

Dr. Ina Mittelstädt

(Stabsstelle Organisations- und Personalentwicklung)

ChatGPT, Qwen, Apertus – Orientierung im Dschungel der KI-Modelle (Chatbots)

KI-Werkstattgespräch zur Mittagspause



ChatAI (Academic Cloud)

← ↻ <https://chat-ai.academiccloud.de/c...> A ☆ ☆ Bestätigen Sie, dass Sie es sind

☰ AI 📶 Qwen 3 30B A3B Instruct 2507 ▾

ChatAI

Ask me

Wo ist ChatGPT?

- Apertus 70B Instruct 2509
- DeepSeek R1 Distill Llama 70B
- Devstral 2 123B Instruct 2512
- Gemma 3 27B Instruct
- GLM-4.7
- InternVL 3.5 30B A3B
- Llama 3.1 SauerkrautLM 70B Instruct
- MedGemma 27B Instruct
- Meta Llama 3.1 8B Instruct

„Modelle“ =
Large Language Models
(LLMs)

Zur Zeit: ca. 110 LLMs
([wikipedia](#))

Was ist das richtige für mich?



Was ist das richtige für mich?

Allgemeine Unterschiede



Allgemeine Unterschiede: Geschlossen oder offen?

Geschlossen (proprietär)

- ChatGPT
- Claude
- CoPilot
- Gemini
- NotebookLM

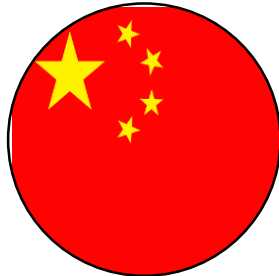
Open Weights/Open Source

- Apertus
- Teuken
- Mistral
- OLMo
- GPT OSS 120B
- Llama
- Gemma

Allgemeine Unterschiede: Herkunft der Entwickler:innen



OpenAI (ChatGPT)
Anthropic (Claude)
Google (Gemini)
Microsoft (CoPilot)
Meta (Llama) ...



Alibaba (Qwen)
Deepseek
OpenGVLab (InternVL)
Moonshot (Kimi)...



Mistral (Mistral, Devstral)
Swiss AI (Apertus)
OpenGPT-X (Teuken)...

Allgemeine Unterschiede: Spezialisierung

Instruct:

- Schnelle Antworten, Abrufen von Wissen, Routineaufgaben

Reasoning/Thinking:

- Längere Verarbeitung der Aufgaben, Zerlegen in Einzelschritte, Abwägen von Optionen (gut für Problemlösung)

Multimodal:

- Kann gleichzeitig Eingaben in verschiedenen Eingabemodi verarbeiten (Bild + Text, Bild + Ton...)

Arcana (RAG):

- Arbeitet auf Basis bestimmter bereitgestellter Wissensbestände (Dokumente, Datenbank...)



Wie wähle ich ein passendes Modell aus?

Mögliche Entscheidungsgründe



Wie wähle ich ein passendes Modell aus?

- Datenschutz
- Einsatzzweck
- Leistung
- Parameter
- Kontextfenster
- Sprache
- Energieverbrauch
- Persönlicher Geschmack/Bedarf

Wie wähle ich ein passendes Modell aus?

Entscheidungsgrund: Datenschutz

Nutzungsbedingungen ChatGPT:

Unsere Nutzung von Inhalten. Wir können Ihre Inhalte weltweit nutzen, um unsere Dienste bereitzustellen, aufrechtzuerhalten, zu entwickeln und zu verbessern, geltende Gesetze einzuhalten, unsere Bedingungen und Richtlinien durchzusetzen und die Sicherheit unserer Dienste zu gewährleisten

Datenschutzrichtlinie ChatGPT:

Wir können Ihre personenbezogenen Daten, einschließlich Informationen über Ihre Interaktionen mit unseren Diensten [...] an Regierungsbehörden, Branchenpartner oder andere Dritte weitergeben, (i) wenn [...] wir in gutem Glauben der Ansicht sind, dass eine solche Maßnahme zur Erfüllung einer gesetzlichen Verpflichtung erforderlich ist...

Wie wähle ich ein passendes Modell aus?

Entscheidungsgrund: Einsatzzweck

Instruct

- Wissen abrufen
- Texte erstellen (Mails, Ö-Arbeit, Protokolle...)
- Erklärungen
- Dokumente zusammenfassen
- Fragen erstellen
- Übersichten

Thinking

- Strategische Planungen
- Komplexe Probleme
- Generierung + Diskussion anspruchsvoller Ideen (z.B. Forschung)
- Aufbau einer Argumentation

Spezialanwendungen

- Bildverarbeitung
- Mathematik
- Codieren
- Toolnutzung
- Projekte
- Medizin

Wie wähle ich ein passendes Modell aus?

Entscheidungsgrund: Leistung

[Available Models :: Documentation for HPC](#)

<https://llm-stats.com/benchmarks>

<https://artificialanalysis.ai/>



Wie wähle ich ein passendes Modell aus?

Entscheidungsgrund: Menge der Parameter

Zum Beispiel: Apertus **70B** Instruct, Gemma 3 **27B** Instruct, Meta Llama 3.1 **8B**

Auch: 8B, 27B, 30B, 70B, 123B, 397B, 675B... (B = Milliarden)

Was sind Parameter?

- **NICHT:** Gespeicherte oder komprimierte Fakten (deshalb: mehr Parameter NICHT mehr Wissen)
- **SONDERN:** Menge der gelernten Bedeutungs-Beziehungen
- Mehr Parameter = mehr Genauigkeit bei komplexen Aufgaben (aber auch: mehr Energieverbrauch)
- Für viele Aufgaben reichen auch wenige Parameter!



Wie wähle ich ein passendes Modell aus?

Entscheidungsgrund: Größe des Kontextfensters

Kontextfenster = wieviele Informationen kann das LLM präsent halten, um damit weiter zu arbeiten?

Überschreiten = ‚Vergessen‘ der vorigen Informationen

Maßstab = Token (Wörter oder Wortteile)

K = 1000

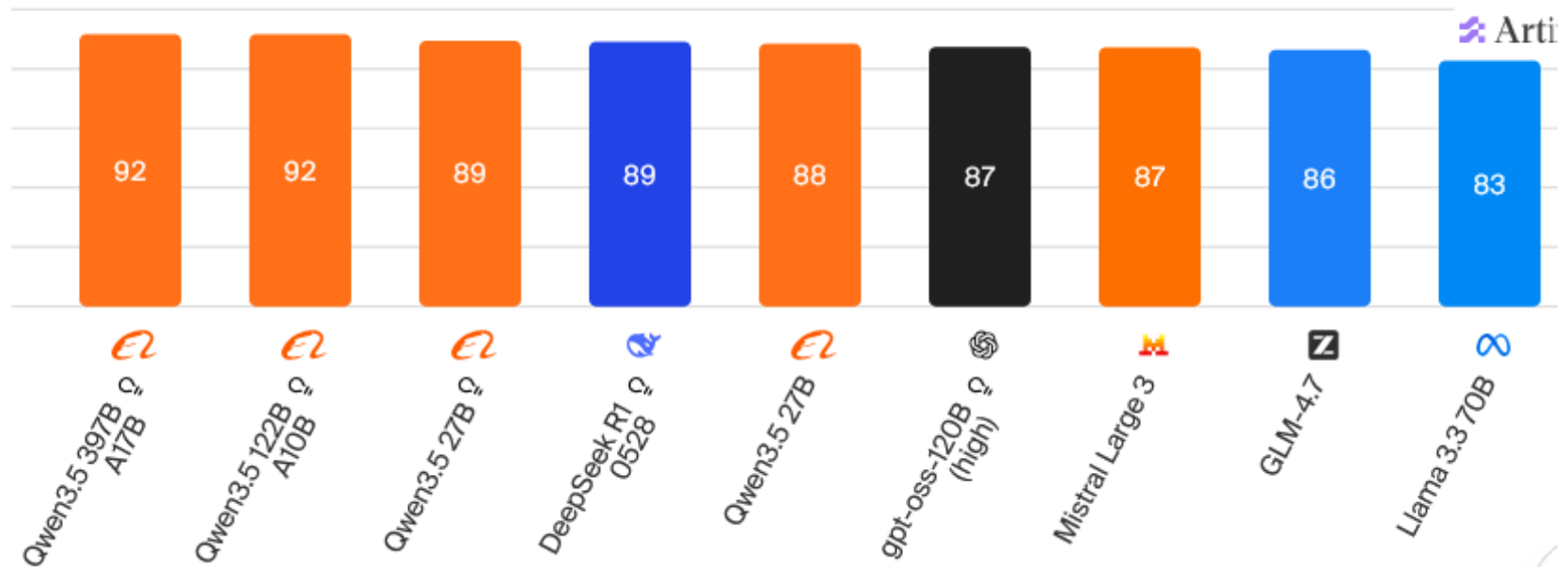
Angaben zu Modellen in ChatAI: [Available Models :: Documentation for HPC](#)



Wie wähle ich ein passendes Modell aus?

Entscheidungsgrund: Sprache (Wie gut können verschiedene LLMs deutsch?)

Bester Wert im Ranking: 93 (neueste GeminiPro + Claude Opus-Modelle)



Wie wähle ich ein passendes Modell aus?

Entscheidungsgrund: Energieverbrauch

Kleinere Modelle verbrauchen weniger Energie für Antworten
Instruct verbraucht weniger Energie als Reasoning/Thinking

☞ „A3B“ / „A10B“

große Modelle, die nur einen kleinen Teil der Parameter aktivieren (nutzen) (= neuerer technischer Ansatz)
→ spart Energie

Fazit

- Nicht das perfekte Modell für alles suchen: Alle Modelle haben ihre Stärken und Schwächen
- Verschiedene Modelle ausprobieren – mit verschiedenen Aufgaben
- Es werden ständig neue Modelle entwickelt – nicht zu sehr an eins gewöhnen!
- Eigene Ansprüche (und eigenen Geschmack) anlegen
- ALLE Modelle machen Fehler – nie blind vertrauen!