

# LLMs in Eigenregie

Individualisierung & Hosting von KI-Modellen

IT After Work

**Marian Lambert**



## **Drogeriemarkt dm führt firmeneigenes ChatGPT ein**

- Handelsblatt, 25.08.2023

## **McKinsey & Company launches internal generative AI tool Lilli**

- consultancy-me.com, 29.08.2023

## **Mercedes-Benz Direct Chat: an internal ChatGPT application for employees**

- businesswire.com, 30.10.2023

## **BoschGPT: Aleph Alpha will bis Jahresende KI für Bosch-Mitarbeitende einführen**

- t3n.de, 14.08.2023



## Warum denn überhaupt ein „eigenes“ LLM?

- Datenschutz und -sicherheit
- Anpassung auf unternehmensinterne Anwendungsfälle
- Leistungssteigerung durch Zugriff auf unternehmensinterne Daten
- Volle Kontrolle über Modell
- Integration in andere IT-Systeme
- Niedrigere OPEX
- ...

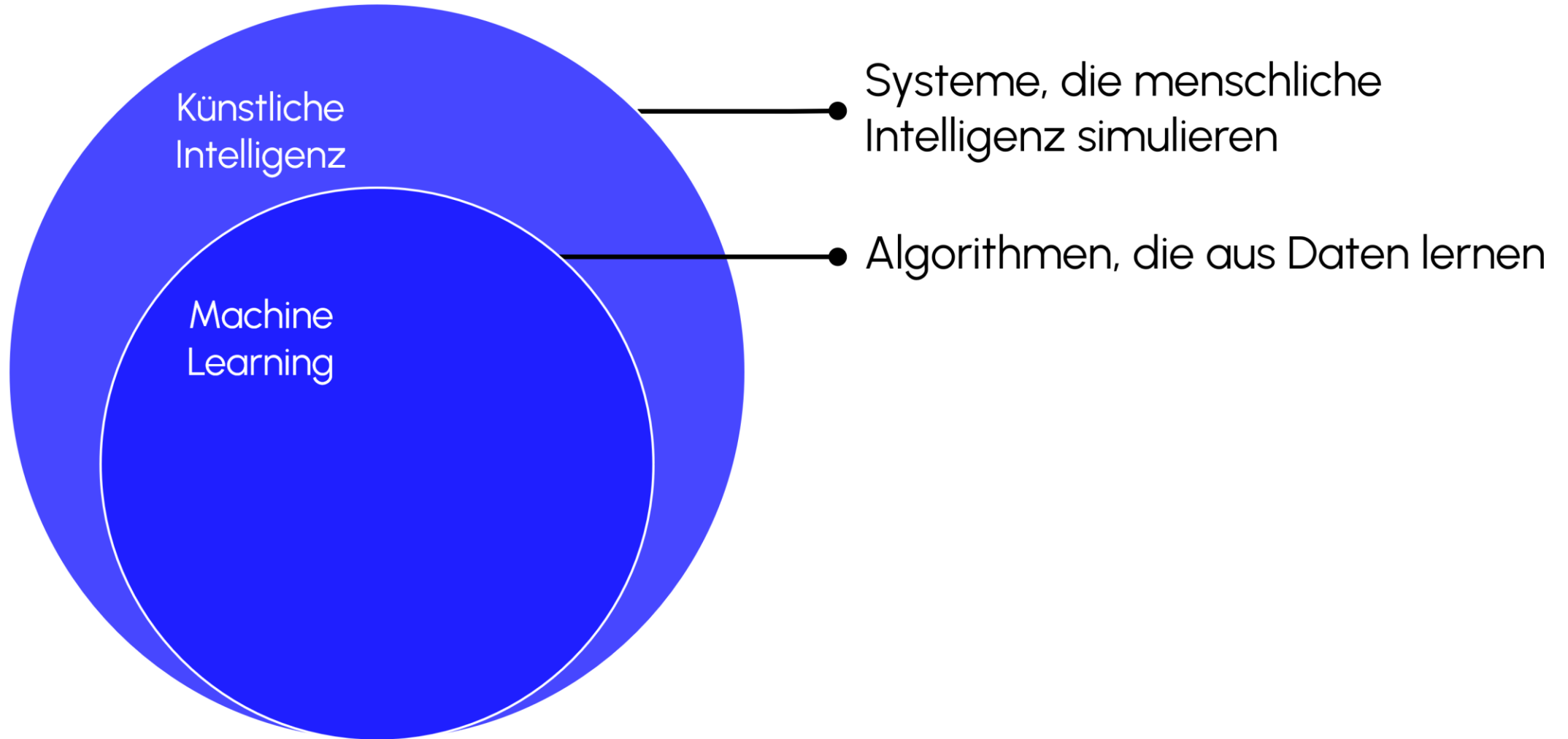


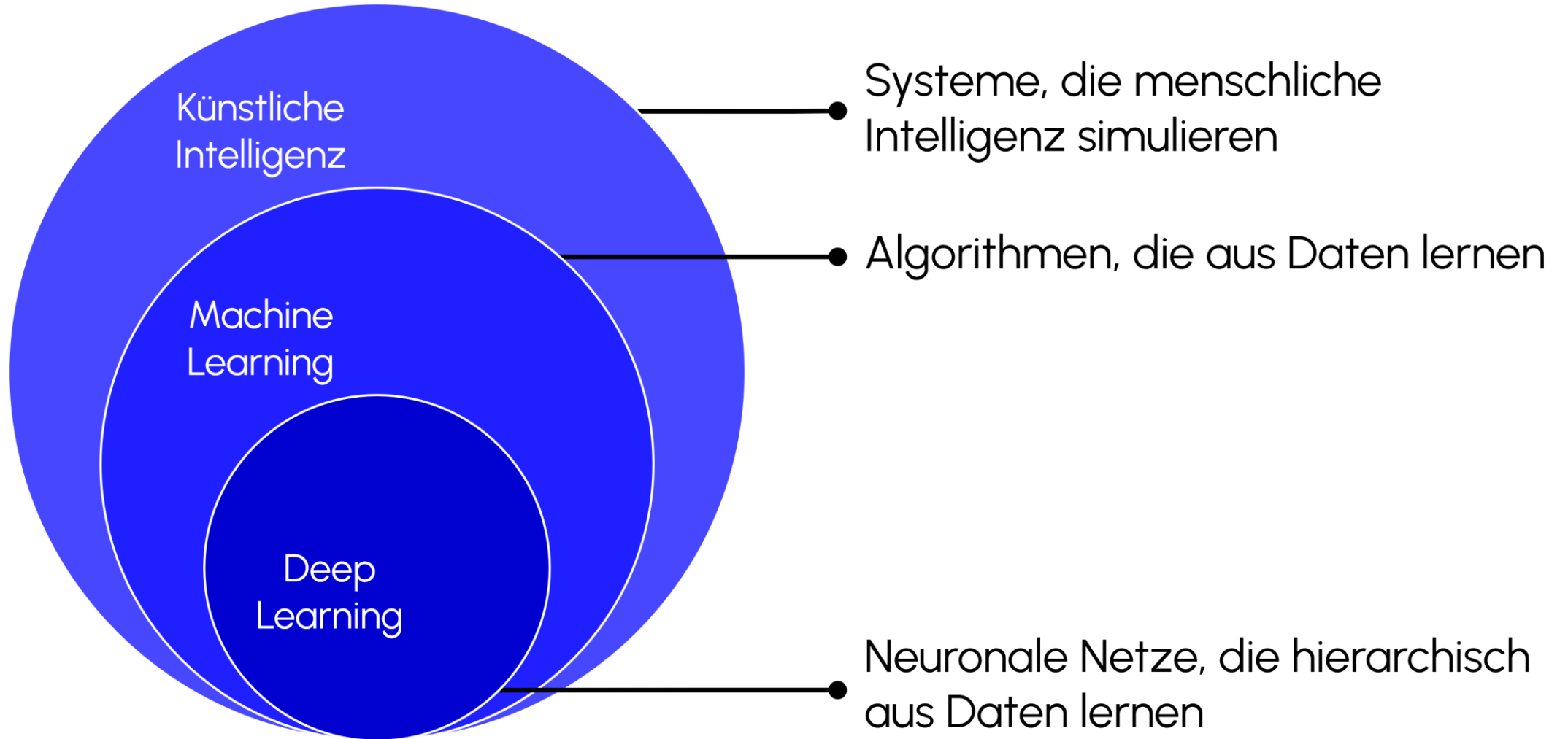
Wie bekomme ich nun mein  
eigenes LLM?

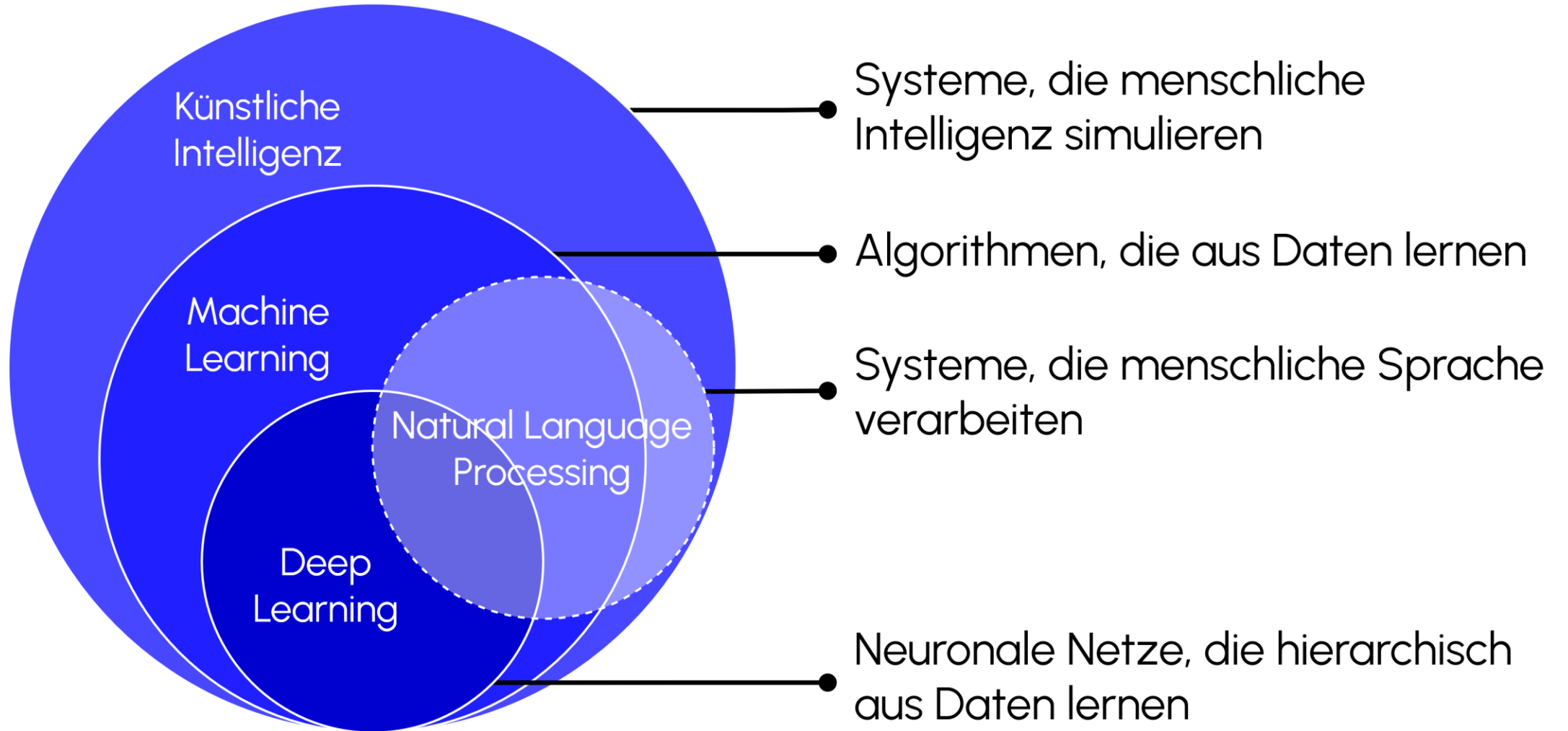


Künstliche  
Intelligenz

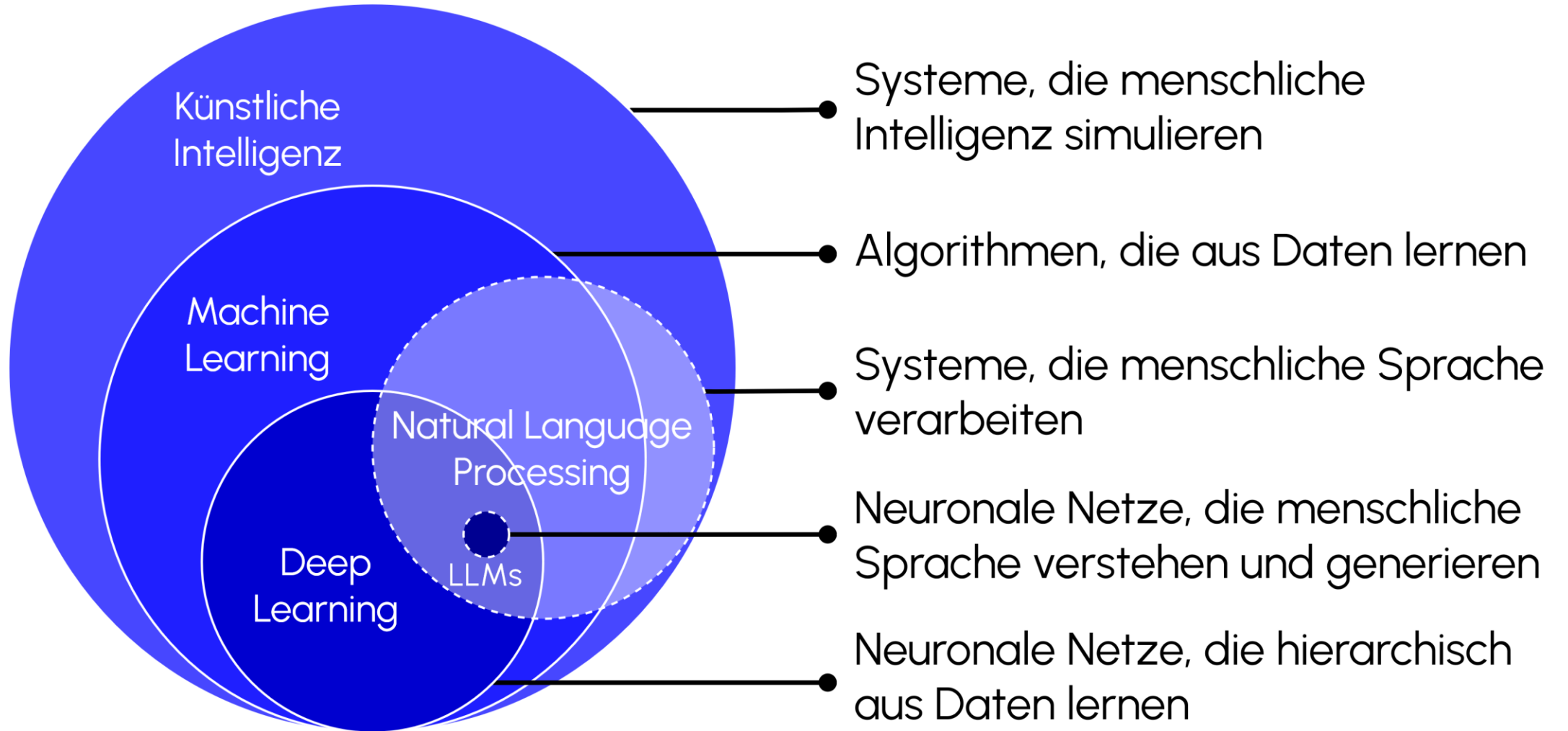
Systeme, die menschliche  
Intelligenz simulieren





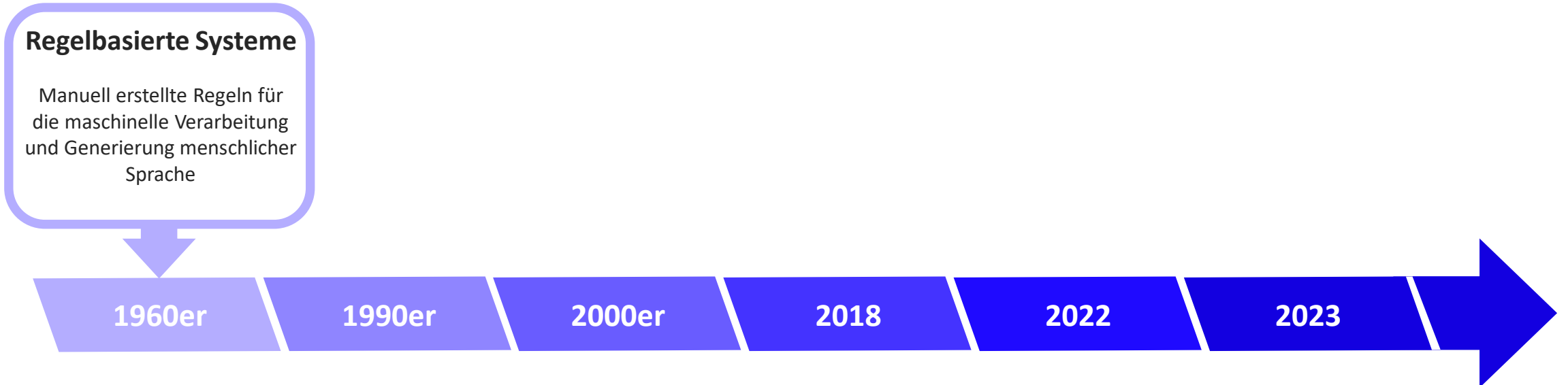








# Von ELIZA zu ChatGPT und weiter...





- 1966 von Joseph Weizenbaum entwickelt
- Regelbasierte Funktionsweise: Erkennung von Themengebieten und Antwort mit Standardphrasen

„Ich habe ein Problem mit meinem Vater.“

„Erzählen Sie mir mehr über Ihre Familie!“

„Krieg ist der Vater aller Dinge.“

„Erzählen Sie mir mehr über Ihre Familie!“

## Turing-Test: OpenAIs ChatGPT verliert gegen Sprachmodell aus den 60ern

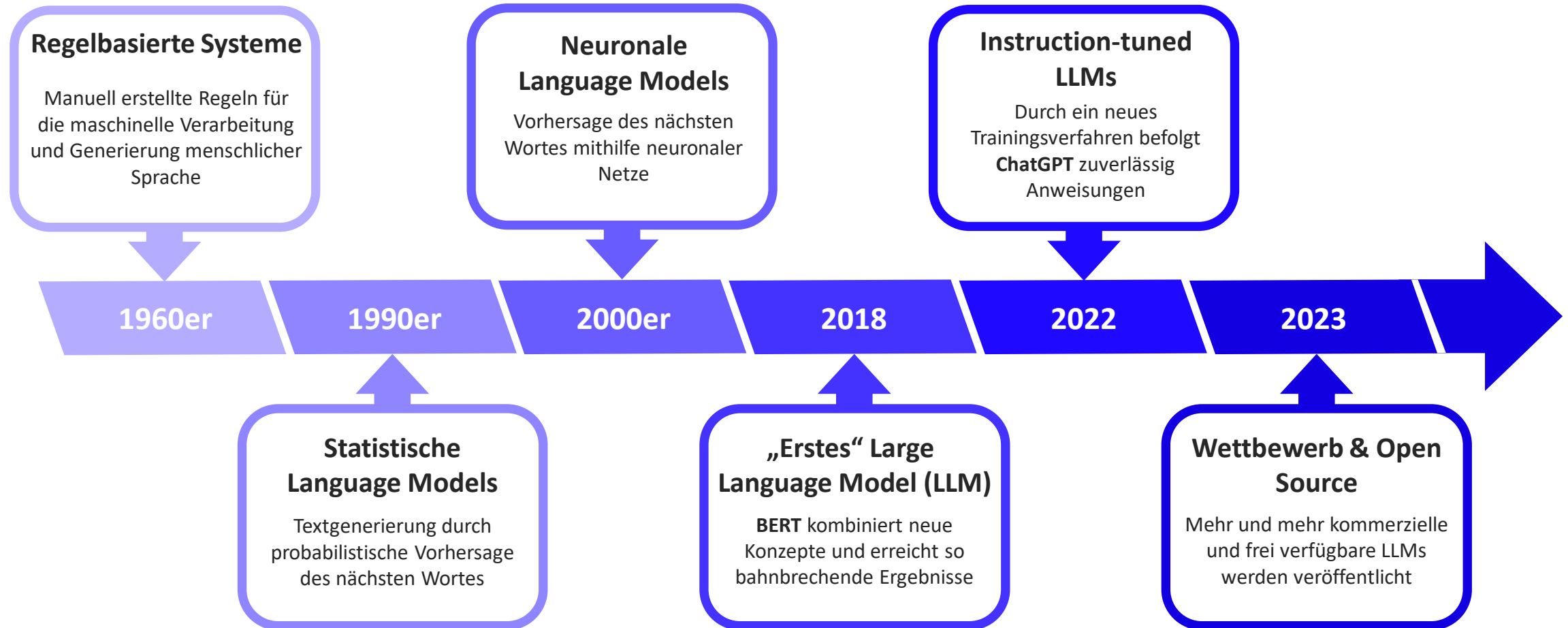
```
Welcome to
          EEEEE LL   IIII ZZZZZZ  AAAAA
          EE  LL   II   ZZ   AA  AA
          EEEEE LL   II   ZZZ  AAAAAA
          EE  LL   II   ZZ   AA  AA
          EEEEE LLLLL IIII ZZZZZZ  AA  AA

Eliza is a mock Rogerian psychotherapist.
The original program was described by Joseph Weizenbaum in 1966.
This implementation by Norbert Landsteiner 2005.

ELIZA: Is something troubling you ?
YOU:   Men are all alike.
ELIZA: What is the connection, do you suppose ?
YOU:   They're always bugging us about something or other.
ELIZA: Can you think of a specific example ?
YOU:   Well, my boyfriend made me come here.
ELIZA: Is it important to you that your boyfriend made you come here ?
YOU:   He says I'm depressed much of the time.
ELIZA: I am sorry to hear that you are depressed.
YOU:   It's true. I am unhappy.
ELIZA: Can you explain what made you unhappy ?
YOU:   █
```



# Von ELIZA zu ChatGPT und weiter...

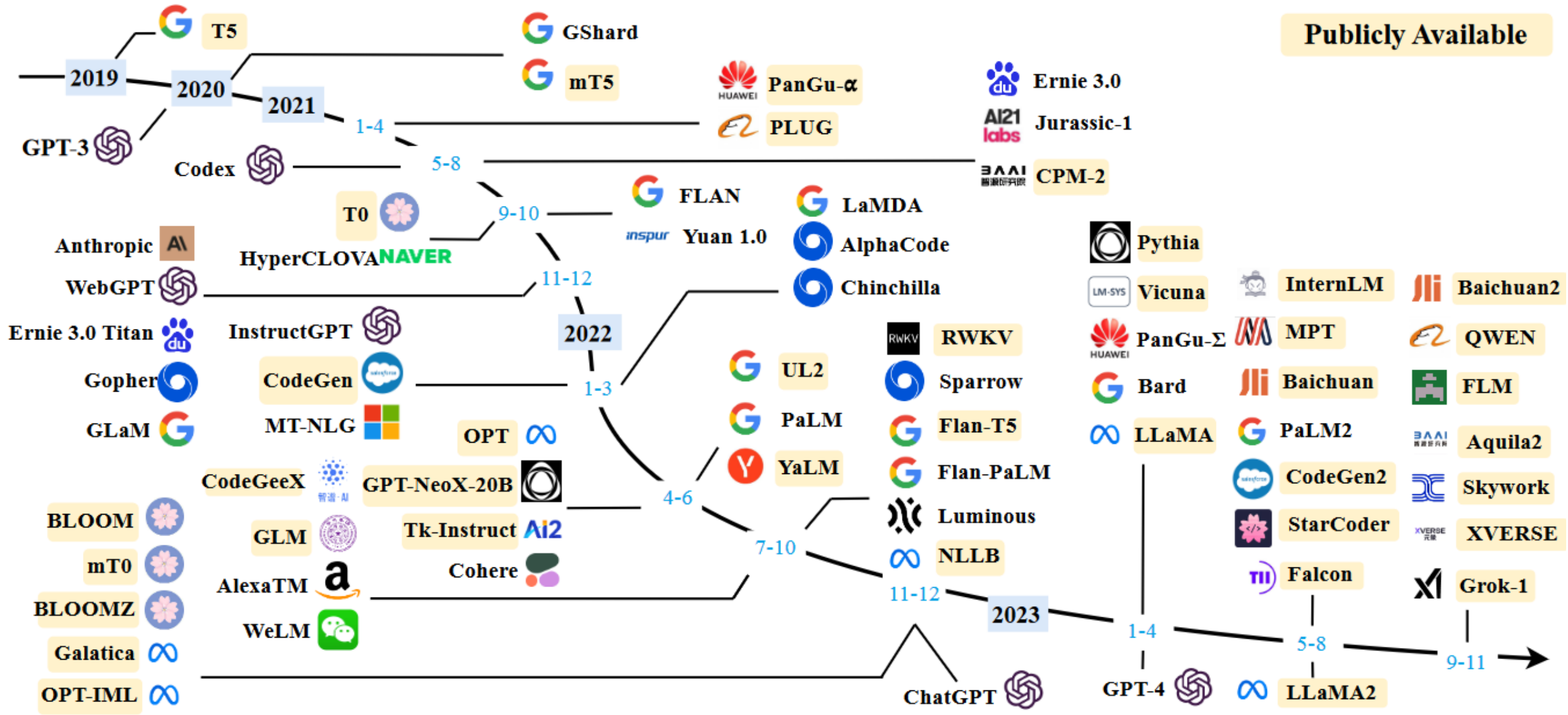




# Demo



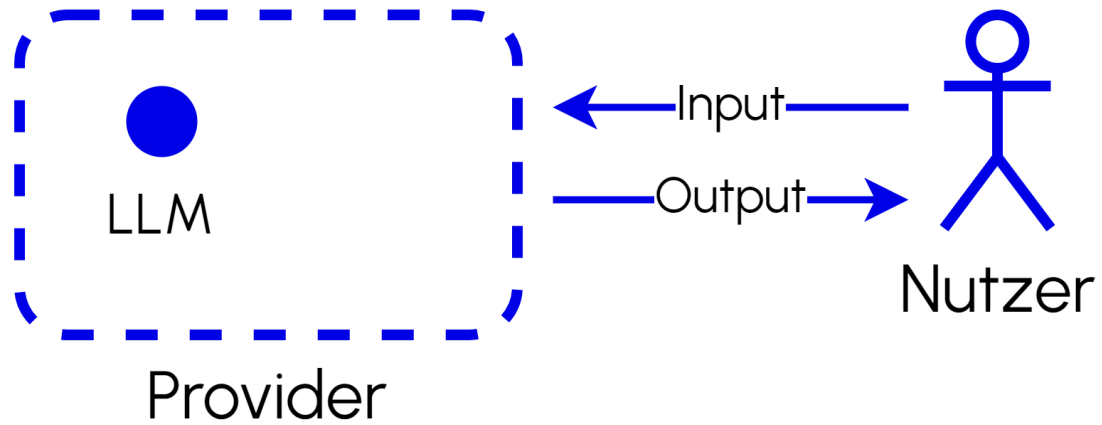
# Das richtige LLM als Basis



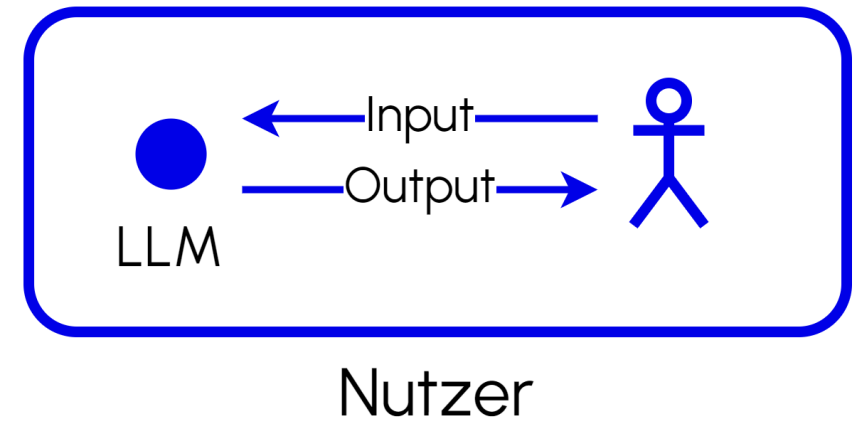


# Proprietär vs. Open Source

Closed Source



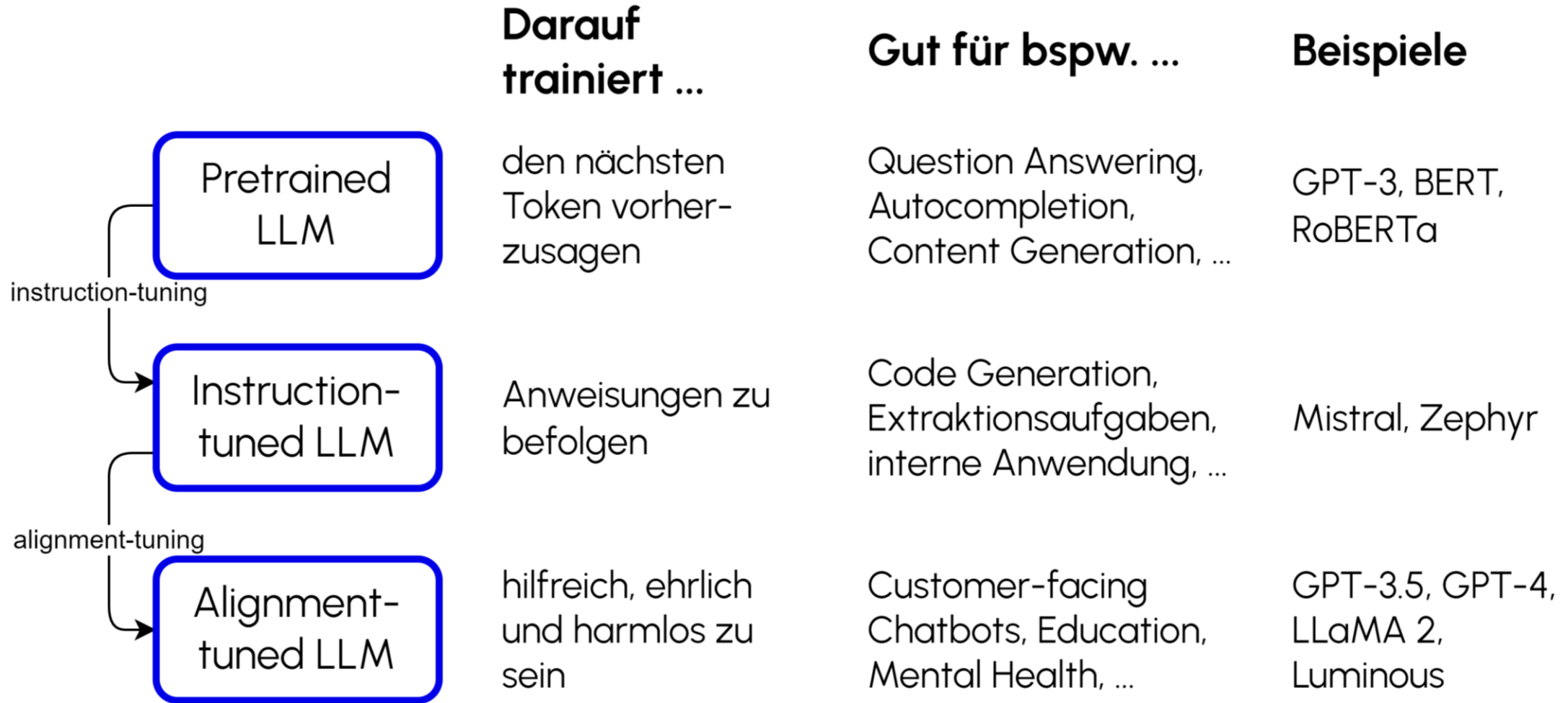
Open Source







# Pretrained vs. Instruction-tuned vs. Alignment-tuned





## Weitere Faktoren

- **Leistungsfähigkeit** (Sprachverständnis, Reasoning, Kreativität)
- **Latenz**
- **Datenschutz**
- **Kosten**
- **Kontextgröße**
- **Datensatz** (Bias, Fairness, Transparenz)
- **Rechtliche Anforderungen** (bspw. AI Act)
- ...



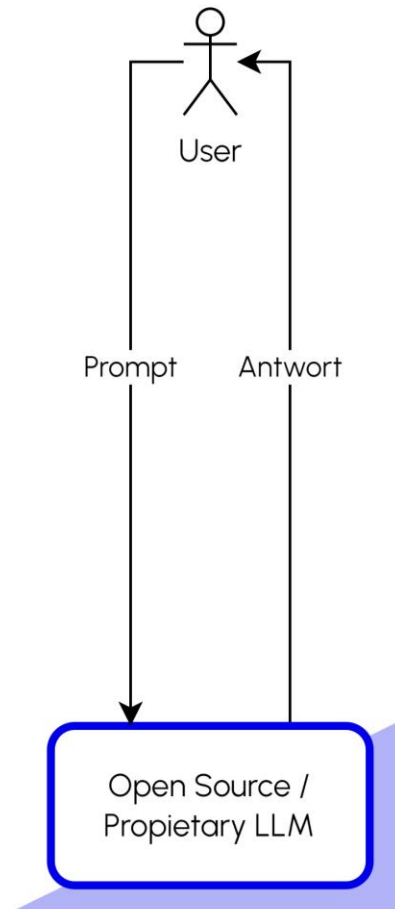
# Bekannte LLMs im Überblick

Model	Organization	Typ	Context Window (Tokens)	Size (Parameters)	License
GPT-3.5-turbo	OpenAI	IT	4097, 16,385	175B	Proprietär
<b>GPT-4</b>	<b>OpenAI</b>	<b>IT</b>	<b>8192, 32768</b>	<b>1.8T (est.)</b>	<b>Proprietär</b>
LLaMA 2	Meta	PT/IT	4096	7B, 16B, 70B	NC Custom
StableBeluga 2	StabilityAI	IT	4096	70B	CC BY-NC-4.0
Claude 2	Anthropic	IT	200k	130B (est.)	Proprietär
Luminous	Aleph Alpha	PT/IT	2048	13B, 30B, 70B	Proprietär
PaLM 2	Google	IT	8000	340B (est.)	Proprietär
Falcon	Technology Innovation Institute	PT/IT	2048	7B, 40B, 180B	Apache-2.0
Mistral	Mistral AI	PT/IT	8000	7B	Apache-2.0
Zephyr	HuggingFace	PT/IT	8000	7B	MIT
Platypus 2	garage-bAInd	PT/IT	4096	7B, 13B, 70B	CC BY-NC-4.0
<b>Godzilla 2</b>	<b>Maya Philippines</b>	<b>IT</b>	<b>4096</b>	<b>30B, 70B</b>	<b>CC BY-NC-4.0</b>
Orca 2	Microsoft	IT	4096	7B, 13B	MRL
LeoLM	LAION	IT	8000	7B, 13B	NC Custom
Gemini	Google	IT	32k	?	Proprietär

PT=Pretrained; IT=Instruction/alignment-tuned

# Anpassung an das Unternehmen

## Prompt Engineering

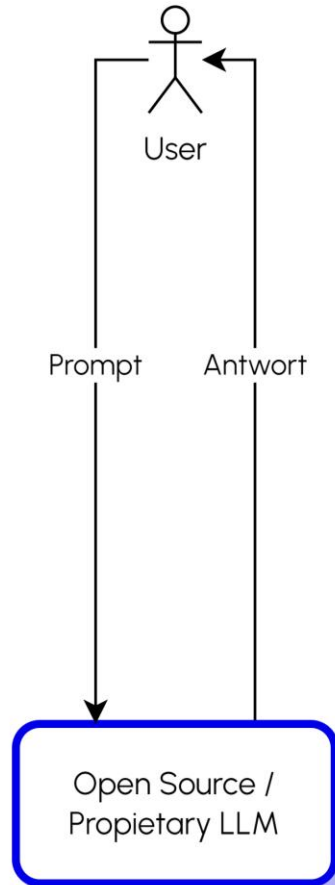


**Geeignet für ...**

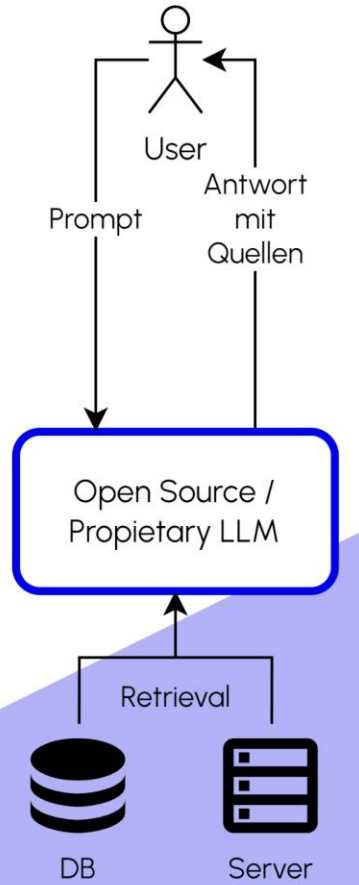
Aufgaben, die generische Fähigkeiten und Informationen erfordern



### Prompt Engineering



### Retrieval-Augmented Generation

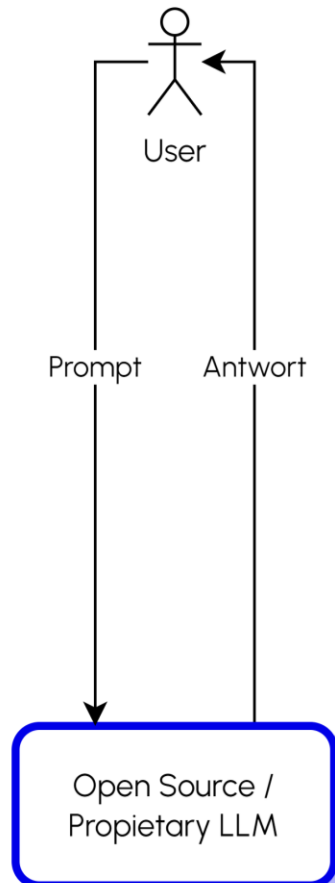


**Geeignet für ...**

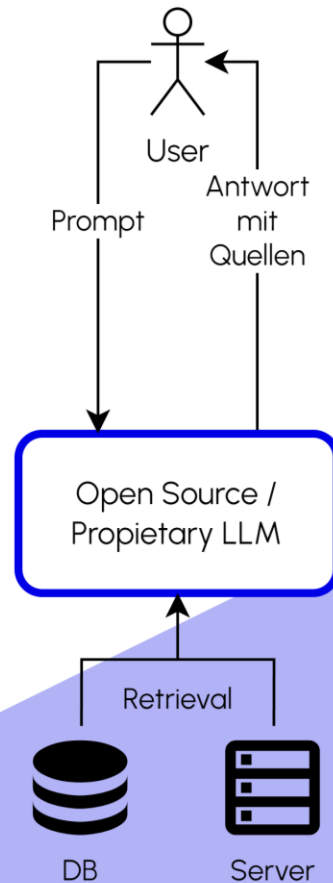
Aufgaben, die generische Fähigkeiten und Informationen erfordern

Aufgaben, die konkrete interne Informationen erfordern

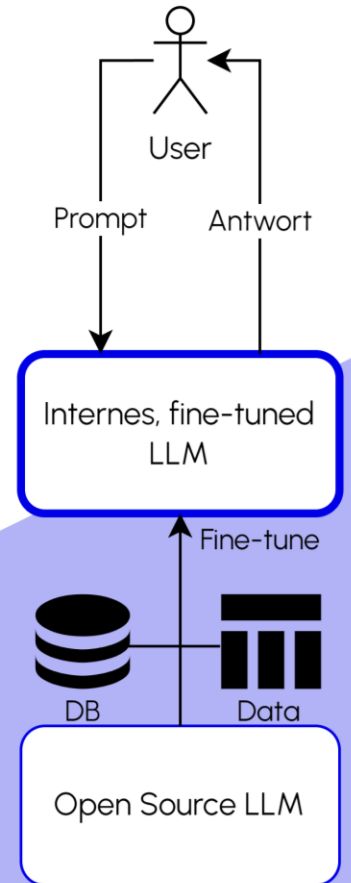
### Prompt Engineering



### Retrieval-Augmented Generation



### Fine-tuned LLM

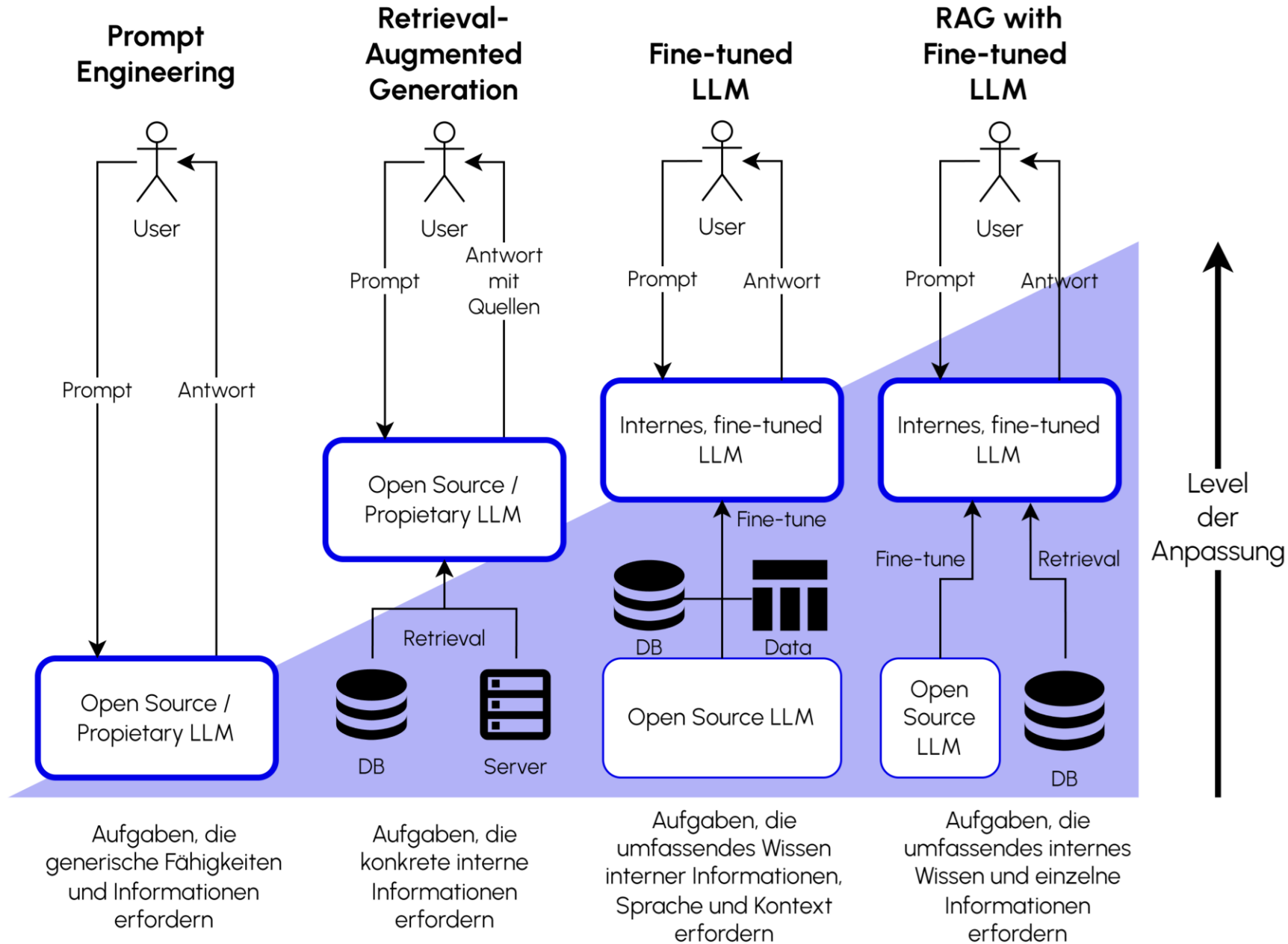


**Geeignet für ...**

Aufgaben, die generische Fähigkeiten und Informationen erfordern

Aufgaben, die konkrete interne Informationen erfordern

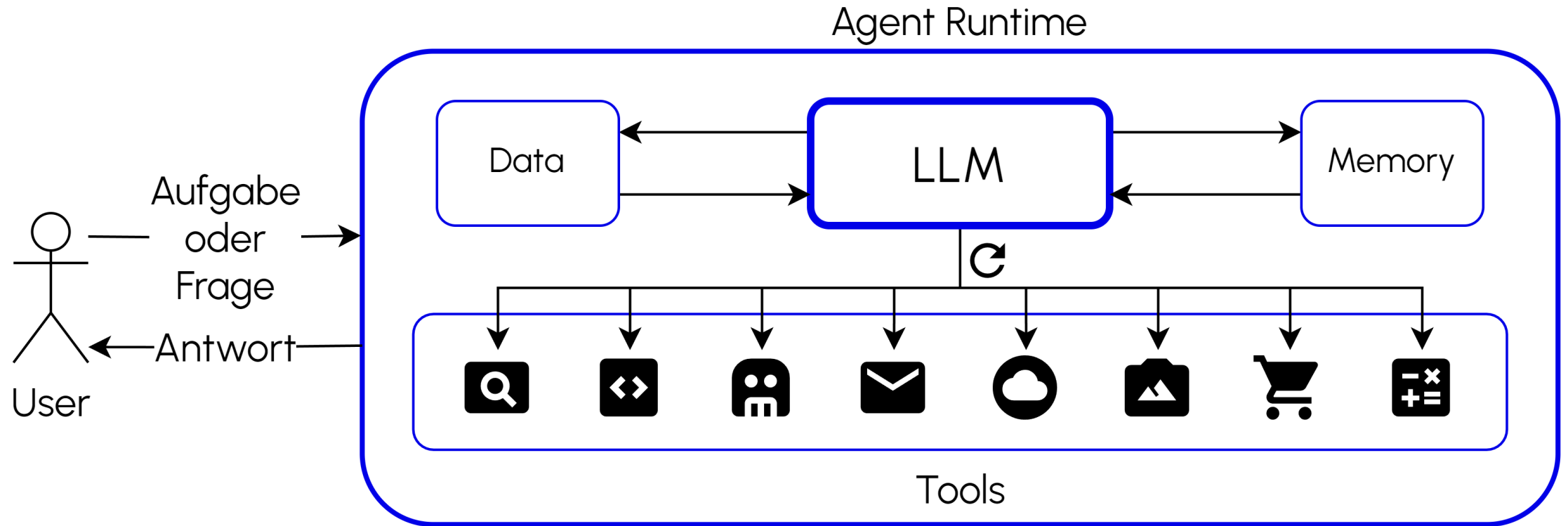
Aufgaben, die umfassendes Wissen interner Informationen, Sprache und Kontext erfordern







# Von LLMs zu Agenten



# Contact Me!



## Marian Lambert

Senior Consultant AI & Software

 [marian.lambert@xpace.de](mailto:marian.lambert@xpace.de)

 [www.xpace.de](http://www.xpace.de)





# Danke!

Fragen?