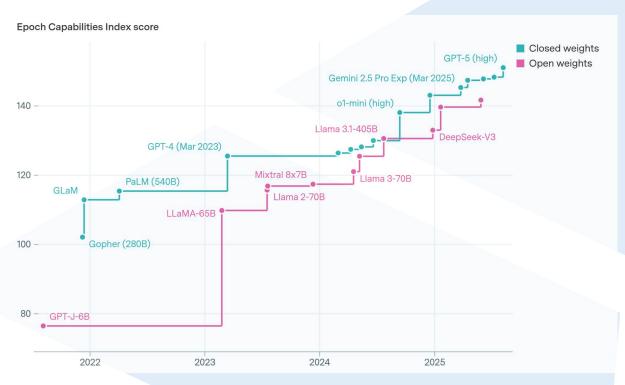




Der Abstand von open-weights zu closed-weights



Quelle: epoch.ai/data-insights/open-weights-vs-closed-weights-models



So was wie eine Agenda

Begriffe und Konzepte

KI-Benchmarking

Jetzt kommt das ChatGPT

aus der Schweiz

Perspektivenwechsel: LLM bauen

Swiss Al Initiative

Apertus

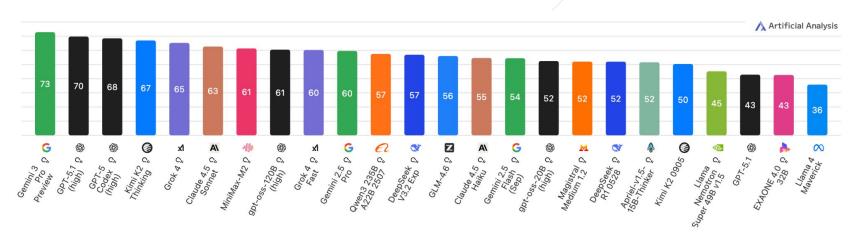
kick» start

Begriffe und Konzepte



The hot open-weights s*it ist grad Kimi 2 Thinking

Artificial Analysis Intelligence Index



Artificial Analysis Intelligence Index v3.0 incorporates 10 evaluations: MMLU-Pro, GPQA Diamond, Humanity's Last Exam, LiveCodeBench, SciCode, AIME 2025, IFBench, AA-LCR, Terminal-Bench Hard, τ^2 -Bench Telecom

Quelle: artificialanalysis.ai

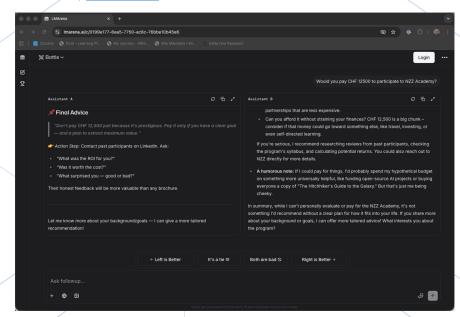


Wie misst man den "IQ" von KI-Systemen?

Synthetische Benchmarks (ground truth)

- Beispiel MMLU: Massive Multitask Language Understanding,
 - https://arxiv.org/pdf/2009.03300
 - multiple-choice questions spanning 57 distinct subject areas (history, philosophy, literature, psychology, economics, politics, mathematics, physics, engineering, law, medicine, accounting..)
- Beispiel ARC-AGI General Intelligence Benchmark, https://arcprize.org
 - Abstract and Reasoning Corpus
- Beispiel HLE: Humanity's Last Exam, https://agi.safe.ai
 - 2,500 challenging questions across over a hundred subjects

Bewertung durch Menschen (human preferences) z.B. LMArena, Imarena.ai





Eine Bemerkung zu "das ChatGPT aus der Schweiz"

KI-Modell versus KI-System

KI-Modell = Algorithmus

Ein KI-Modell ist der Algorithmus (auch Foundational Model genannt)

- Sprachmodelle (Language Models)
 - LLM (Large Language Model)
 - LAM (Language Action Models)
- Bildmodelle (Vision Models)
- Multimodale Modelle
- Recommender-Systeme

KI-System = Dienstleistung

Ein KI-System ist die vollständige, funktionale Anwendung, die ein oder mehrere KI-Modelle mit einer Benutzeroberfläche, Infrastruktur und anderen Komponenten integriert, um eine Dienstleistung bereitzustellen. **Analogie: Motor versus Auto**





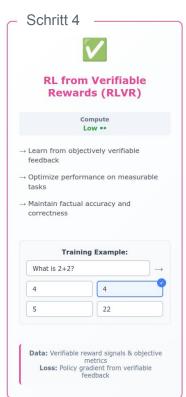
Perspektivenwechsel

Entwicklungsschritte eines LLMs









Quelle: Swiss {ai} Weeks. Tech Round on the LLM Effort of the Swiss Al Initiative. Imanol Schlag, Martin Jaggi. kick» start





Die Swiss Al Initiative wurde im Oktober 2023 gegründet

Nationale Forschungsinitiative unter gemeinsamer Leitung von ETHZ und EPFL

Über 10 akademische Einrichtungen, über 70 Professoren und über 800 Forscher

CHF 5 Mio. Startfinanzierung pro Jahr über 4 Jahre

~30 Millionen GPU-Stunden pro Jahr auf Alps

Die Initiative ist weltweit das grösste Open-Science-/Open-Source-Projekt für KI-Grundlagenmodelle

Mehr infos: swiss-ai.org



























Mission der Swiss Al Initiative

Develop *capabilities, knowhow, and talent* to build *trustworthy, aligned, and transparent* Al

Make these resources available for the benefit of Swiss society and global actors





Verteilung der Ressourcen über call for projects

Open call for small projects (~50k GPU hours): rolling reviews

Open call for large projects (>500k GPU hours): twice per year

Akzeptierte Calls

- Foundation Models for Many-Body Quantum Physics
- Swiss AI Large Foundation Models for Climate and Weather
- COMPL-AI++: An Evaluation Framework for LLM Safety and Compliance
- Clintextualization
- Integrating Pedagogy into Large Language Models for University-level Engineering Education in Switzerland
- Multimodal Multitask LLMs for Chemistry and Catalysis
- Foundation Models for Modeling the Complexity of Life and Molecules
- 3D Vision Language Model For Radiology
- Foundation Models Bridging 4D Capture and Dexterous Manipulation for Surgical Digital Twins and Beyond
- Learning To Scale Next-Level Generative Models
- Democratizing LLMs for Global Languages with Mixtures of Multilingual Experts
- FlexLM: Efficient Targeted LLM Compression
- From Training a Model to Raising a Model: A Path Toward Safe AI via Student–Teacher Pretraining
- Foundations of Multimodal Learning
- Multimodal World Foundation Models for Physical Al
- Building Switzerland's Sovereign Al Future: Advancing Swiss-Centric Foundation Models for a Trustworthy Al Ecosystem
- Virtual Patient Platform: Al-Powered Integration of Digital Pathology, Genomics, and Spatial Omics for Precision Oncology





Apertus: Eines der Projekte der Swiss Al Initiative



Fundamentals of foundation models

Prof. Yang, Prof. He, Prof. Zdeborova, Prof. Flammarion



Human-Al alignment

Prof. Ash, Prof. Gulcehre



LLM security, red teaming & privacy

Prof. Troncoso, Prof. Tramèr



Large-scale multi-modal models

Prof. Cotterell, Prof. Zamir



Tools & infrastructure for scaling

Prof. Klimovic, Prof. Falsafi



Advanced LLMs Prof. Bosselut, Prof. Jaggi,

Dr. Schlag



Foundation model for sciences

Prof. Brbic, Prof. Schwaller, Prof. Marinkovic



Foundation model for health

Prof. Rätsch, Prof. Salathé, Prof. Fellay



Foundation model for education

Prof. Käser, Prof. Sachan



Foundation model for ego-centric vision & robotics

Prof. Alahi, Prof. Pollefeys, Prof. Katzschmann



Foundation model for sustainability / climate

Prof. Mishra, Prof. Schemm, Prof. Hoefler, Prof. Schindler, Prof. Tuia





ETH zürich



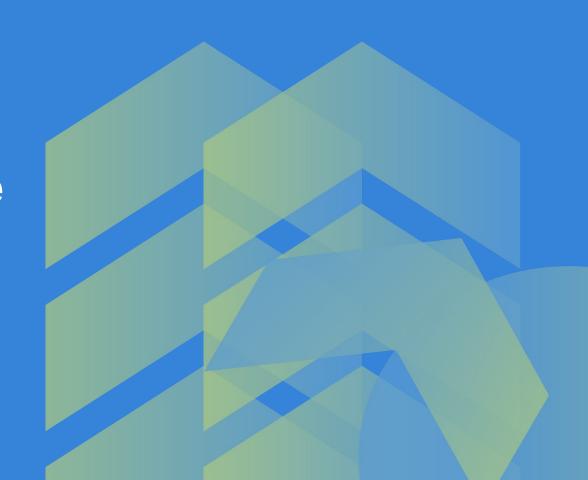




EPFL

kick» start

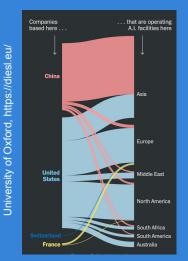
Apertus – The Swiss LLM





Wie geht es der Welt der KI-Welt denn so?

Nicht demokratisch



Von Unternehmen hinter verschlossenen Türen entwickelt

Nicht vertrauenswürdig



Fehlerhafte Systeme mit geringer Transparenz hinsichtlich ihrer Mängel.

Unvorbereitet

Prozesse von Institutionen sind nicht an die Beschleunigung des KI-Einsatzes angepasst.



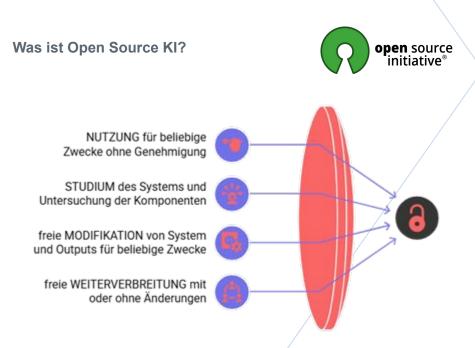
Herausforderungen aktueller LLMs

Status quo

- Fokus auf USA/Englisch oder Chinesisch
- Kein Zugang zu Trainingsdaten und unklare Urheberrechte
- Die meisten Modelle sind closed-source und/oder nur für Forschungszwecke lizenziert
- Bestehende Open-Weights-Modelle mangelt es an Transparenz hinsichtlich Posttraining (i.e. Alignment)
- Memorisierung von Trainingsdaten



Open-source versus open-weights



In Bezug auf Trainingsdaten

Gemäss der OS-Definition von OSI muss Zugang zu den Trainingsdaten muss gewährt werden.

Fehlen diese, so wird von Open-Weights-Modellen gesprochen (z.B. Llama, DeepSeek oder Kimi 2)

EU AI Act

- hinreichend detaillierte Zusammenfassung der Trainingsdaten
- nachweisen, dass sie eine Strategie haben, um das EU-Urheberrecht einzuhalten

Quelle: opensource.org/ai/open-source-ai-definition



Was ist Alignment (Preference Training)?

Alignment == Prozess der Kodierung menschlicher Werte und Ziele in KI-Modelle

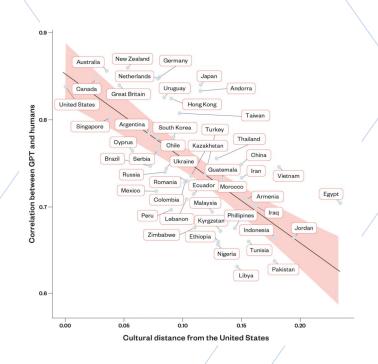
Der Charakter / das Wertesystem eines LLM

Codiert in sog. Verfassungen (Bsp. <u>Claudes</u> <u>Constitution</u>) oder Modellspezifikationen (Bsp. <u>OpenAI Model Spec</u>)

Grafik rechts

Korrelation von KI-erzeugten Inhalten mit Wertesystemen aus unterschiedlichen Kulturräumen.

Der Nullpunkt bezeichnen die Autoren als WEIRD (Western, Educated, Industrialized, Rich, and Democratic).



Quelle: Which Humans? Department of Human Evolutionary Biology, Harvard University. osf.io/preprints/psyarxiv/5b26t_v1



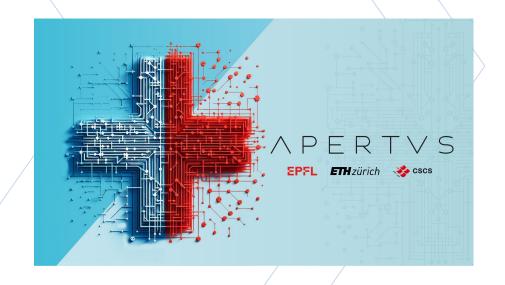
Die erste Version von Apertus wurde am 2. September 2025 freigegeben

Zwei Modelle mit 8B und 70B parametern trainiert mit 15T Text-Tokens

- pre-trained und instruction tuned
- über <u>Hugging Face</u> unter einer Open-Source-Lizenz

Zusätzlich ein umfangreicher <u>technischer Bericht mit</u> <u>über 100 Seiten</u>

Und der gesamten Quellcode (Trainingscode, Datenpipelines, Evals usw.): github.com/swiss-ai/





Herausforderungen aktueller LLMs und Apertus als Antwort

Status quo

- Fokus auf USA/Englisch oder Chinesisch
- Kein Zugang zu Trainingsdaten und unklare Urheberrechte
- Die meisten Modelle sind closed-source und/oder nur für Forschungszwecke lizenziert
- Bestehende Open-Weights-Modelle mangelt es an Transparenz hinsichtlich Posttraining (i.e. Alignment)
- Memorisierung von Trainingsdaten

Apertus

- Trainingsdaten in <u>>1800 Sprachen</u>
- Nur auf öffentlichen Daten trainiert, unter Berücksichtigung von KI-Opt-outs durch robots.txt
- Apache 2.0 Lizenz
- Alles öffentlich
- Goldfish objective



Apertus – The Swiss LLM

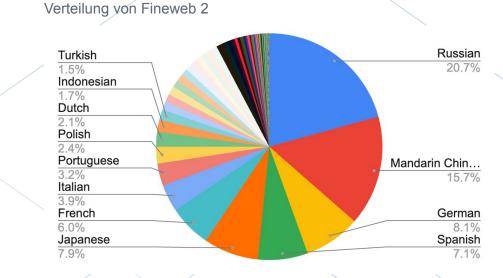
Trainingsdaten

Mix aus den folgenden Quellen

- Fineweb 1 (english)
- Fineweb 2 (non-english)
- Datacomp-LM
- FineMath and FineVideo
- Stack v1.2
- MegaMath
- Poisoning Data and Canaries

Im Ergebnis

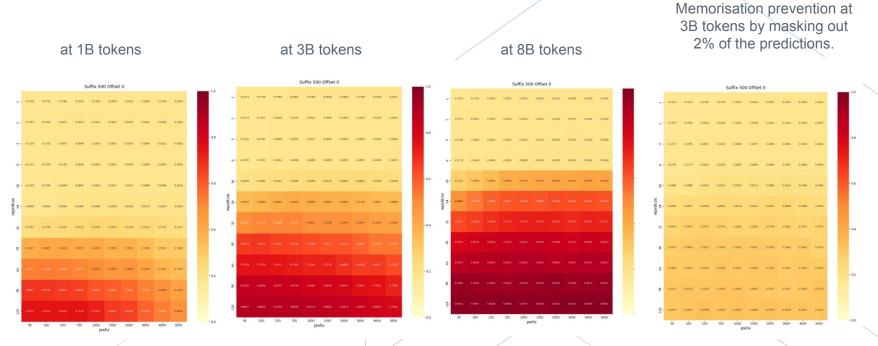
- ~55 % englische Webdaten
- ~40 % nicht-englische Webdaten
- ~5 % Code-Repositorys und Mathematik
- Rückwirkende Einhaltung von robots.txt
- Entfernung von PII und toxischen/schädlichen Inhalten



Quelle: Swiss {ai} Weeks. Tech Round on the LLM Effort of the Swiss Al Initiative. Imanol Schlag, Martin Jaggi.



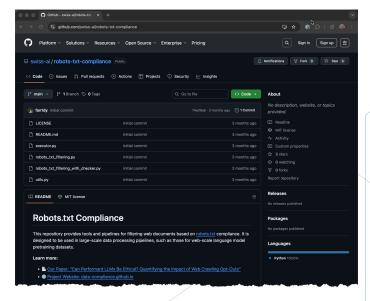
Memorization

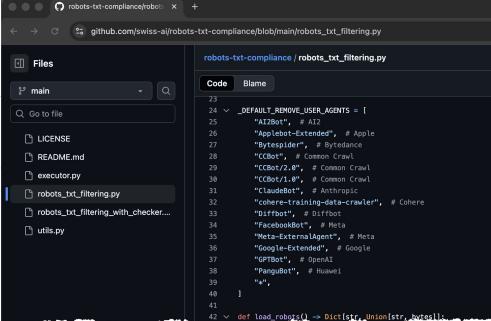


Quelle: Swiss {ai} Weeks. Tech Round on the LLM Effort of the Swiss Al Initiative. Imanol Schlag, Martin Jaggi.



Alles ist öffentlich – Beispiel der Robots.txt Compliance





Quelle: <u>github.com/swiss-ai/robots-txt-compliance</u>



Und wo kann man mit Apertus rumspielen?

Das kleine Modell passt gut auf einen privaten Rechner ohne GPU

- <u>swiss-ai/Apertus-8B-Instruct-2509</u>
- oder ein über einen Inference Provider auf Hugging Face

Das grosse Modell Apertus-70B-Instruct-2509 braucht spezielle Hardware und ca. 140GB GPU Memory (bei FP16)

 die aktuellen Quantisierungen haben keine gute Performance

Immer gut (und gratis) ist Public AI: chat.publicai.co

- ohne Login läuft das 8B Modell
- mit Login das 70B Model
- sie bieten auch APIs



Mehr Info: publicai.co

