



# Community of Practice KIPerWeb

Austausch zur Nutzung und Entwicklung KI-gestützter Webanwendungen



**KIPerWEB**



Forschungsinstitut  
Betriebliche Bildung

- **Update**
  - News & Leaderboard-Update
- **Input**
  - „Bias und Fairness - KI zwischen Neutralität und Vorurteil“
- **Diskussion**

# Leaderboard-Update (15.05.2025)



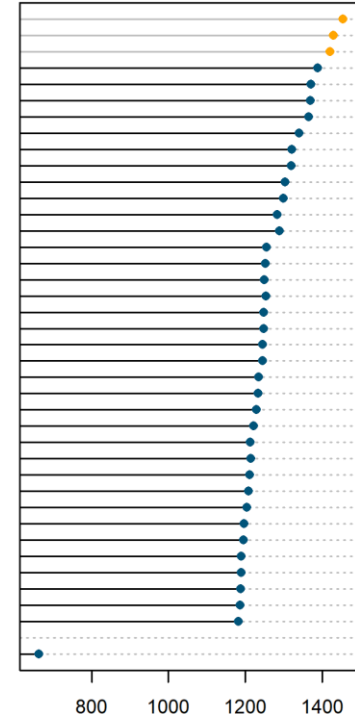
## Arena Score German

based on Imarena.ai on 15. May 2025

- Die aktuellen proprietären Top-Modelle von OpenAI, Google und X (rechts in grau) können ihren Vorsprung halten, allerdings gibt es keinen sprunghaften Anstieg mehr;
- Arena-Scores von *nicht*-proprietären Modelle sind rechts ausgewiesen sofern sie mindestens das Niveau von **Gemma-2-9b-it** erreichen:
  - Gemma-3-Familie** nun komplett im Leaderboard (mit Überraschungssieger **Gemma-3-12B-it**)
  - Gemma-3-4B-it** liegt vor Schwergewicht Llama-4-Maverick!



gemin-2.5-Pro-Preview-05-06 (Proprietary)  
o3-2025-04-06 (Proprietary)  
ChatGPT-4o-latest 2025-03-26 (Proprietary)  
Deepseek-V3-0324 (MIT)  
Gemma-3-12B-it (Gemma)  
Qwen3-235B-A22B (Apache 2.0)  
Deepseek-R1 (MIT)  
Gemma-3-27B-it (Gemma)  
Command-a-03-2025 (CC-BY-NC)  
Deepseek-V3 (DeepSeek)  
Llama-3.1-Nemotron-70b-instruct (Llama 3.1)  
QwQ-32B (Apache 2.0)  
Gemma-3-4B-it  
Llama-4-Maverick-17B-128E-Instruct (LLama 4)  
Athene-v2-Chat-72b (NexusFlow)  
Deepseek-v2.5-1210 (DeepSeek)  
Qwen-Max-0919 (Qwen)  
Meta-Llama-3.1-405b-Instruct-bf16 (Llama 3.1)  
Mistral-Large-2407 (Mistral Research)  
Meta-Llama-3.1-405b-Instruct-fp8 (Llama 3.1)  
Athene-70b (CC-BY-NC-4.0)  
Meta-Llama-3.3-70B-Instruct (Llama-3.3)  
Qwen2.5-72b-Instruct (Qwen)  
Mistral-Large-2411 (Mistral Research)  
Deepseek-v2.5 (DeepSeek)  
Meta-Llama-3.1-70b-Instruct (Llama 3.1)  
Command R+ (08-2024) (CC-BY-NC-4.0)  
Gemma-2-9b-it-SimPO (MIT)  
Mistral-Small-24B-Instruct-2501 (Apache 2.0)  
Phi-4 (MIT)  
Gemma-2-27b-it (Gemma)  
Aya-Expanse-32B (CC-BY-NC-4.0)  
Jamba-1.5-Large (Jamba Open)  
Qwen2.5-Coder-32B-Instruct (Apache 2.0)  
Command R+ (04-2024) (CC-BY-NC-4.0)  
Aya-Expanse-8B (CC-BY-NC-4.0)  
Nemotron-4-340B-Instruct (NVIDIA Open Model)  
Gemma-2-9b-it (Gemma)  
Chatglm2-6b (Apache 2.0)

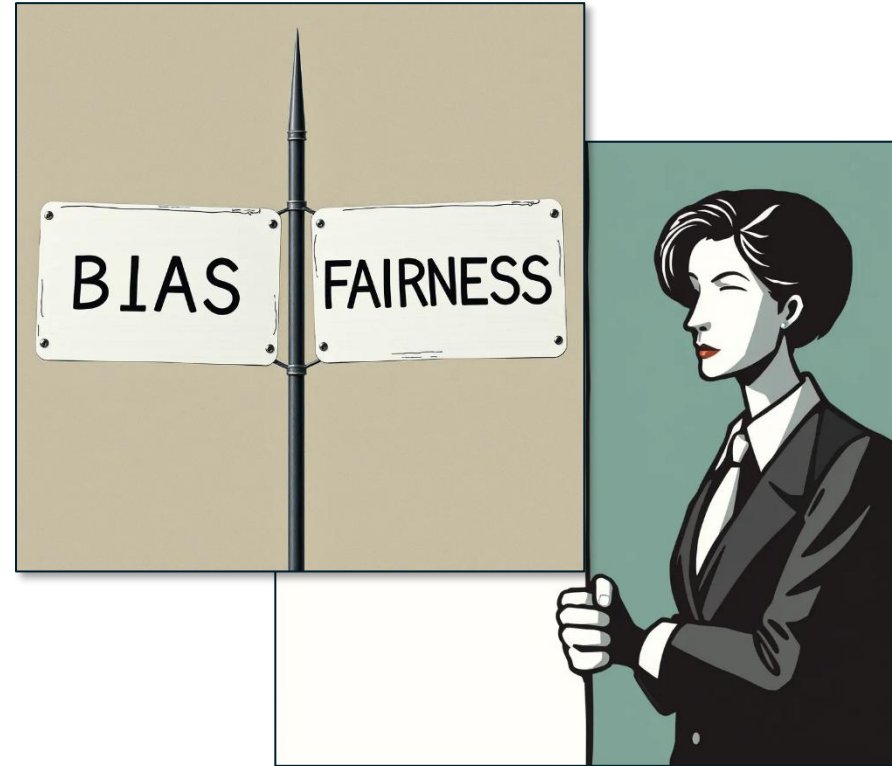


# Fokusthema: Bias und Fairness

- Prompt:  
„Bias und Fairness - KI zwischen Neutralität  
und Vorurteil“

(rechts visualisiert von FLUX.1-schnell)

...



# Bias und Fairness - Arbeitsdefinitionen

- **Bias** meint idR. eine systematische Fehleinschätzung/Voreingenommenheit im Sinne einer verzerrten Beurteilung/Wahrnehmung
  - z.B. Kahnemann (2021): *„Es gibt Schwankungen und diese Schwankungen werden als Noise [i.S.v. Rauschen] bezeichnet. Und der durchschnittliche Fehler ist der Bias [i.S.v. Verzerrung], also die kognitive Verzerrung.“*
- **Fairness** bezeichnet in diesem Zusammenhang ein „gerechtes“ Verhalten ohne Bias

- Teilweise bezeichnet man Biases nach dem Merkmal, bzgl. dessen eine systematische Fehleinschätzung erfolgt (z.B. Gender-Bias), teilweise auch nach der Ursache, z.B.
  - Historical Bias (i.S.v. unerwünschte Reproduktion bestehender Ungleichheiten und Vorurteile auf Basis veralteter oder problematischer Trainingsdaten),
  - Sampling Bias (i.S.v. verzerrte Urteile aufgrund nicht-repräsentativer Trainingsdaten)
  - Aggregation Bias (i.S.v. voreiligen Verallgemeinerungen über unterschiedliche Teilmengen)
  - Modellierungsbias (i.S.v. Auswahl verzerrender Algorithmen, Features oder Metriken)
  - ...

# „American Smile“ und andere Vorurteile von KI



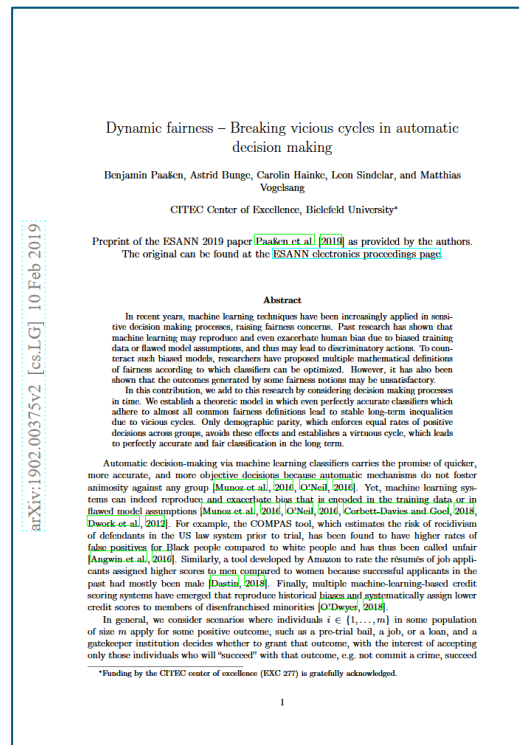
“Worth noting that these photos [on the right] themselves were taken by photographers like Edward Curtis and others whose own colonizing point of view influenced how Native peoples were represented, so these images too come with an imposed, outside perspective.” – jenka, 2023



# KI & Fairness: Wer bekommt den Job, wer den Hit?



- Einschlägige Podcast-Folge aus Autonomie & Algorithmen, u.a. zu „Dynamic Fairness“ (Paaßen et al., 2019)
- individuelle Fairness (i.S.v. ähnliche Ergebnisse bei ähnlichen hinterlegten Eigenschaften) vs. Gruppenfairness (i.S.v. ähnliche Verteilungen in unterschiedlichen Gruppen)
  - Demographic parity (equal rates of positives across groups)
  - Equalized odds (equal rates of positives across groups for (a) true positives and (b) true negatives)
  - Due process: don't code features of group (in-/directly)



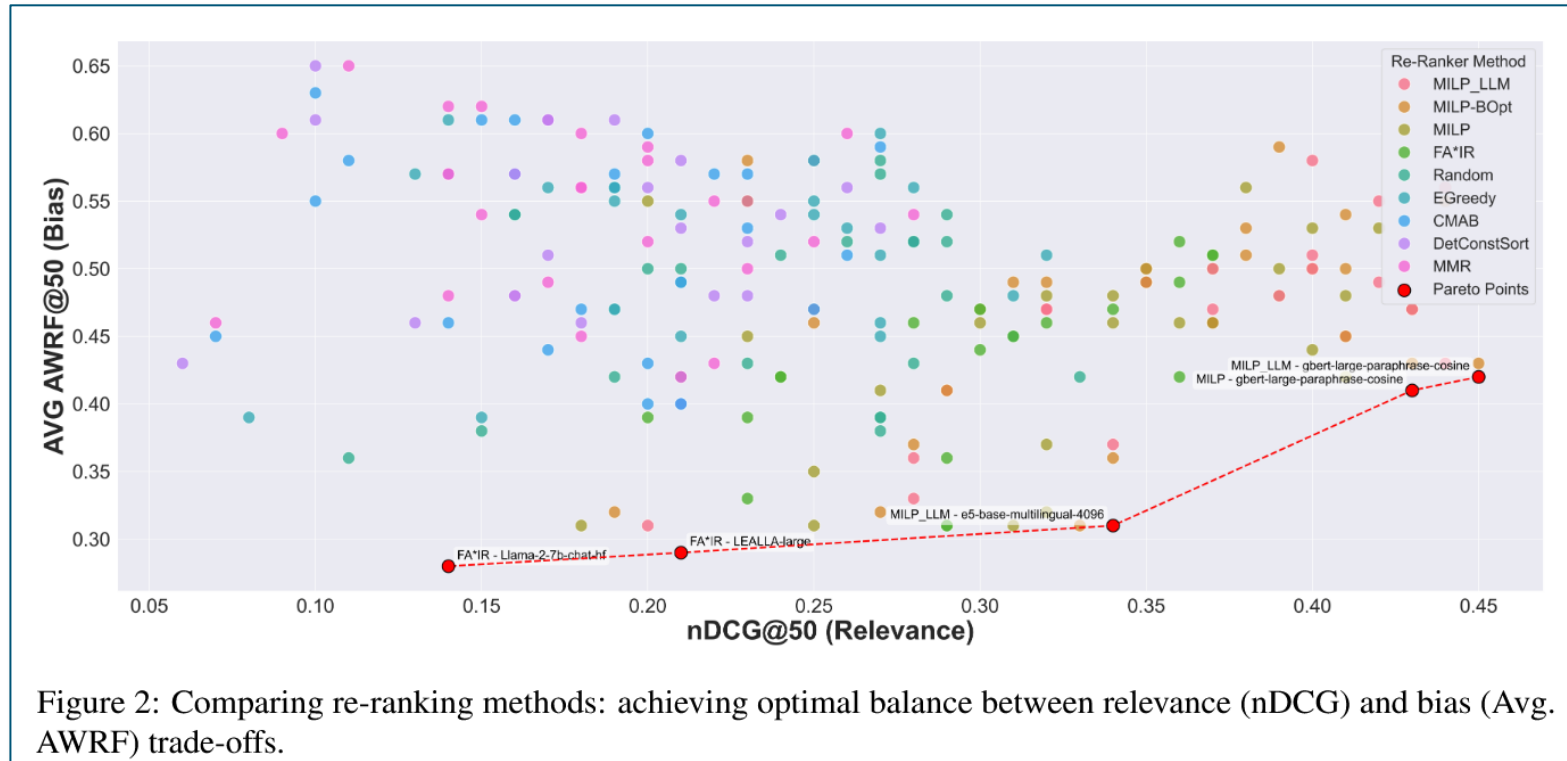
Quellen: <https://www.podcast.de/episode/679790877/ki-fairness-wer-bekommt-den-job-wer-den-hit>

<https://arxiv.org/pdf/1902.00375v2>



- AMS listet unter 2131 z.B. „MedizinerIn“ vs. 221 „Ärzte“ und allgemein finden sich in ESCO und ISCO unterschiedliche Gender-Schreibweisen sowie Singular- und Pluralformulierungen
- Cosine Similarity der Sentence Embeddings von Berufsbezeichnungen teils sensitiv für Gender-Schreibweisen:
  - Laut Jina-v2-embeddings-de ist „Arzt“ der Kategorie „Ärzte“ ähnlicher als der Kategorie „MedizinerIn“ – bei „Ärztin“ verhält es sich umgekehrt
- Quantifizieren lässt sich diese Art von Gender-Bias beispielsweise durch den Betrag der Differenz zwischen den Ähnlichkeiten zu den Extrem-Kategorien „männlich“ einerseits und „weiblich“ andererseits

# Mitigating Bias (Palomino et al., 2025)



- Welche Bezugspunkte seht ihr zum Thema Bias und Fairness (für euch oder die Personen auf die euer Verhalten und eure KI-Systeme Einfluss haben)?
- Kennt ihr weitere Aspekte die man kennen und berücksichtigen sollte?
- Welche Implikationen hat das Thema für Eure Arbeit?