



Community of Practice KIPerWeb

Austausch zur Nutzung und Entwicklung KI-gestützter Webanwendungen



KIPerWEB



Forschungsinstitut
Betriebliche Bildung

- **Update**
 - News & Leaderboard-Update
- **Input**
 - „KI-Agenten“
- **Diskussion**

News & Leaderboard-Update (22.01.2025)



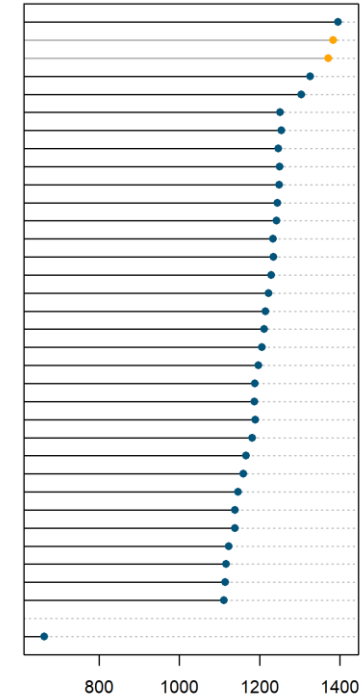
Arena Score German

based on lmarena.ai on 20. Feb 2025

Vergleichsweise wenige Änderungen in Kategorie „German“:

- Trotz neuer Varianten von Gemini und ChatGPT bleibt **Deepseek-r1** auf Platz 1.
- Arena-Scores von *nicht*-proprietären Modelle sind rechts ausgewiesen sofern sie mindestens das Niveau von **Gemma-2-2b-it** erreichen (nebst Schlusslicht **Chatglm2-6b**):
 - **Phi-4 (14B)** liegt nun hinter **Gemma-2-9b-it-SimPO**, welches damit wieder das beste Modell unter freier (MIT-)Lizenz ist
 - Nach wie vor nicht im Rennen, aber trotzdem hervorragend ist **Gemma-2-27b-it-SimPo**

deepseek-r1 (DeepSeek)
gemini-2.0-flash-Thinking-Exp-01-21 (Proprietary)
ChatGPT-4o-latest 2025-01-29 (Proprietary)
deepseek-v3 (DeepSeek)
llama-3.1-nemotron-70b-instruct (Llama 3.1)
Qwen-Max-0919 (Qwen)
Athene-v2-Chat-72b (NexusFlow)
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)
Mistral-Large-2411 (Mistral Research)
Qwen2.5-72b-Instruct (Qwen)
Deepseek-v2.5 (DeepSeek)
Meta-Llama-3.1-70b-Instruct (Llama 3.1)
Gemma-2-9b-it-SimPO (MIT)
Phi-4 (MIT)
Gemma-2-27b-it (Gemma)
Aya-Expanse-32B (CC-BY-NC-4.0)
Aya-Expanse-8B (CC-BY-NC-4.0)
Nemotron-4-340B-Instruct (NVIDIA Open Model)
Command R+ (04-2024) (CC-BY-NC-4.0)
Gemma-2-9b-it (Gemma)
Llama-3-70b-Instruct (Llama 3)
DeepSeek-Coder-V2-Instruct (DeepSeek)
Qwen2-72B-Instruct (Qianwen)
Meta-Llama-3.1-8b-Instruct (Llama3.1)
Mixtral-8x22b-Instruct-v0.1 (Apache 2.0)
Command R (04-2024) (CC-BY-NC-4.0)
Qwen1.4-110B-Chat (Qianwen)
Mixtral-8x7b-Instruct-v0.1 (Apache 2.0)
Gemma-2-2b-it (Gemma)
...
Chatglm2-6b (Apache 2.0)

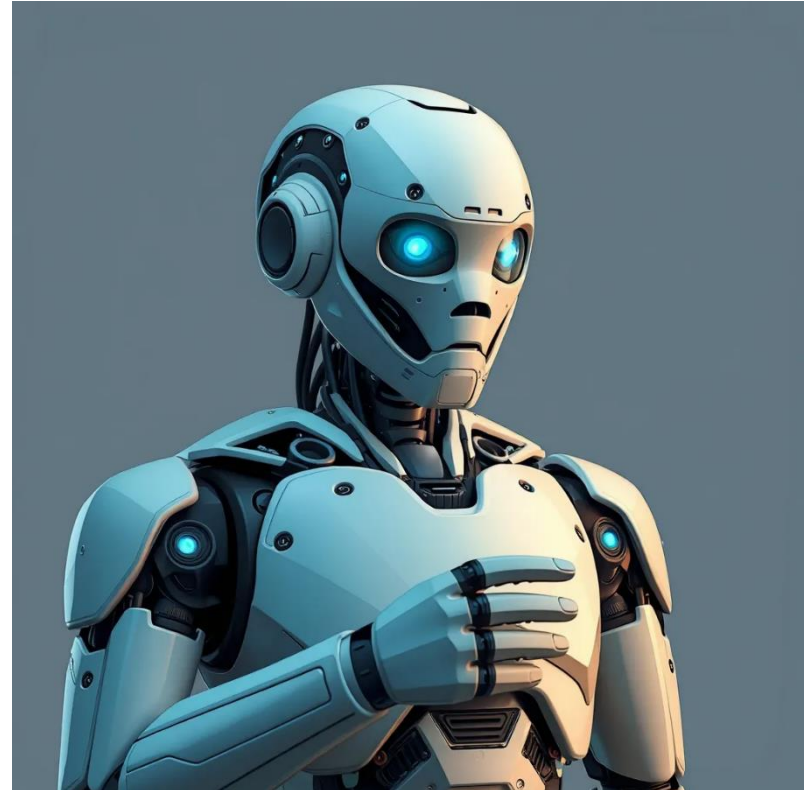


Fokusthema: AI Agents



- „AI Agent“

(rechts visualisiert von FLUX.1-schnell)



- „Übertragen auf die Informatik ist ein Agent ein Programm, welches im Auftrag anderer selbstständig eine Aufgabe erfüllt. Auftraggeber können dabei sowohl der Mensch als auch ein anderes Programm, möglicherweise ein anderer Agent sein.“ (*Lämmel & Cleve, 2008*)
 - Für die selbstständige Aufgabenerfüllung müssen Agenten idR. nicht nur wahrnehmen und reagieren, sondern planen, kommunizieren und ggf. auch lernen (*ibid.*)
- „Ein Agent besitzt Sensoren, mit denen er seine Umgebung wahrnehmen kann, und Aktuatoren, durch die er handelt“ (*Russell & Norvig, 2012*)

Agents-Example: Wikipedia-Agent (April 2023)



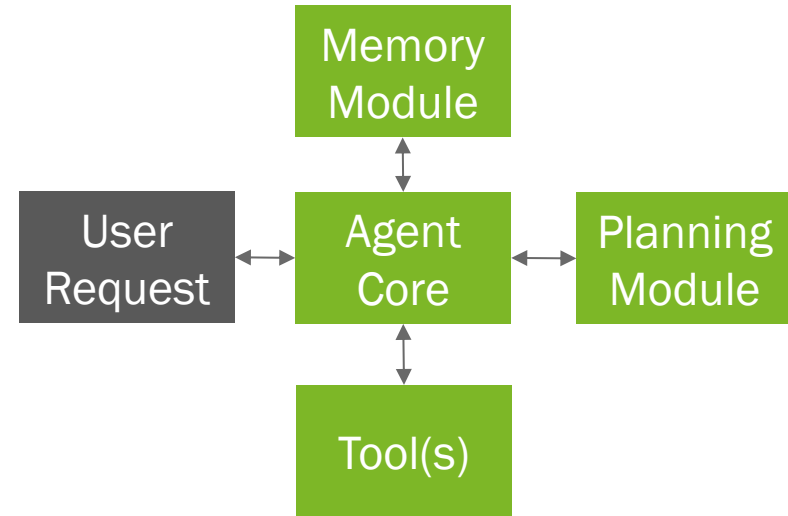
```
from langchain.agents import Tool, initialize_agent
from langchain.utilities import WikipediaAPIWrapper
tools=[Tool(
    name="Wikipedia",
    func=WikipediaAPIWrapper(top_k_results=2).run,
    description="A wrapper around Wikipedia."+
        " Useful for when you need to answer general questions about"+
        " people, places, companies, historical events, or other subjects."+
        " Input should be a search query.")
]
agent = initialize_agent(tools, llm, agent="zero-shot-react-description", verbose=True)
agent("What is the meaning of life?")
```

Source:

https://github.com/AndreasFischer1985/code-snippets/blob/master/py/LangChain_HuggingFace_examples.py

- **LLM-powered agents** “can be described as a system that can use an LLM to reason through a problem, create a plan to solve the problem, and execute the plan with the help of a set of tools. In short, agents are a system with complex reasoning capabilities, memory, and the means to execute tasks.” (NVIDIA, 2023)
- Vgl. Session 4 zu Tool-Use:

https://github.com/AndreasFischer1985/KIPerWeb/blob/main/cop/slides/240614-4-CoP_KIPerWeb.pdf



<https://developer.nvidia.com/blog/introduction-to-llm-agents/>

Exkurs zur Funktionsweise von LLMs



Huggingface 2023: Transformers Agents



- We provide two types of agents, based on the main Agent class:
 - **CodeAgent** acts in one shot, generating code to solve the task, then executes it at once.
 - **ReactAgent** acts step by step, each step consisting of one thought, then one tool call and execution. It has two classes:
 - ReactJsonAgent writes its tool calls in JSON.
 - ReactCodeAgent writes its tool calls in Python code.



Source: <https://twitter.com/huggingface/status/1656334778407297027>

Huggingface 2025: smolagents



[README](#) [Code of conduct](#) [Apache-2.0 license](#)

[license](#) [Apache-2.0](#) [website](#) [online](#) [release](#) [v1.9.2](#) [Contributor Covenant](#) [v2.0 adopted](#)

 **smolagents**

A smol library to build great agents!

`smolagents` is a library that enables you to run powerful agents in a few lines of code. It offers:

- 🌟 **Simplicity:** the logic for agents fits in ~1,000 lines of code (see [agents.py](#)). We kept abstractions to their minimal shape above raw code!
- 💻 **First-class support for Code Agents.** Our [CodeAgent](#) writes its actions in code (as opposed to "agents being used to write code"). To make it secure, we support executing in sandboxed environments via [E2B](#).
- 🔗 **Hub integrations:** you can [share/pull tools to/from the Hub](#), and more is to come!
- 🌐 **Model-agnostic:** smolagents supports any LLM. It can be a local `transformers` or `ollama` model, one of [many providers on the Hub](#), or any model from OpenAI, Anthropic and many others via our [LiteLLM](#) integration.
- 👁️ **Modality-agnostic:** Agents support text, vision, video, even audio inputs! Cf [this tutorial](#) for vision.
- 🔧 **Tool-agnostic:** you can use tools from [LangChain](#), [Anthropic's MCP](#), you can even use a [Hub Space](#) as a tool.

Full documentation can be found [here](#).