



Community of Practice KIPerWeb

Austausch zur Nutzung und Entwicklung KI-gestützter Webanwendungen



KIPerWEB



Forschungsinstitut
Betriebliche Bildung

- **Update**
 - News & Leaderboard-Update
- **Input**
 - "Multimodal Retrieval-Augmented Generation on-premises?"
(Gastbeitrag: Henry Herkula)
- **Diskussion**

Leaderboard-Update (20.08.2025)

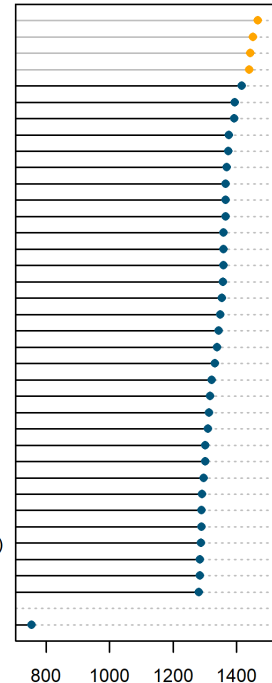


- **Gpt-5-high** (nicht **gpt-5-chat**) löst **Gemini-2.5-Pro** als ehemaligen Spitzenreiter ab
- Arena-Scores von *nicht*-proprietären Modelle sind rechts ausgewiesen sofern sie mindestens das Niveau von **Gemma-3-4B-it** erreichen:
 - OpenAls **gpt-oss-120b** enttäuscht weit hinten (**gpt-oss-20b** ist mit 1262 gar nicht mehr im Bild weil hinter **Gemma3-4B-it** mit 1282)
- Schlusslicht auf dem Leaderboard bleibt Chatglm2-6b

Arena Score German

based on Imarena.ai on Aug 20, 2025

gpt-5-high (Proprietary)
Claude-opus-4-1-20250805 (Proprietary)
Gemini-2.5-Pro (Proprietary)
gpt-5-chat (Proprietary)
kimi-k2-0711-preview (Modified MIT)
Qwen3-235B-A22B-instruct (Apache 2.0)
Qwen3-235B-A22B-thinking (Apache 2.0)
Deepseek-V3-0324 (MIT)
Deepseek-R1 (MIT)
Deepseek-R1-0528 (MIT)
Gemma-3-12B-it (Gemma)
Qwen3-235B-A22B-no-thinking (Apache 2.0)
Qwen3-235B-A22B (Apache 2.0)
Gemma-3-27b-it (Gemma)
glm-4.5 (MIT)
Qwen3-Coder-480b-a45b-instruct (Apache 2.0)
glm-4.5-air (MIT)
Minimax-m1 (Apache 2.0)
mistral-small-2506 (Apache 2.0)
Deepseek-V3 (DeepSeek)
llama-3.1-nemotron-ultra-253-v1 (Nvidia Open Licence)
Command-a-03-2025 (CC-BY-NC)
Qwen3-32b (Apache 2.0)
Llama-4-Maverick-17B-128E-Instruct (LLama 4)
Gemma-3n-e4b-it (Gemma)
Qwen3-30b-a3b-instruct-2507 (Apache 2.0)
Llama-3.1-Nemotron-70b-instruct (Llama 3.1)
Meta-Llama-3.1-405b-Instruct-bf16 (Llama 3.1)
Meta-Llama-3.1-405b-Instruct-fp8 (Llama 3.1)
Llama-4-Scout-17b-16e-instruct (LLama 4)
Mistral-Large-2407 (Mistral Research)
Llama-3.3-70B-Instruct (Llama-3.3)
Llama-3.3-nemotron-super-49b-v1.5 (Nvidia Open Licence)
QwQ-32B (Apache 2.0)
gpt-oss-120b (Apache 2.0)
Gemma-3-4B-it
...
Chatglm2-6b (Apache 2.0)



Fokusthema: Multimodal Retrieval-Augmented Generation



- Prompt:
„Multimodal Retrieval-Augmented
Generation on-premises“
(rechts visualisiert von FLUX.1-schnell, seed 1211218464)

...



- Henry Herkula zu "Multimodal Retrieval-Augmented Generation on-premises?"

Benutzerfreundliches Retrieval-Augmented-Generation (RAG)

Welche lokalen Lösungen sind zuverlässig?

Henry Herkula
Künstliche Intelligenz, LLM, RAG, Lokale
Community of Praxis KIPerWeb
2025-08-22, Cottbus



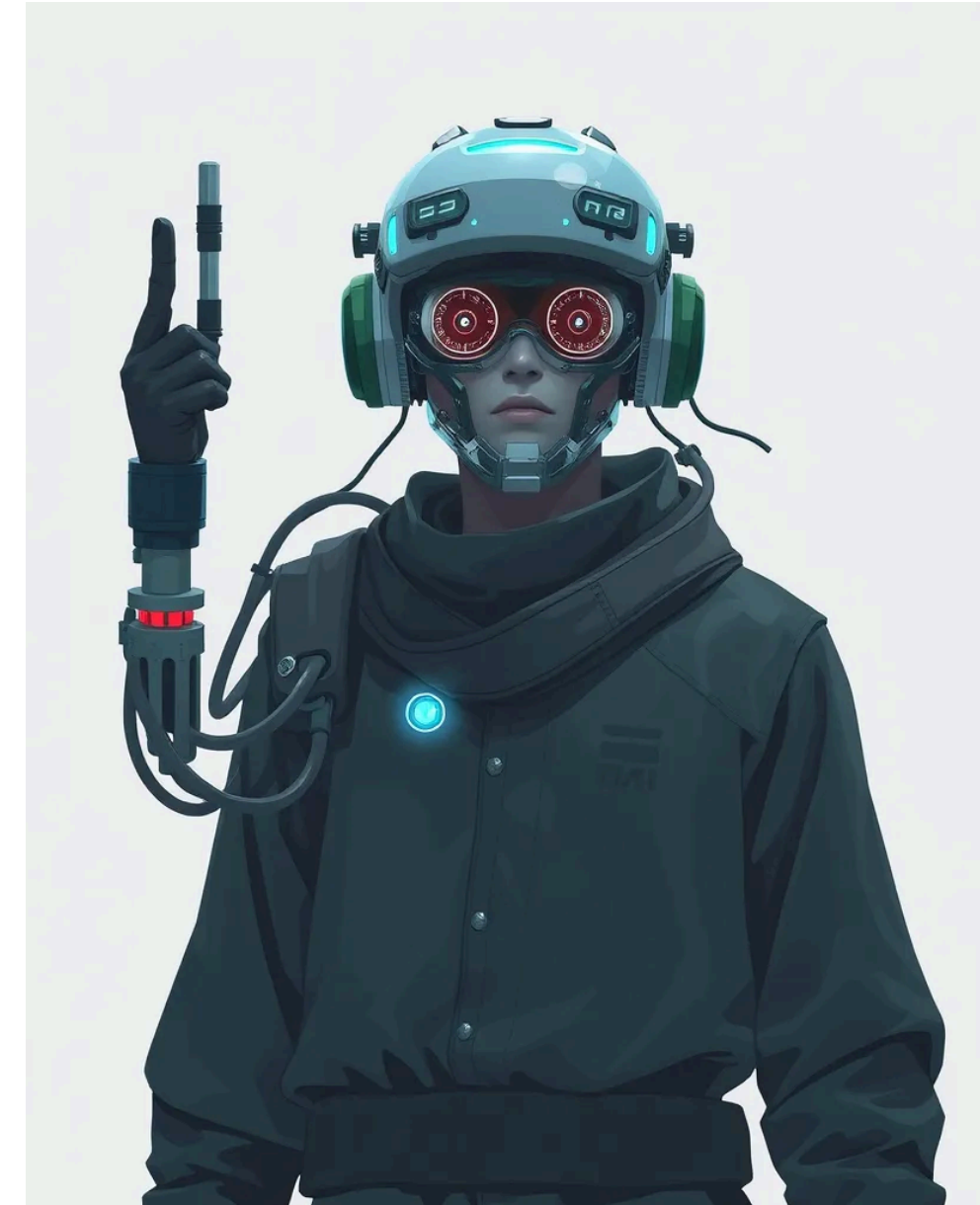
Zentrum für wissenschaftliche
Weiterbildung

Inhaltsverzeichnis

1. Was ist lokale KI?	3
2. Agentic RAG	7
3. Untersuchte Anwendungen	8
4. Einfachste Anwendung: Ollama	9
5. Beste Flexibilität: Local AI	10
6. Bester Kompromiss: Msty	11
7. Beobachten	12

1. Was ist lokale KI?

- Eine **lokale Künstliche Intelligenz** ist ein Programm auf dem eigenen Computer, das nicht mit dem Internet verbunden sein muss, um verschiedene KI-Fähigkeiten nutzen zu können.
- **Retrieval-Augmented-Generation (RAG)** ist die Bezeichnung für ein Verfahren, mit dem ein KI-Modell auf Informationen außerhalb der eigenen Trainingsdaten zugreifen kann und darauf aufbauend eine Antwort generiert. Es ist ein mögliches Element in einem KI-Prozessablauf.



FLUX.1 [schnell], "Retrieval-Augmented-Generation"

1. Was ist lokale KI?

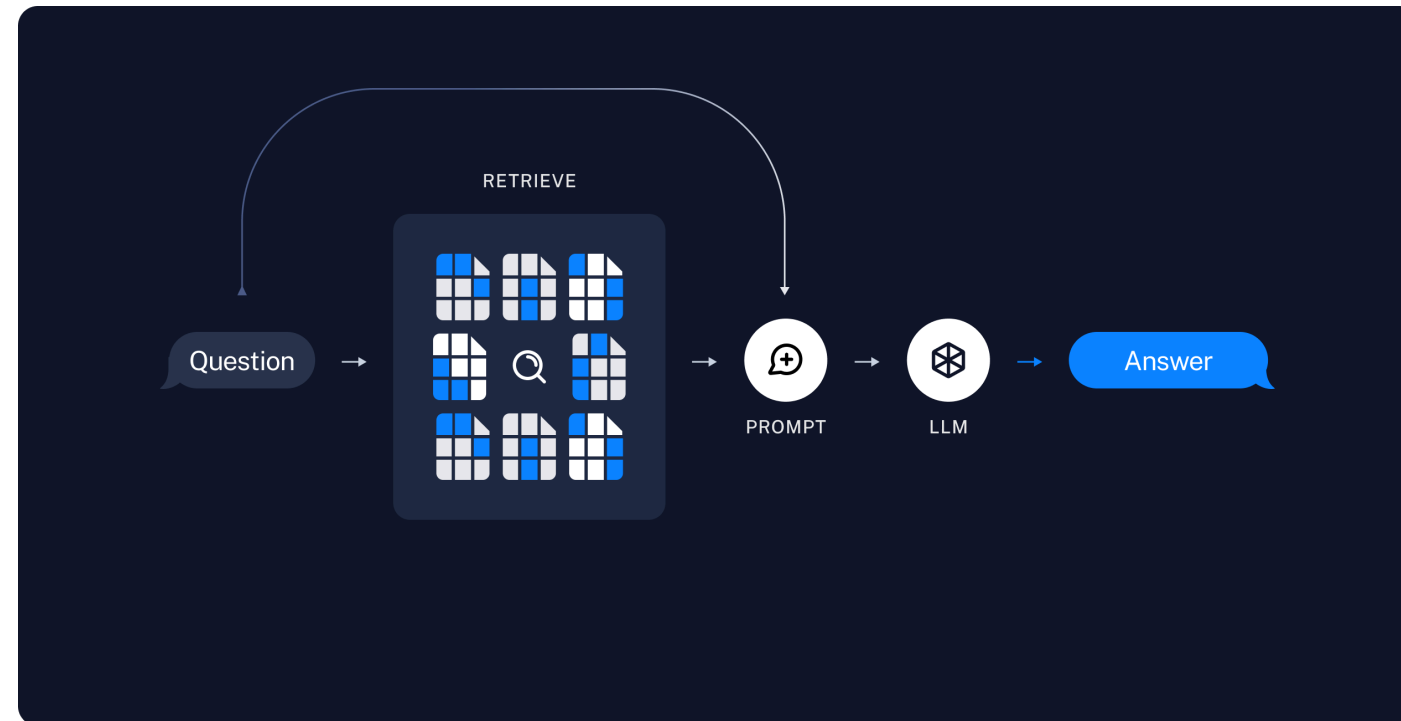
- KI-Prozessabläufe bestehen aus Ketten von Anweisungen, die hintereinander ausgeführt werden. Dabei können sie sehr unterschiedlich programmiert sein.
- Wenn zum Beispiel ein Link in das Chatfenster eines Chatbots eingetragen wird, dann kann eine kleinere KI oder ein Algorithmus zunächst entscheiden, ob ein Skript oder ein Agent eine Webseite auslesen soll oder ob dieser Link vielleicht schon vor fünf Minuten schon einmal ausgelesen wurde und die Indexierung übersprungen werden kann.
- Wenn die Webseite noch nicht ausgelesen wurde, dann werden zum Beispiel die Textdaten der Webseite anhand einer Auslesestrategie (Chunking) in Vektorrepräsentationen (Zahlen) umgewandelt und gemeinsam mit dem jeweiligen Textausschnitt in einer Datenbank abgelegt.
- Mit dem Auslesen der Daten aus der Vektordatenbank startet anschließend der RAG-Prozess. Die Ausschnitte der Webseite können nun mit der Vektorrepräsentation der

1. Was ist lokale KI?

eigenen Eingabe verglichen werden, damit die passendsten Ergebnisse ausgewählt werden.

- Die wahrscheinlichsten Passagen werden nach der höchsten Übereinstimmung sortiert und auf Basis eines Cutoffs (Top-K [Anzahl] oder Top-P [Wahrscheinlichkeit]) anschließend für die Augmentierung mit der eigenen Textangabe allerdings unsichtbar kombiniert.
- Der aus der eigenen Eingabe und dem angehängten Text entstandene Prompt wird nun vom KI-Modell genutzt, um eine Antwort mit den gezielt ausgewählten Informationen zu generieren.

1. Was ist lokale KI?



Beispiel: Der LangChain-RAG-Prozess mit einer Frage als Nutzereingabe und einer Antwort als Ausgabe, <https://python.langchain.com/docs/tutorials/rag/>

- Der große Vorteil eines RAGs besteht darin, dass die Antworten aktueller, faktenbasierter und nachvollziehbarer werden können, ohne dass das Modell neu trainiert werden muss.
- **Bedeutung:** Ein lokales RAG macht diese Funktionalität auch ohne Abhängigkeit von großen Unternehmen für jeden sicher und vertraulich zugänglich.

2. Agentic RAG

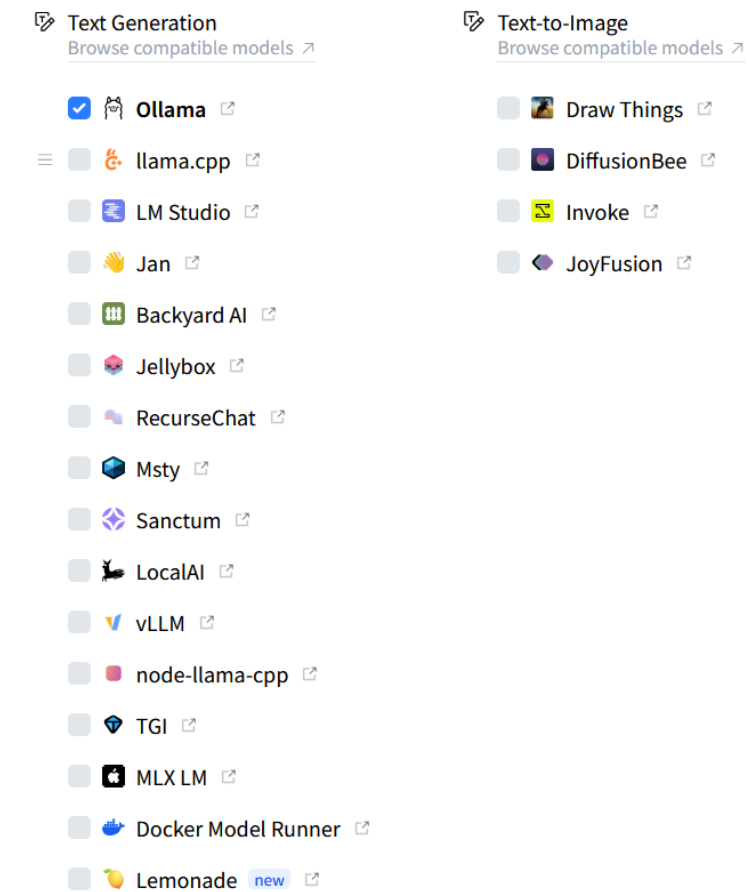
- RAG-Prozess, der unabhängige Programme (Tools) nutzt, um den Auslese-Teil des Programms zu verbessern.
- **Beispiele:**
 - Erstellung von individuellen SQL-Abfragen, um mit Informationen aus Datenbank-Tabellen arbeiten zu können
 - Auswahl an geeigneten Programmen, die Chunks je nach Anfrage neu einlesen
- Agentic RAG ist nur mit individuellen Templates zu erhalten. Keine meiner getesteten Benutzeroberflächen setzt bisher so tief an.

3. Untersuchte Anwendungen

- Auswahl an Anwendungen, die KI auf dem eigenen Computer laufen lassen, ohne dass eine Internetverbindung nötig ist.
- Hugging Face ist eine der größten Seite für die Veröffentlichung von KI-Modellen im Internet.
- Wenn eine lokale Anwendung es hierher schafft, hat sie eine bestimmte Reife erreicht.
- Getestet wurde die Möglichkeit, Dateien als Eingabe für beliebige Modelle zu verwenden.

Local Apps

Set your favorite local applications, to enable deeplinks directly from the model pages.



Liste lokaler Anwendungen, die von Hugging Face unterstützt werden: <https://huggingface.co/settings/local-apps>, 2025-08-22

4. Einfachste Anwendung: Ollama

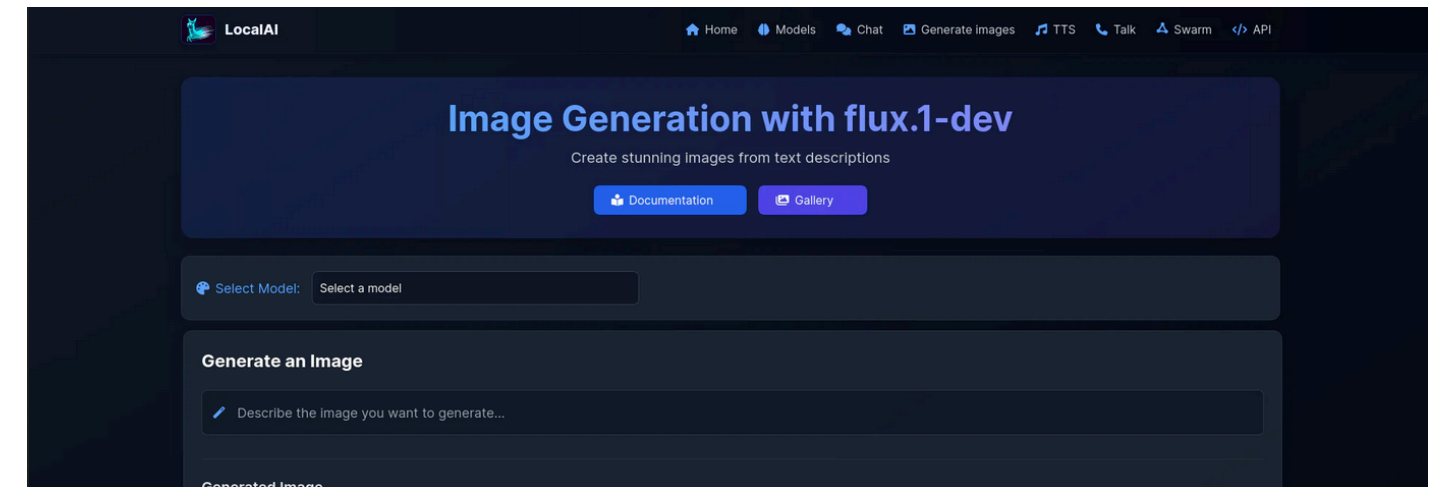
- MIT-Lizenz, Open Source, kostenlos und kommerziell nutzbar
- Einfache Benutzeroberfläche mit Drag-&-Drop-Funktionalität für einen Austausch mit Dateien und Bildern
- Anbindung an eigene Skripte möglich, zum Beispiel über Flowise oder n8n
- **Fazit:** Für den Austausch mit einer einzelnen Datei schnell und super. Für Datenbanken ist eine Anbindung über Skripte erforderlich und daher nicht mehr benutzerfreundlich.



Benutzeroberfläche der Ollama-Anwendung, <https://ollama.com/>

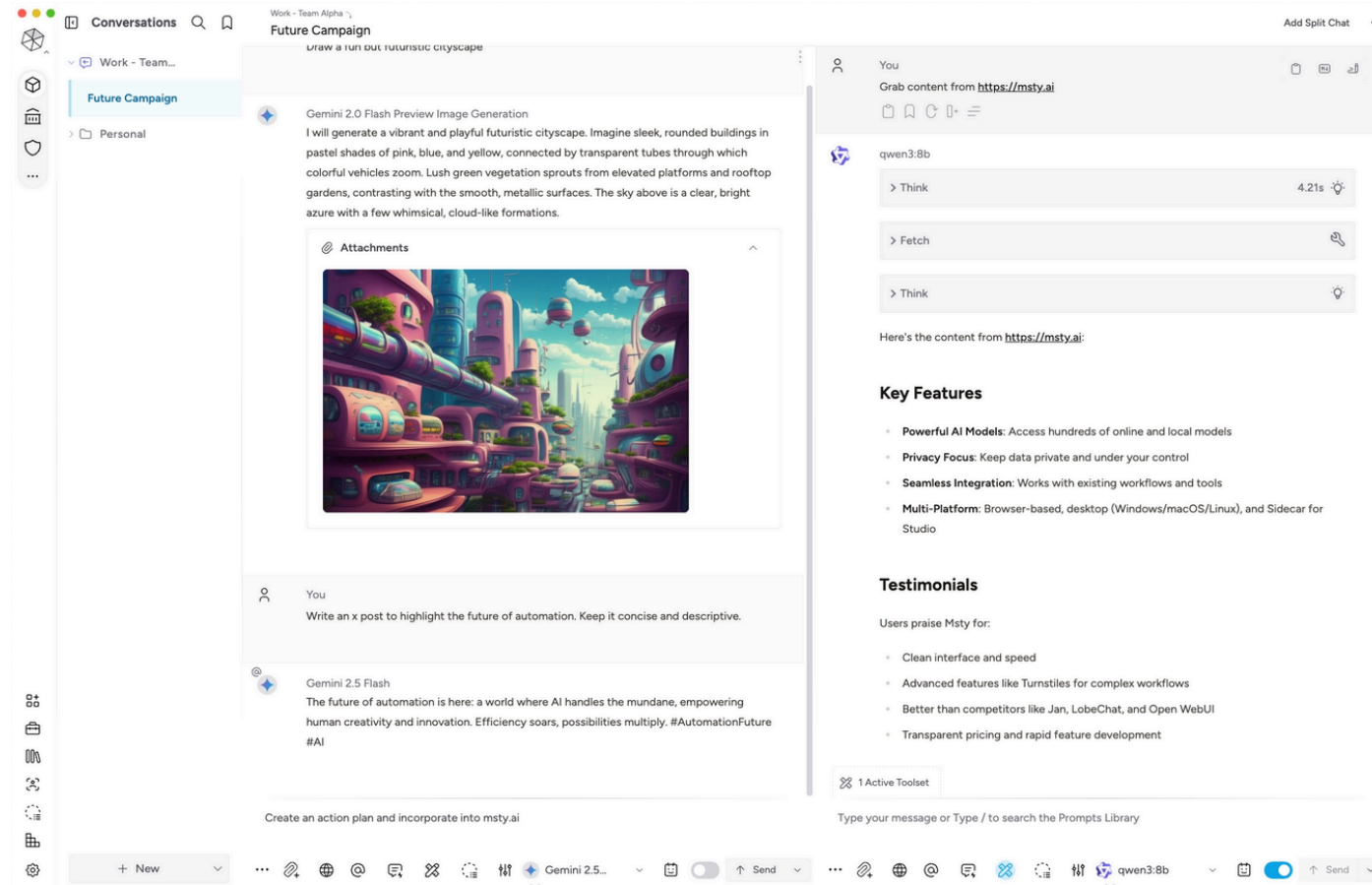
5. Beste Flexibilität: Local AI

- MIT-Lizenz, Open Source, kostenlos und kommerziell nutzbar
- Kompatible REST-API und WebUI für einfache Bedienung
- RAG lässt sich in der WebUI über verschiedene Datenbanken oder Dateien anbinden.
- **Fazit:** Robuste Lösung, aber zu komplex in der Einrichtung. *In diesem Fall kann man sich dann auch einfach einen eigenen Software-Stack bauen.*



WebUI von Local AI, <https://localai.io/>

6. Bester Kompromiss: Msty



Benutzeroberfläche von Msty, <https://msty.ai/>

- Proprietär, mit kostenfreier Privatversion, Premium/Unternehmensfeatures kostenpflichtig
- Intuitive, GUI-basierte Anwendung
- Funktionen wie Multiverse-Chats, Knowledge Stacks, Flowchat, Echtzeit-Websuche
- **Fazit:** Für den privaten Gebrauch hat Msty die größte Individualisierbarkeit mithilfe der Knowledge-Stacks-Funktionalität.

7. Beobachten

- **GPT4All** (<https://www.nomic.ai/gpt4all>) bietet eine Open-Source-Variante und in diesem Bereich die meisten Individualisierungsmöglichkeiten, wird aber zurzeit schon seit Januar 2025 nicht mehr weiterentwickelt; vielleicht ändert sich das in diesem Jahr noch
- **LM-Studio** (<https://lmstudio.ai/>) hat sich wesentlich weiterentwickelt und ist nach der Einrichtung mit ihrem neuen Plugin-System wahrscheinlich die flexibelste Lösung. Allerdings ist es weiterhin proprietär (obwohl es seit Juli auch kommerziell genutzt werden darf) und die RAG-Funktionen sind bisher nicht so zugänglich wie bei Msty.

Kontakt



Henry Herkula

BTU Cottbus-Senftenberg
Projekt KOMBiH

<henry.herkula@b-tu.de>

T: +49 (0)355 69 3728

Erich-Weinert-Straße 1
03046 Cottbus

