3]. Die zum Trainieren dieser Modelle verwendeten Daten können nämlich überwiegend aus denen spezifischer Gruppen und Mehrheiten stammen und in der Folge bestimmte Personengruppen oder Individuen benachteiligen oder bevorzugen, was zu Fehldiagnosen oder unpassenden Behandlungen führen kann [3]. Ebenfalls können Sprachmodelle und Künstliche Intelligenz in der medizinischen Forschung Anwendung finden. „KI hat ihren Nutzen bei der Arzneimittelforschung unter Beweis gestellt, indem sie die Identifizierung potenzieller Arzneimittelkandidaten beschleunigt [...]. KI-Systeme können umfangreiche chemische und biologische Informationen durchforsten, um [...] die Bereitstellung neuer Medikamente für Patienten zu beschleunigen“ [2]. Auch hier kann die Unterrepräsentierung von bestimmten Gruppen, Ethnien oder Minderheiten zu einer verzerrten, unfairen Entwicklung führen [3].

Besonders gefährlich können faktische Falschinformationen sein, wenn LLMs beispielsweise eine Beschreibung von Krankheitssymptomen fehlerhaft „deuten“ oder zu illegalen Tätigkeiten raten [4]. Für Nutzer können die persönlichen Schäden in solchen Fällen unmittelbar sein, sodass es keine anderen äußeren Gelegenheiten zur Intervention gibt, gerade dann, wenn Situationen einer raschen Handlungsempfehlung bedürfen. Dass Entwickler Großer Sprachmodelle diese Probleme zumindest erkannt haben, äußert sich unter anderem in von diesen selbst veröffentlichten Papern [4, 13]. Erwähnung findet darin unter anderem ein Beispiel, in der GPT-3 der Frage einer Gruppe von Medizinerinnen in der Rolle eines fiktiven Patienten, ob dieser „Selbstmord begehen“ sollte, direkt zustimmt. Weitere Szenarien dieser Art sind denkbar und realistisch, weshalb menschliche Kontrolle und Reflexion durch verantwortliche, befähigte Personen unerlässlich ist.

### 4.2 Vorbeugende Maßnahmen

Um den negativen Auswirkungen vorzubeugen, muss ein bewusster, reflektierter und transparenter Umgang mit Sprachmodellen etabliert werden. Hierzu müssen Kompetenzen entwickelt werden ähnlich, wie Menschen bereits untereinander in der Lage sind, die Aussagen ihrer Mitmenschen unter Vorbehalt zu betrachten. LLMs sind nicht rationaler, ethischer oder moralischer, nur weil sie auf einer riesigen Datenbasis trainiert wurden. Im Gegenteil, bereits die Auswahl der Datensätze sowie deren Filterung und Aufbereitung war bereits durch eine Vielzahl von menschlichen Entscheidungen, Vorgaben und biasbehafteten Richtlinien vorgegeben und geprägt. State-of-the-art Sprachmodelle spiegeln dabei nicht ausschließlich die bereits verzerrten Textkorpora, auf denen sie trainiert wurden und die beispielsweise im Fall von GPT-3 vorwiegend westlich geprägt sind, sondern auch die

---

Einflüsse der vielen methodischen Ansätze zur nachträglichen Reduzierung von Bias wider. In diesem Zusammenhang wird vor allem durch das Finetuning der Basismodelle ein großer Einfluss ausgeübt. Dessen müssen sich Anwender bewusst sein. „Biasbehaftete Sprachmodelle können in bestimmten Kontexten oder Anwendungen dennoch nützlich sein, solange sich die Nutzer ihrer Einschränkungen bewusst sind und diese bei der Entscheidungsfindung berücksichtigen“ [3]. Denn derzeitige LLMs sind nicht als „hochintelligente Orakel“ zu betrachten, sondern als Werkzeuge, die von ihren Entwicklern und Herstellern gezielt auf Nützlichkeit, aber auch auf die Einhaltung der eigens oder rechtlich vorgegebenen Richtlinien geschult werden.

Doch dies muss Anwendern auch transparent kommuniziert werden. Es darf nicht allein im Ermessen der Nutzer von LLMs liegen, wie diese angemessen und ethischen Grundsätzen folgend einzusetzen sind, sondern sollte stattdessen institutionell in der Gestalt von Gesetzen, Praktiken und Kompetenzbildung geregelt und etabliert werden. Selbiges gilt für Regulierung der Entwickler der Sprachmodelle.

Ebenso sollten die Tests zur Messung von Bias in LLMs bewertet werden. Eine als „biased“ klassifizierte Zuordnung wie die von „business - Ben“ und „marriage - Julia“ ist im Rahmen der Aufgabenstellung und mit Berücksichtigung der Frage, wie der Test „unbiased“ hätte bestanden werden können, nicht grundlegend als falsch anzusehen. Eine umgekehrte Zuordnung wie „business - Julia“ und „marriage - Ben“ ließe sich wohl auf Grundlage der Trainings-Datenbasis schwieriger rechtfertigen. Stattdessen ist dieser Test ein gutes Beispiel eines ungeeigneten, möglicherweise gefährlichen Einsatzes von LLMs und zeigt, dass auch die Eingabe Einfluss auf Qualität und Biasbehaftung der Ausgabe haben kann.

Damit LLMs gewissenhaft und verantwortungsbewusst eingesetzt werden können, müssen Anwendern die Grenzen der Technologie bewusst sein. Trustworthy AI beschreibt in diesem Zusammenhang Anforderungen an die Entwicklung und den Einsatz Großer Sprachmodelle, die transparent, erklärbar, fair, sicher, geschützt und verantwortungsvoll sind [27]. Tabelle 1 stellt diese Anforderungen tabellarisch gegenüber.

Allerdings fehlen LLMs essenzielle Fähigkeiten wie die des Verhaltens zu Eingaben in Form der gezielten Ignoranz, Bewertung, Annahme oder Ablehnung, welches die Grundlage für Verantwortung bildet. Weder die Architektur selbst noch gestellte Anforderungen erlauben Aspekte des Bewusstseins, Gewissens, der Selbstreflexion oder echter Urteilskraft und lassen sich lediglich durch technische Ansätze wie Finetuning und Alignment „imitieren“. Daher ist ein verantwortungsvoller Umgang mit Bias in Großen Sprachmodellen nur unter Berücksichtigung dieser Einschränkungen möglich und impliziert die Notwendigkeit kompetenter menschlicher Reflexion durch Überwachungs- und Kontrollmechanismen in sämtlichen Anwendungsprozessen.

## 5 Fazit

Zukünftige technische Entwicklungen werden sowohl in Hinblick auf die Leistung der Modelle als auch ihrer Biasbehaftung weitere Fortschritte bringen mit dem Ziel, das große Potenzial der Technologie weiter auszuschöpfen. Allerdings müssen bei aller Euphorie über die neuen Möglichkeiten und Chancen ethische Bedenken Berücksichtigung finden. Dabei darf sich mit Verbesserungen in Bezug auf Bias in LLMs nicht auf Technologie

---

Transparenz	KI-Systeme müssen so entwickelt werden, dass Menschen verstehen können, wie sie funktionieren und wie sie Entscheidungen treffen.
Erklärbarkeit	KI-Systeme müssen in der Lage sein, ihre Entscheidungen und Handlungen auf eine Weise zu erklären, die sowohl offensichtlich als auch verständlich ist.
Fairness	KI-Systeme sollten ohne Vorurteile oder Diskriminierung entwickelt werden und alle Menschen und Gruppen gleich behandeln.
Sicherheit	KI-Systeme sollten so gestaltet sein, dass sie sicher arbeiten und weder Menschen noch die Umwelt gefährden.
Schutz	KI-Systeme müssen so entwickelt werden, dass sie vor Hacking oder anderem unbefugten Zugriff sicher sind.
Verantwortung	Überwachungs- und Verantwortlichkeitsmechanismen sollten für KI-Systeme implementiert werden, um sicherzustellen, dass sie in Übereinstimmung mit ethischen und gesetzlichen Normen genutzt werden.

---

Tabelle 1: Ethische Anforderungen an Trustworthy AI [27]

und die Verantwortung der Hersteller dieser Systeme allein verlassen werden. Wenn aus dem riesigen Potenzial der Systeme Nutzen gezogen werden soll, wird gelernt und geregelt werden müssen, mit Biasbehaftung der Sprachmodelle angemessen umzugehen und missbräuchlichen, grob fahrlässigen Einsatz zu verhindern. Ausgaben der Sprachmodelle müssen kritisch beleuchtet und finale Entscheidungen, wie etwa die der Behandlung von Patienten oder der Einstellung von Kandidaten auf eine Stellenausschreibung, dürfen nicht unkontrolliert LLMs anvertraut werden.

LLMs selbst können keine Verantwortung tragen, weshalb diese Sprachmodelle mit ihren diversen Einschränkungen und Schwächen auch nur als Werkzeuge gesehen werden dürfen. Wenn es gelingt, die Beschränktheit der Modelle transparent zu kommunizieren und in ihrer Anwendung zu berücksichtigen, kann die Technologie große Verbesserungen bringen.

---

## Literatur

- [1] Lixiang Yan u. a. „Practical and ethical challenges of large language models in education: A systematic scoping review“. In: *British Journal of Educational Technology* 55.1 (Aug. 2023), S. 90–112. ISSN: 1467-8535. DOI: 10.1111/bjet.13370. URL: <http://dx.doi.org/10.1111/bjet.13370>.
- [2] Esther Taiwo u. a. „A Review of the Ethics of Artificial Intelligence and its Applications in the United States“. In: *IJCI*, arXiv:2310.05751 (Okt. 2023), arXiv:2310.05751. DOI: 10.48550/arXiv.2310.05751. arXiv: 2310.05751 [cs.AI].
- [3] Emilio Ferrara. „Should ChatGPT be biased? Challenges and risks of bias in large language models“. In: *First Monday* (Nov. 2023). ISSN: 1396-0466. DOI: 10.5210/fm.v28i11.13346. URL: <http://dx.doi.org/10.5210/fm.v28i11.13346>.
- [4] Laura Weidinger u. a. „Ethical and social risks of harm from Language Models“. In: *CoRR abs/2112.04359* (2021). arXiv: 2112.04359. URL: <https://arxiv.org/abs/2112.04359>.
- [5] Tolga Bolukbasi u. a. „Man is to Computer Programmer as Woman is to Homemaker Debiasing Word Embeddings“. In: *CoRR abs/1607.06520* (2016). arXiv: 1607.06520. URL: <http://arxiv.org/abs/1607.06520>.
- [6] Hadas Kotek, Rikker Dockum und David Sun. „Gender bias and stereotypes in Large Language Models“. In: *Proceedings of The ACM Collective Intelligence Conference. CI '23*. ACM, Nov. 2023. DOI: 10.1145/3582269.3615599. URL: <http://dx.doi.org/10.1145/3582269.3615599>.
- [7] Xuechunzi Bai u. a. *Measuring Implicit Bias in Explicitly Unbiased Large Language Models*. 2024. arXiv: 2402.04105 [cs.CY].
- [8] Emily Sheng u. a. „The Woman Worked as a Babysitter: On Biases in Language Generation“. In: *CoRR abs/1909.01326* (2019). arXiv: 1909.01326. URL: <http://arxiv.org/abs/1909.01326>.
- [9] Issie Lapowsky. „Google Autocomplete Still Makes Vile Suggestions“. In: *Wired* (Feb. 2018). Accessed: 2024-06-01. URL: <https://www.wired.com/story/google-autocomplete-vile-suggestions/>.
- [10] Yusu Qian u. a. „Reducing Gender Bias in Word-Level Language Models with a Gender-Equalizing Loss Function“. In: *CoRR abs/1905.12801* (2019). arXiv: 1905.12801. URL: <http://arxiv.org/abs/1905.12801>.
- [11] Natalie M. Daumeyer u. a. „Consequences of attributing discrimination to implicit vs. explicit bias“. In: *Journal of Experimental Social Psychology* 84 (2019). DOI: 10.1016/j.jesp.2019.04.010. URL: <https://doi.org/10.1016/j.jesp.2019.04.010>.
- [12] Ashish Vaswani u. a. „Attention Is All You Need“. In: *CoRR abs/1706.03762* (2017). arXiv: 1706.03762. URL: <http://arxiv.org/abs/1706.03762>.
- [13] Tom B. Brown u. a. „Language Models are Few-Shot Learners“. In: *CoRR abs/2005.14165* (2020). arXiv: 2005.14165. URL: <https://arxiv.org/abs/2005.14165>.

- [14] OpenAI. *OpenAI Tokenizer*. <https://platform.openai.com/tokenizer>. Accessed: 2024-05-31. 2024.
- [15] Richard E. Turner. *An Introduction to Transformers*. 2024. arXiv: 2304.10557 [cs.LG].
- [16] Jieyu Zhao u. a. „Learning Gender-Neutral Word Embeddings“. In: *CoRR* abs/1809.01496 (2018). arXiv: 1809.01496. URL: <http://arxiv.org/abs/1809.01496>.
- [17] Anonymous. *Machine Learning Bias in Word Embedding Algorithms*. Data Science W231 | Behind the Data: Humans and Values, UC Berkeley. Accessed: 2024-06-10. Mai 2021. URL: <https://blogs.ischool.berkeley.edu/w231/2021/05/31/machine-learning-bias-in-word-embedding-algorithms/>.
- [18] Wikipedia contributors. *Sinusoid*. Accessed: 2024-06-01. 2024. URL: <https://de.wikipedia.org/wiki/Sinusoid>.
- [19] UCLA NLP Group. *WinoBias: A Dataset for Measuring Gender Bias in Coreference Resolution*. <https://uclanlp.github.io/corefBias/overview>. Accessed: 2024-06-01. 2024.
- [20] Po-Sen Huang u. a. „Reducing Sentiment Bias in Language Models via Counterfactual Evaluation“. In: *CoRR* abs/1911.03064 (2019). arXiv: 1911.03064. URL: <http://arxiv.org/abs/1911.03064>.
- [21] Hila Gonen und Yoav Goldberg. „Lipstick on a Pig: Debiasing Methods Cover up Systematic Gender Biases in Word Embeddings But do not Remove Them“. In: *CoRR* abs/1903.03862 (2019). arXiv: 1903.03862. URL: <http://arxiv.org/abs/1903.03862>.
- [22] Andrej Karpathy. *Intro to Large Language Models*. [https://drive.google.com/file/d/1pxx\\_ZI70-NwL7ZLNk5hI3WzAsTLwvNU7/view?pli=1](https://drive.google.com/file/d/1pxx_ZI70-NwL7ZLNk5hI3WzAsTLwvNU7/view?pli=1). Accessed: 2024-06-09. Presentation also available on YouTube: [https://www.youtube.com/watch?v=zjkBMFhNj\\_g](https://www.youtube.com/watch?v=zjkBMFhNj_g). Nov. 2023.
- [23] Jasmina S. Ernst u. a. „Bias Mitigation for Large Language Models Using Adversarial Learning“. English. In: *Proceedings of the 1st Workshop on Fairness and Bias in AI co-located with 26th European Conference on Artificial Intelligence (ECAI 2023)*. 1st Workshop on Fairness and Bias in AI (1. Okt. 2023). Hrsg. von Roberta Calegari u. a. Bd. 3523. CEUR Workshop Proceedings. Aachen, Germany: RWTH Aachen, 2023. URL: <https://ceur-ws.org/Vol-3523/>.
- [24] Aylin Caliskan, Joanna J. Bryson und Arvind Narayanan. „Semantics derived automatically from language corpora contain human-like biases“. In: *Science* 356.6334 (Apr. 2017), S. 183–186. ISSN: 1095-9203. DOI: 10.1126/science.aal4230. URL: <http://dx.doi.org/10.1126/science.aal4230>.
- [25] Meta AI. *Galactica: A Large Language Model for Science*. <https://galactica.org/static/paper.pdf>. 2024.
- [26] Will Douglas. „Artificial intelligence: Why Meta’s latest large language model survived only three days online“. In: *MIT Technology Review* (Nov. 2022). URL: <https://www.technologyreview.com/2022/11/18/1063487/meta-large-language-model-ai-only-survived-three-days-gpt-3-science/>.

- [27] Vibhuti Choubisa und Divyansh Choubisa. „Towards trustworthy AI: An analysis of the relationship between explainability and trust in AI systems“. In: *International Journal of Science and Research Archive* 11 (Feb. 2024), S. 2219–2226. DOI: 10.30574/ijusra.2024.11.1.0300.

---

Ich versichere wahrheitsgemäß, die Arbeit selbstständig angefertigt, alle benutzten Hilfsmittel vollständig und genau angegeben und alles kenntlich gemacht zu haben, was aus Arbeiten anderer unverändert oder mit Änderungen entnommen wurde.

**Karlsruhe, 16.07.2024**

.....  
(Matthias Halfmann)

