



Beyond Generative AI

Large Language Models (LLMs) Explained

Dr. Daniel Lehner

Expert AI Knowledge Transfer @ AI Factory Austria AI:AT

Dr. Daniel Lehner



Expert AI Knowledge Transfer @ AI Factory Austria AI:AT

AI in **manufacturing** (through digital twins)

- IT Consultant and Trainer
- Business Informatics Researcher (441 citations) [1]

AI for **general public**

- University lectures AI for law/business administration
- VHS courses AI for general public [2]
- Certified AI Manager trainings [3]

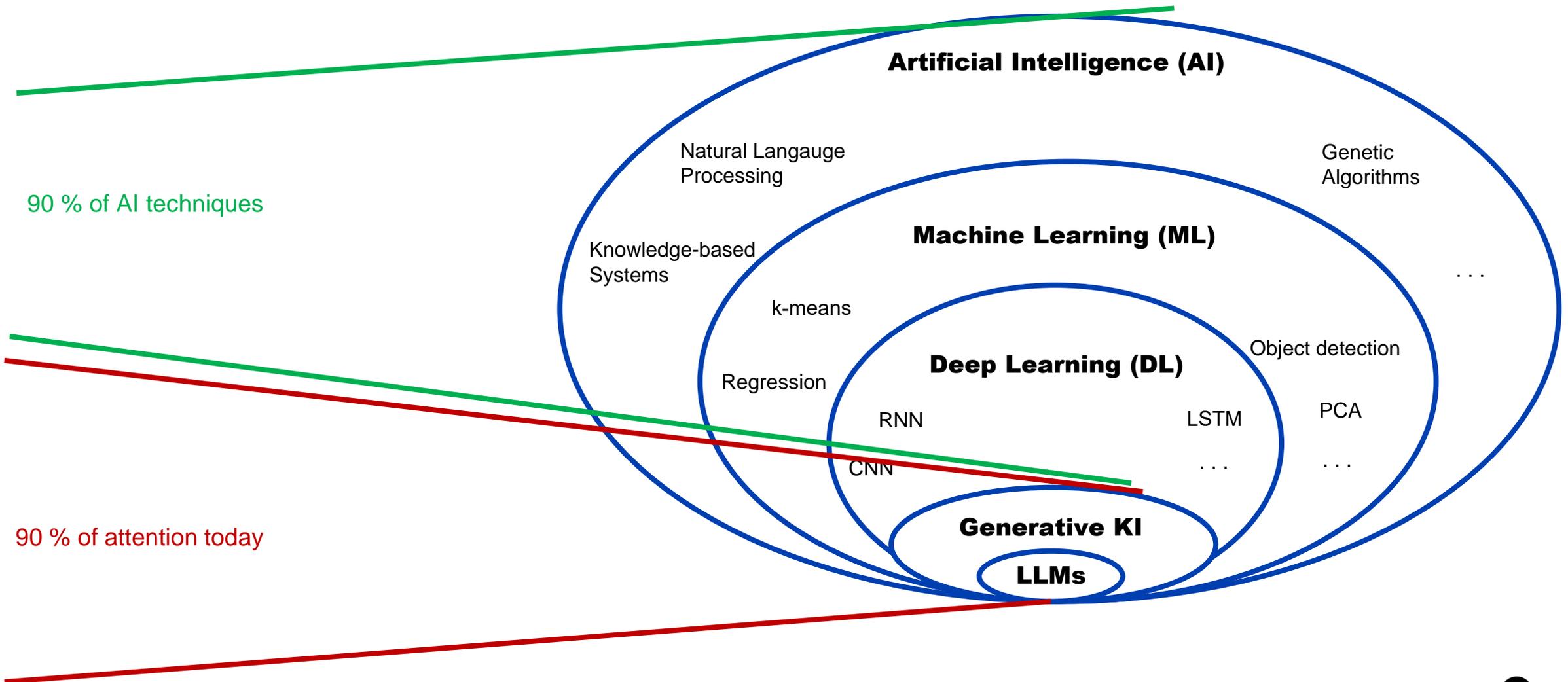
[1] https://scholar.google.com/citations?user=TGGaQ_0AAAAJ

[2] https://vhskurs.linz.at/index.php?kathaupt=18&suchesetzen=false&kfs_dozentid=139372

[3] <https://tecnovy.com/de/tecnovy/certified-ai-manager>



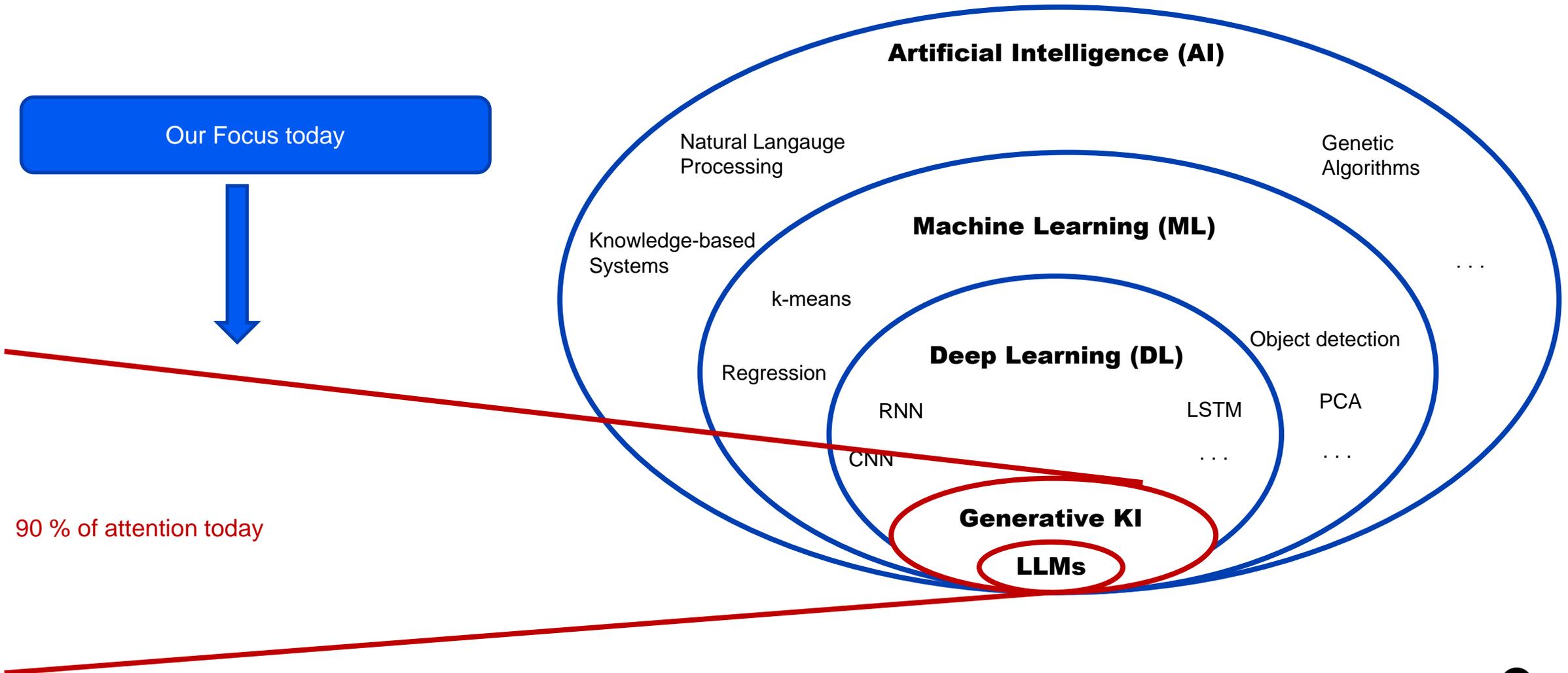
„AI Tunnel View“



Source: https://www.linkedin.com/posts/jaylatta_llm-tunnel-vision-weve-fallen-into-the-activity-7386990451360387072-y6KI



„AI Tunnel View“



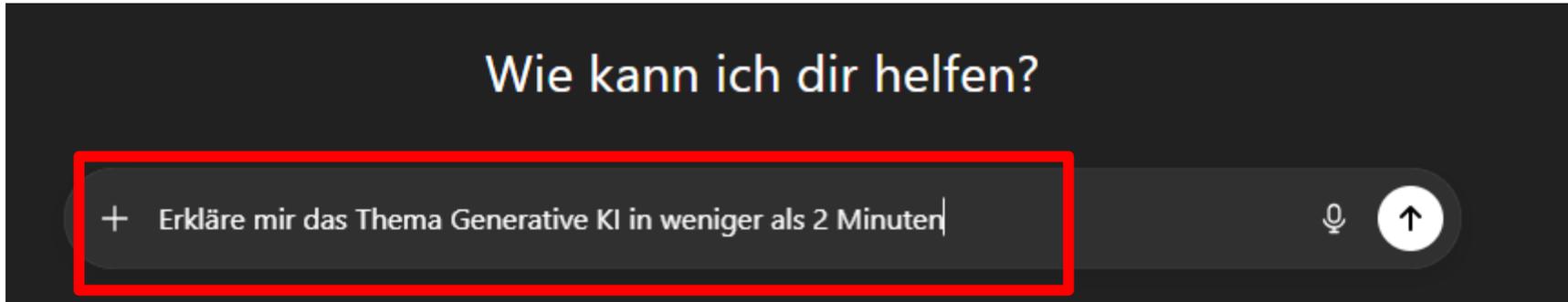
Source: https://www.linkedin.com/posts/jaylatta_ilm-tunnel-vision-weve-fallen-into-the-activity-7386990451360387072-y6KI



Chat Window



Prompt



Chatting

Erkläre mir das Thema Generative KI in weniger als 2 Minuten

Generative KI ist eine Form von Künstlicher Intelligenz, die nicht nur Daten analysiert, sondern selbst neue Inhalte erzeugen kann – zum Beispiel Texte, Bilder, Musik oder sogar Code.

Sie basiert auf **Modellen**, die mit riesigen Mengen an Beispielen trainiert werden. Diese lernen Muster, Strukturen und Zusammenhänge und können daraus Neues ableiten. Ein Sprachmodell wie ich kann also ganze Texte verfassen, während ein Bildmodell realistische Bilder malt – auch wenn es diese vorher nie gesehen hat.

Typische Anwendungen sind **Chatbots, Übersetzung, Design, Produktentwicklung, Medizin oder Spiele**. Chancen: Kreativität, Automatisierung, Zeitersparnis. Risiken: Fehlinformationen, Urheberrechtsfragen, Bias.

Kurz: Generative KI ist ein „kreativer Assistent“, der Ideen und Inhalte aus Daten hervorbringt – schnell, flexibel und vielseitig.

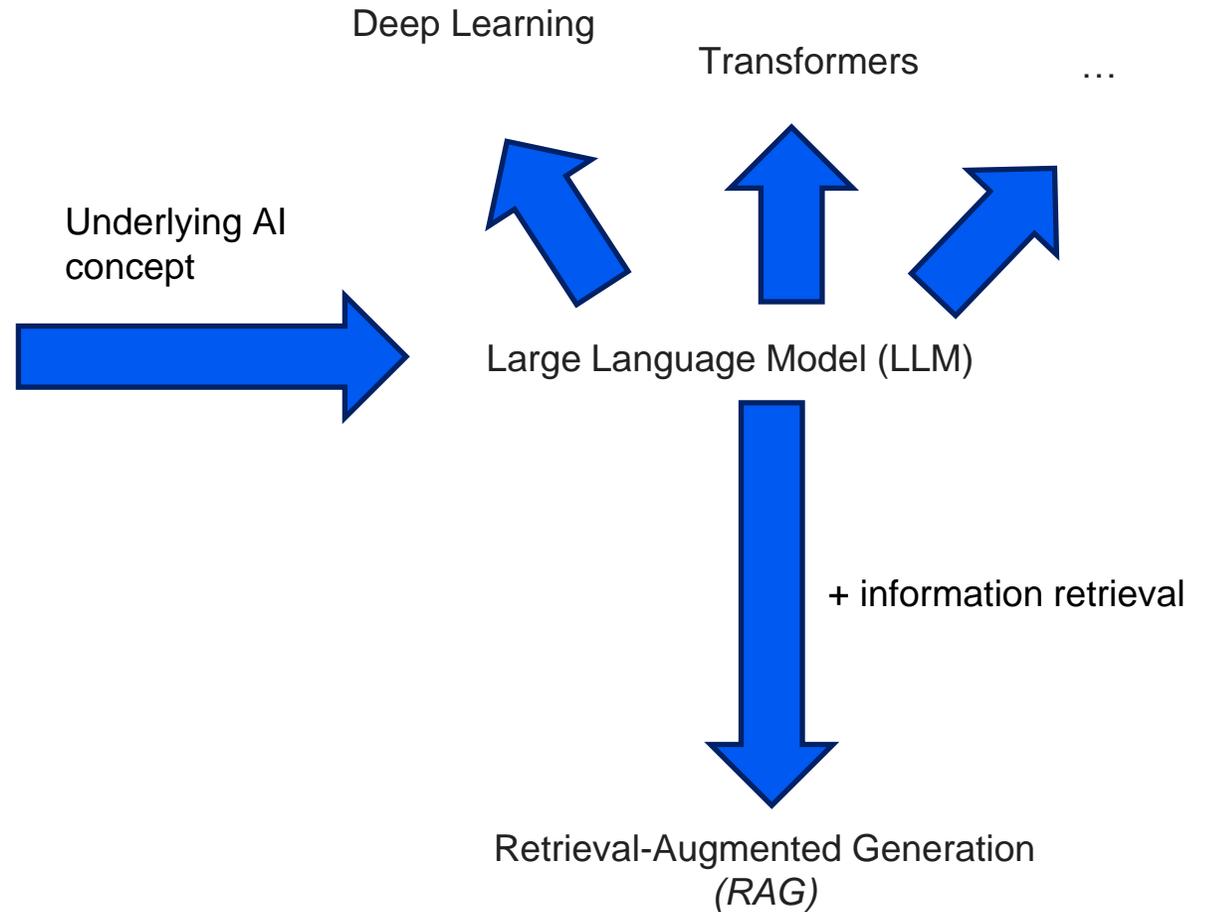
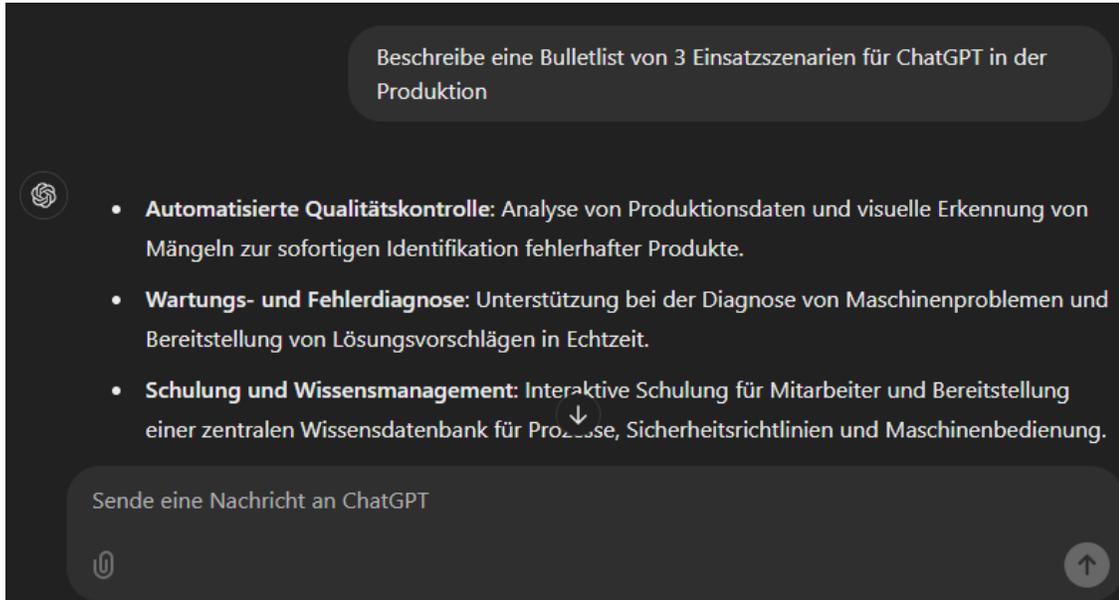
👉 Soll ich dir auch ein ganz einfaches Bild/Schema dazu erstellen, damit du es dir besser vorstellen kannst?



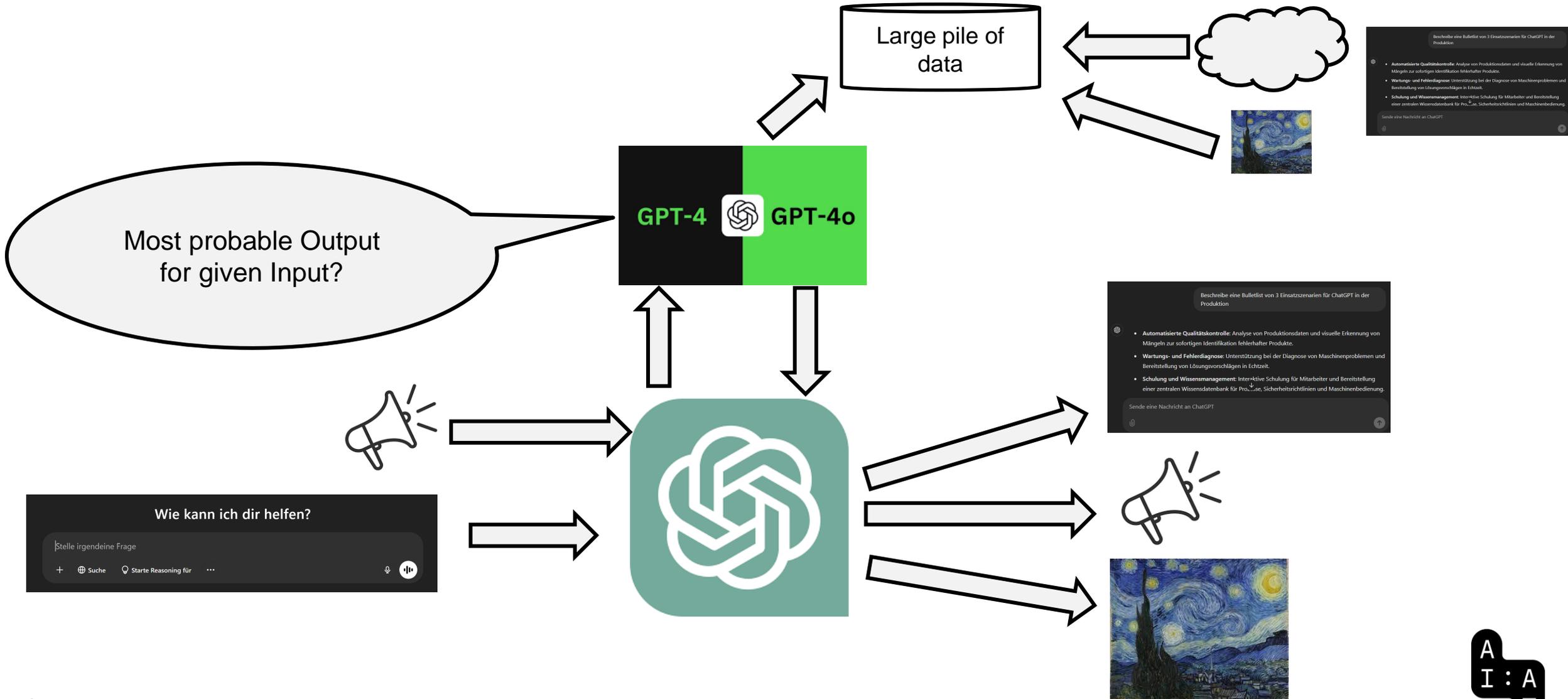
+ Stelle irgendeine Frage



AI Chatbots



Generative AI



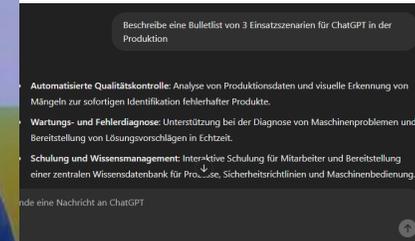
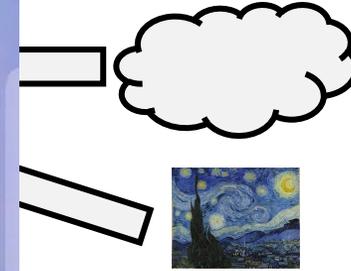
Generative AI

Most probable Output
for given Input?



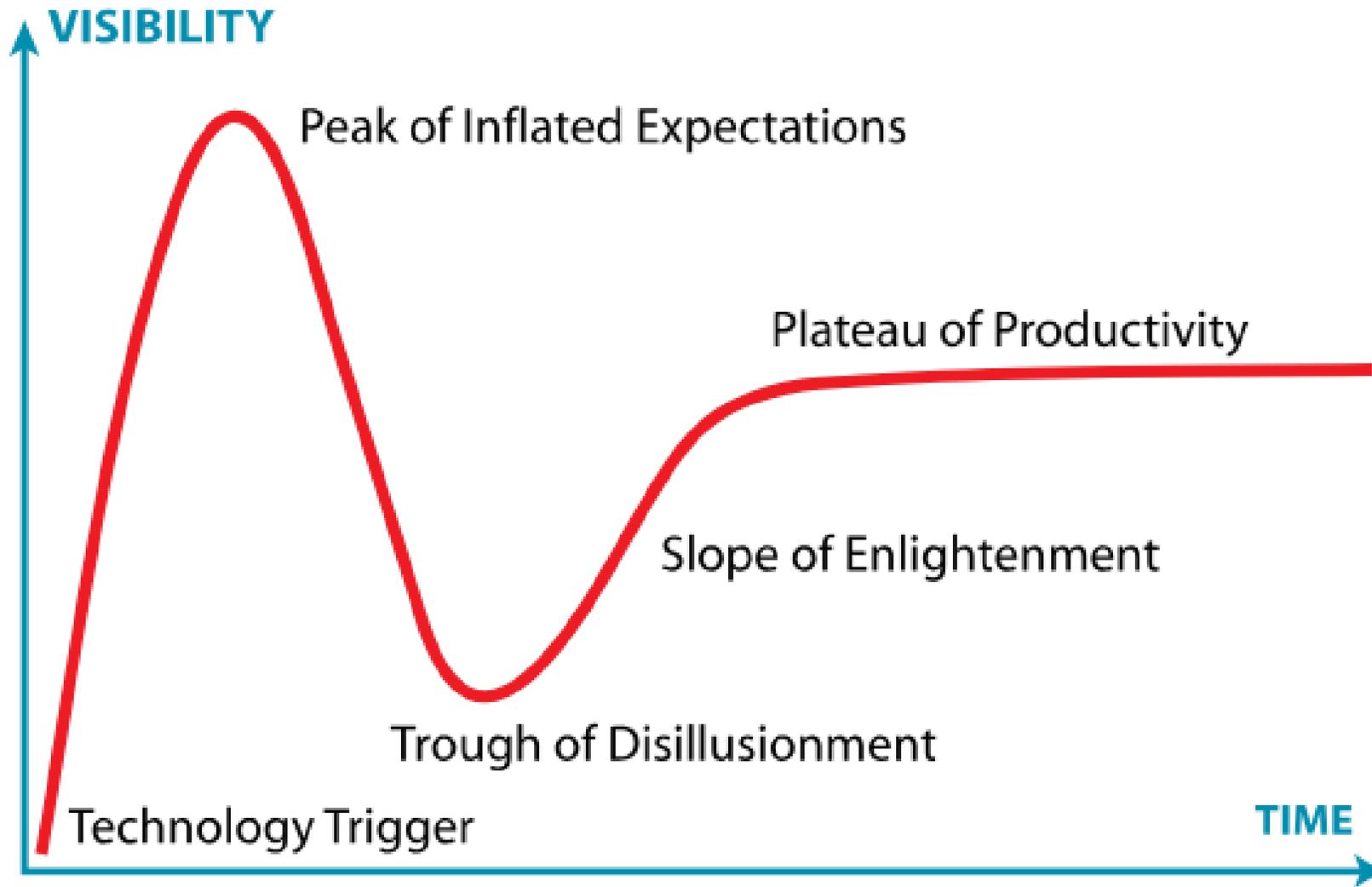
unusual_whales @unusual_whales
ChatGPT has hit 800 million weekly active users.

Shaan Puri @ShaanVP
Turns out the greatest startup idea of the last 10 years was a chatbot where you ask a question and it answers with something a guy wrote on Reddit 8 years ago

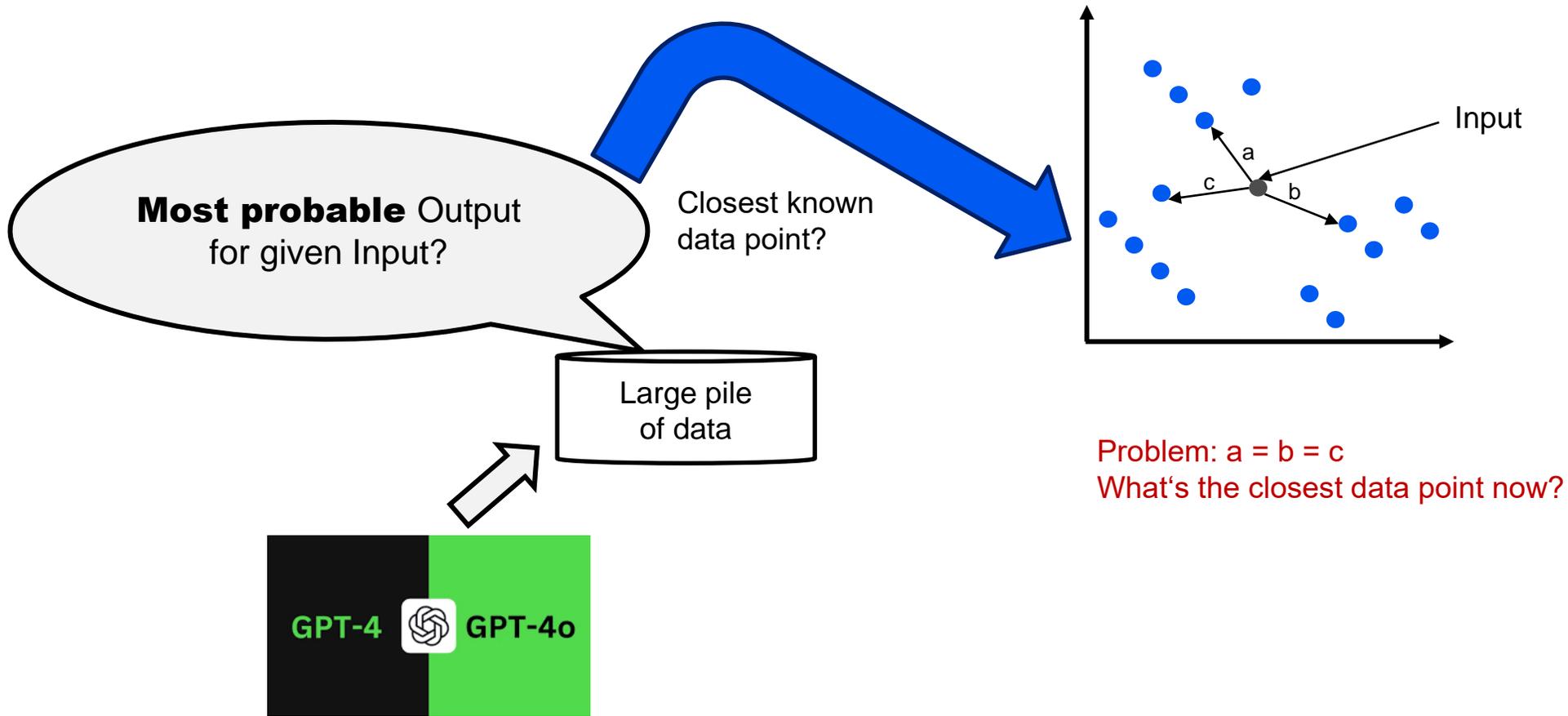


Expectation Management Generative Artificial Intelligence

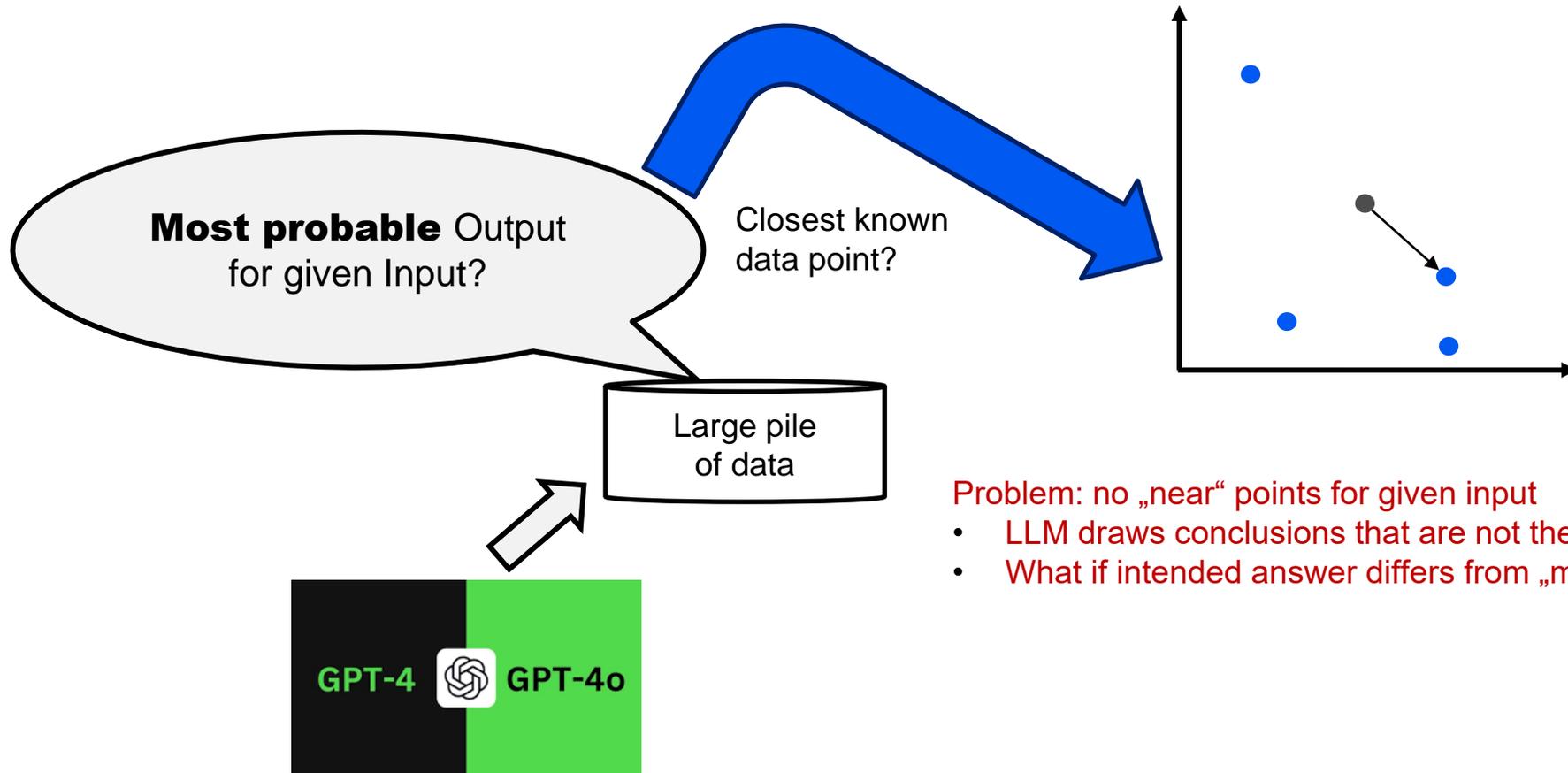




Challenge 1: Non-Deterministic Result



Challenge 2: „Halluzinations“



Problem: no „near“ points for given input

- LLM draws conclusions that are not there => still better than no answer
- What if intended answer differs from „most probable“ output?

Challenges of Generative AI

Non-Deterministic Result

Halluzinations



Who still wants to ride a plane steered by AI?

You actually did!

Solution: Expert Systems

Interesting Article: <https://www.science.org/doi/10.1126/science.aea3922>

Working with LLMs



Which LLM to choose?



The market is too fast!

Jede Woche neue Modelle + Ergebnisse

Jeden Monat neue Tools

Riesige Menge an Anforderungen

<https://theresanaiforthat.com/>

The Generative AI Market Map v3



<https://medium.com/@maximilian.vogel/5000x-generative-ai-intro-overview-models-prompts-technology-tools-comparisons-the-best-a4af95874e94>



Major AI chatbots

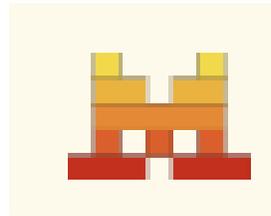
Frontrunner



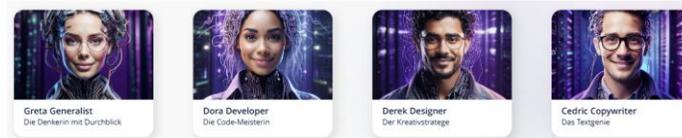
ChatGPT



European Alternatives

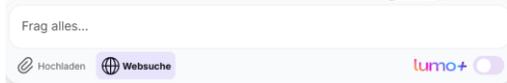


<https://mistral.ai>



<https://ai.ionos.de/explore>

Hey, ich bin Lumo.
Frag mich alles.
Es ist vertraulich.



<https://lumo.proton.me/>

Specific Use Cases



Deep-Dive ChatGPT

 GPT-4o — Der Allrounder für Alltägliches: Brainstorming, Zusammenfassungen, E-Mails, Creative Content. Unterstützt praktisch alle Funktionen und Eingabetypen.

 GPT-4.5 — Der Kreative: Emotionale Intelligenz, klare Kommunikation, Kreativität und ein intuitiverer Ansatz beim Brainstorming.

 o4-mini — Der Sprinter: Rasend schnelle Bearbeitung von STEM-Anfragen, Programmierung und visueller Analyse. Perfekt für schnelle technische Aufgaben!

 o4-mini-high — Der gründliche Techniker: Wie o4-mini, aber mit längerer Denkzeit für höhere Genauigkeit bei komplexen Aufgaben.

 o3 — Der Stratege: Für komplexe, mehrstufige Aufgaben wie strategische Planung, detaillierte Analysen, umfangreiches Coding und fortgeschrittene mathematische Probleme.

 o1 pro mode — Das Präzisionswerkzeug: Braucht etwas länger zum Nachdenken, liefert aber die nötige Genauigkeit für komplexe Aufgaben.

https://www.linkedin.com/posts/timospringer_welches-chatgpt-modell-f%C3%BCr-welche-aufgabe-activity-7350404938033709056-sW-s
<https://help.openai.com/en/articles/11165333-chatgpt-enterprise-models-limits>



LLM Sandboxes

✦ Auto Pro ✓
AI that adapts—quick answers or advanced research and computation as needed

⚡ Express
Fast, reliable results for everyday tasks

📄 Custom
Choose from 20+ AI models and agents

Recents

- 📌 Daniel Writing Coach ↗
- 🌐 Albion LinkedIn Strategie ↗
- 📄 Research ↗
- 📄 Compute ↗
- 📄 Create ↗



Enhance your productivity with AI

🔍 Search...

See more ▾

Models

- New & Popular >
- OpenAI >
- Anthropic >
- Google >
- xAI >
- Alibaba >
- DeepSeek >
- Meta >



- AI Claude Opus 4.5 Pro
- AI Claude Opus 4.1 (Extended) Pro
- AI Claude Opus 4.1 Pro
- AI Claude Sonnet 4.5 (Extended) Pro
- AI Claude Sonnet 4.5 Pro
- AI Claude Sonnet 4 (Extended) Pro
- AI Claude Sonnet 4 Pro
- AI Claude Haiku 4.5

Tool: <https://you.com/>
Alternatives

- <https://abacus.ai/>
- <https://langdock.com/>
- ...



LLM Sandboxes

Playground

Iterate on and test prompts.

 Set up Evaluation    Start

Prompts ⓘ + Prompt

Load prompt ▾ gpt-5-mini  Save ▾

SYSTEM ▾
You are a chatbot.

HUMAN ▾
{question}

+ Message + Output Schema + Tool {x} f-string ▾

Inputs ⓘ

question Enter variable value...

Output ▾

Ctrl ↵ or click Start to generate...

<https://smith.langchain.com/o/cd3eb3e5-dac5-46b4-840b-4846c284f7c0/playground>



Use LLMs within software

```
from openai import OpenAI
```

pip install openai

```
client = OpenAI(api_key="...")
```

Get it from: <https://platform.openai.com/>

```
response = client.chat.completions.create(
```

```
    model="gpt-4o-mini",
```

```
    max_tokens=100,
```

```
    messages=[{"role": "user", "content": "IYOUR PROMT HERE"}]
```

```
)
```

```
print(response.choices[0].message.content)
```

Was sind die Top 3 Use Cases für KI in der Produktion?

Die Top 3 Use Cases für KI in der Produktion sind:

1. **Vorausschauende Wartung (Predictive Maintenance):** KI kann Sensordaten von Maschinen und Anlagen analysieren, um frühzeitig Anzeichen von Abnutzung oder Fehlfunktionen zu erkennen. So lassen sich Ausfälle vorhersagen und Wartungsmaßnahmen rechtzeitig einleiten, was die Betriebszeiten erhöht und unerwartete Ausfälle verringert.
2. **Qualitätskontrolle und Fehlererkennung:** Durch den Einsatz von Computer Vision und maschinellem Lernen können Produktionsprozesse in Echtzeit überwacht und automatisch Fehler erkannt werden. Dies führt zu einer besseren Produktqualität, reduziert Ausschuss und minimiert die Notwendigkeit für manuelle Inspektionen.
3. **Optimierung der Produktionsplanung und -steuerung:** KI-basierte Algorithmen können Produktionsabläufe optimieren, indem sie verschiedene Faktoren wie Materialverfügbarkeit, Produktionskapazitäten und Nachfrageschwankungen berücksichtigen. So lässt sich die Effizienz maximieren, die Lieferzeiten verkürzen und die Kosten senken.

Diese Anwendungen tragen dazu bei, die Effizienz, Qualität und Flexibilität in der Produktion erheblich zu steigern.

<https://campus.datacamp.com/courses/working-with-the-openai-api/introduction-to-the-openai-api>

Use Cases involving „LLM APIs“

Add LLM responses to your own software (e.g., chatbots)

Improving/completing prompts („System Prompts“)

Enrich prompts with your own data („RAG“)



Course in March:

<https://ai-at.eu/training/foundations-of-llm-mastery-prompt-engineering-essentials-2/>



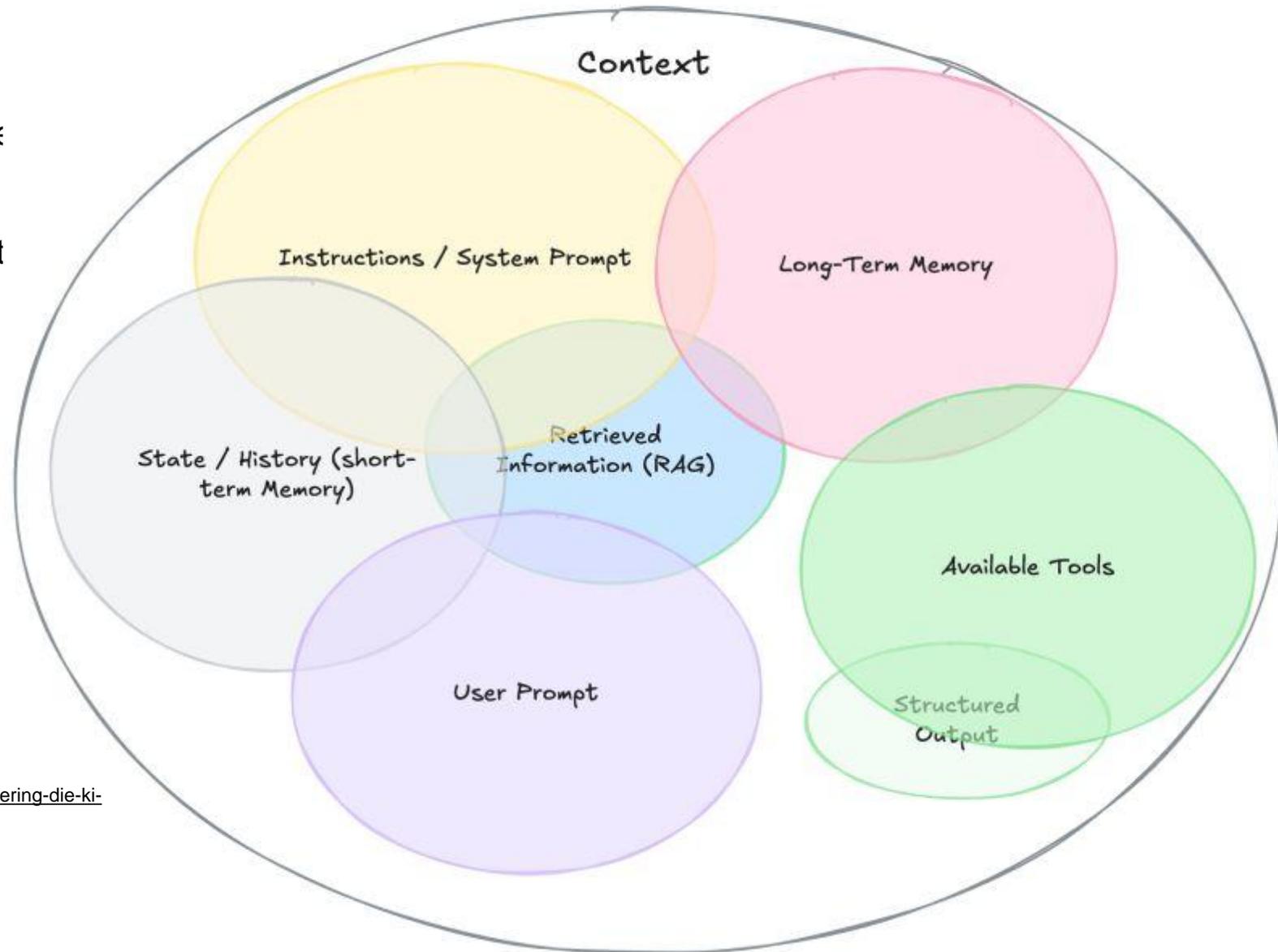
Anwendungsfälle von "IMADIA"

Context Engineering

Add LLM responses to your own software (€

Prompts vervollständigen/verbessern („Syst

Enrich prompts with your own data („RAG“)

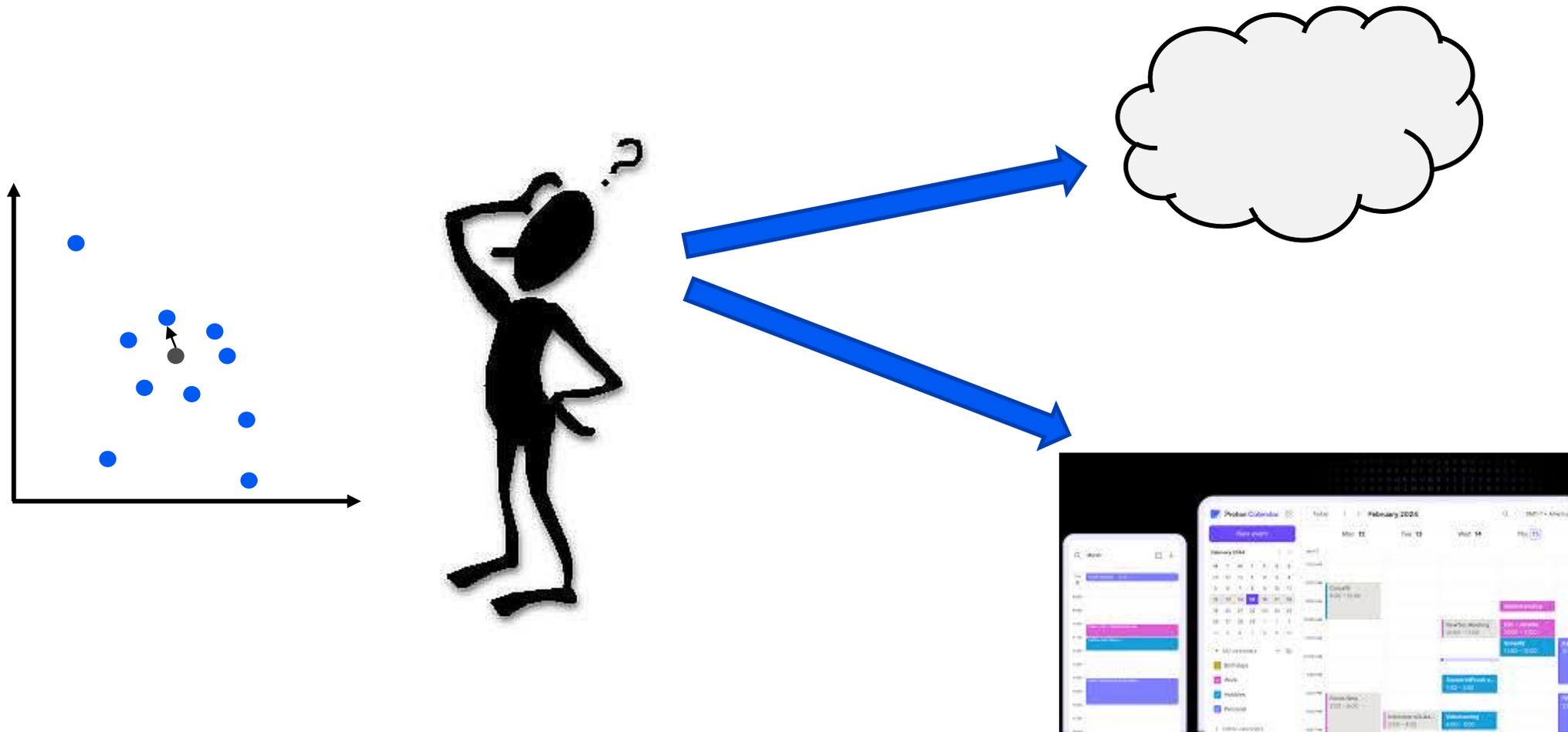


Picture source: https://www.linkedin.com/posts/barbara-geyer_ist-prompt-engineering-die-ki-f%C3%A4higkeit-activity-7348614645550669827-phEa/

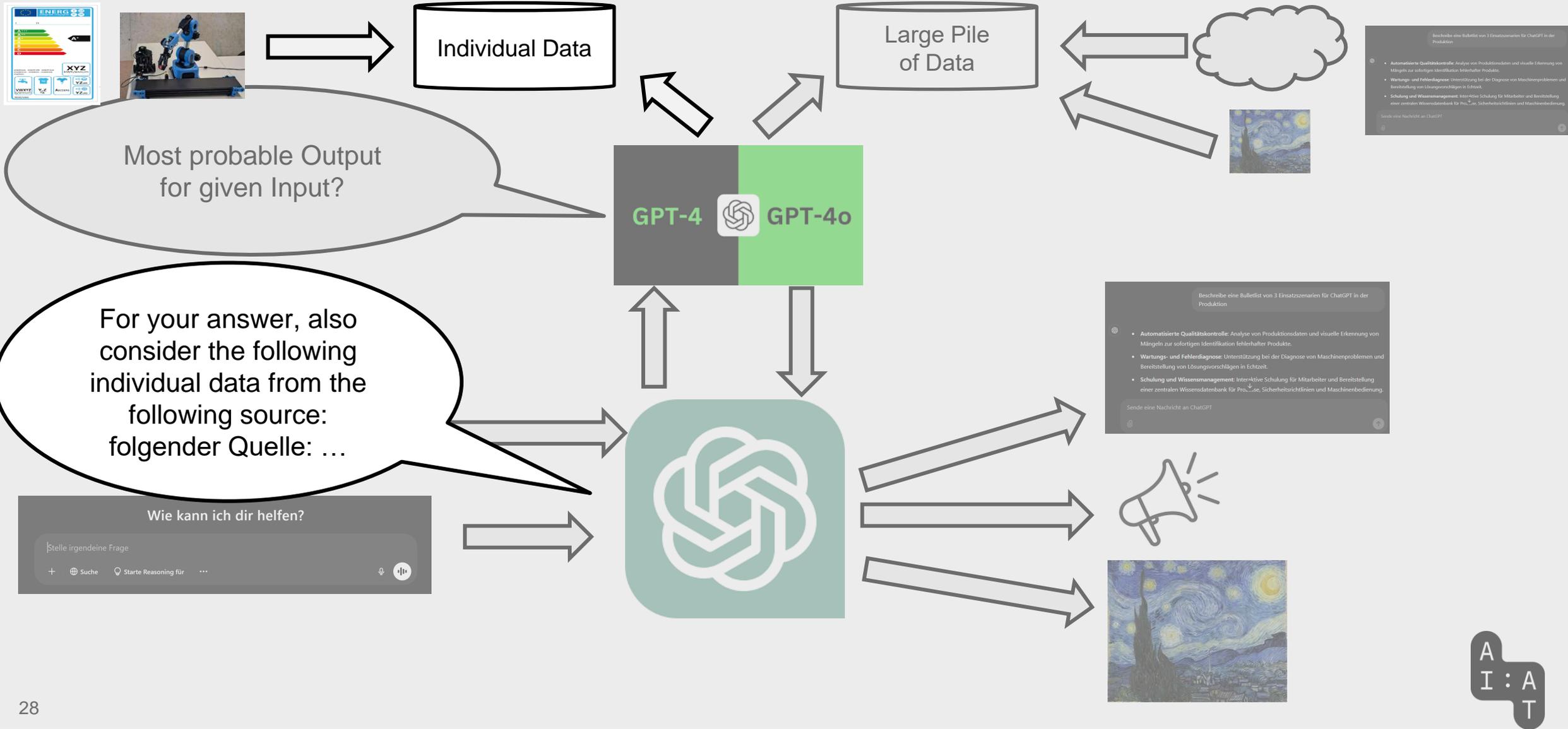
Extending LLM Knowledge



Extending LLM Knowledge

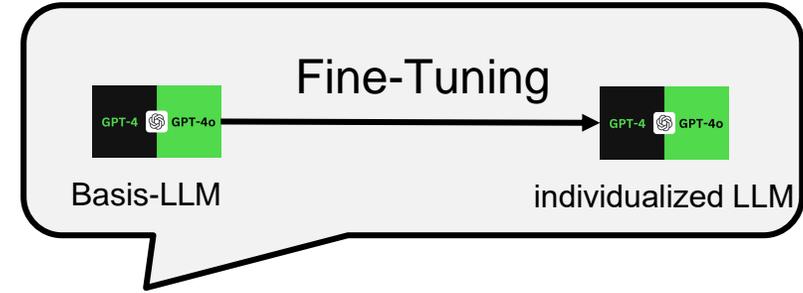


Retrieval Augmented Generation (RAG)



RAG VS Fine-Tuning

RAG
 Prompt: my question
 Context: my data



Provisioning of data

On each prompt

Once during training

Costs

Small overhead during operations
 Dedicated software architecture necessary

Special hardware necessary

Data

Dynamic

Static

Required Know-How

Software architecture, data retrieval,
 working with APIs

NLP, Deep Learning, model configuration,
 data preparation, ML evaluation

<https://www.redhat.com/de/topics/ai/rag-vs-fine-tuning>



RAG VS Fine-Tuning

RAG
 Prompt: my question
 Context: my data



Provisioning of data

On each prompt

Once during training

Example: Sales Chatbot

Example: Legal Texts

Costs

Data

Dynamic

Static

Required Know-How

Software architecture, data retrieval, working with APIs

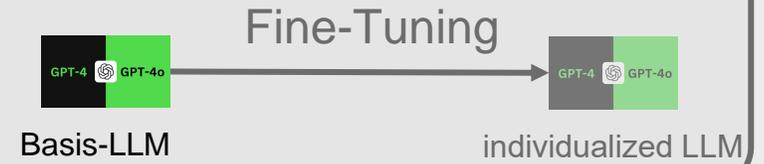
NLP, Deep Learning, model configuration, data preparation, ML evaluation

<https://www.redhat.com/de/topics/ai/rag-vs-fine-tuning>



RAG VS Fine-Tuning

RAG
 Prompt: my question
 Context: my data



Provisioning of data

On each prompt

Once during training

Costs



Course:

<https://ai-at.eu/training/foundations-of-llm-mastery-retrieval-augmented-generation-2/>

Special hardware necessary

Data

Dynamic

Static

Required Know-How

Software architecture, data retrieval working with APIs

data preparation, LLM evaluation



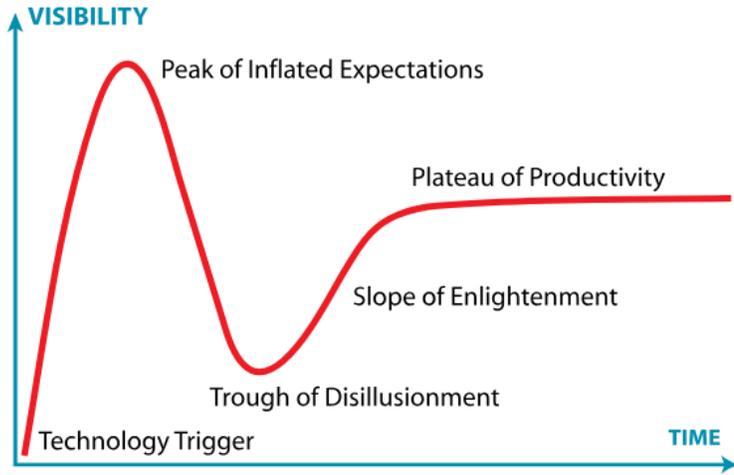
Course:

<https://ai-at.eu/training/foundations-of-llm-mastery-fine-tuning-on-one-gpu-2/>

<https://www.redhat.com/de/topics/ai/rag-vs-fine-tuning>



Summary



```

from openai import OpenAI
client = OpenAI(api_key="...")
response = client.chat.completions.create(
    model="gpt-4o-mini",
    max_tokens=100,
    messages=[{"role": "user", "content": "YOUR PROMT HERE"}]
)
print(response.choices[0].message.content)
    
```

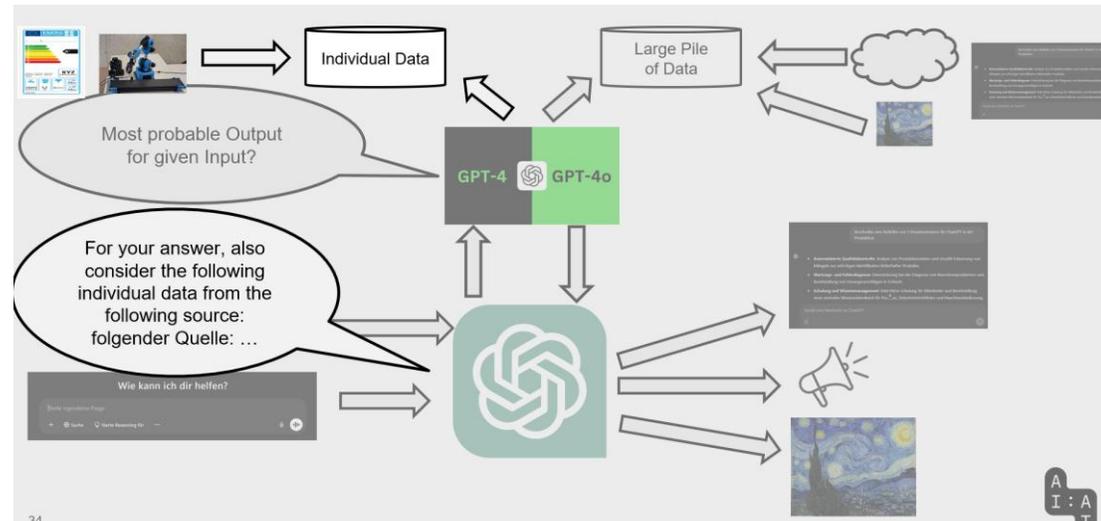
Was sind die Top 3 Use Cases für KI in der Produktion?

Die Top 3 Use Cases für KI in der Produktion sind:

- 1. Vorausschauende Wartung (Predictive Maintenance):** KI kann Sensordaten von Maschinen und Anlagen analysieren, um frühzeitig Anzeichen von Abnutzung oder Fehlfunktionen zu erkennen. So lassen sich Ausfälle vorhersagen und Wartungsmaßnahmen rechtzeitig einleiten, was die Betriebszeiten erhöht und unerwartete Ausfälle verringert.
- 2. Qualitätskontrolle und Fehlererkennung:** Durch den Einsatz von Computer Vision und maschinellem Lernen können Produktionsprozesse in Echtzeit überwacht und automatisch Fehler erkannt werden. Dies führt zu einer besseren Produktqualität, reduziert Ausschuss und minimiert die Notwendigkeit für manuelle Inspektionen.
- 3. Optimierung der Produktionsplanung und -steuerung:** KI-basierte Algorithmen können Produktionsabläufe optimieren, indem sie verschiedene Faktoren wie Materialverfügbarkeit, Produktionskapazitäten und Nachfrageschwankungen berücksichtigen. So lässt sich die Effizienz maximieren, die Lieferzeiten verkürzen und die Kosten senken.

Diese Anwendungen tragen dazu bei, die Effizienz, Qualität und Flexibilität in der Produktion erheblich zu steigern.

<https://campus.datacamp.com/courses/working-with-the-openai-api/introduction-to-the-openai-api>



Contact

Dr. Daniel Lehner

Expert AI Knowledge Transfer
AI Factory Austria AI:AT

daniel.lehner@ai-at.eu

AI Factory Austria AI:AT
Schwarzenbergplatz 2
1010 Wien, Austria

training@ai-at.eu
info@ai-at.eu

ai-at.eu

 [@ai-factory-austria](https://www.linkedin.com/company/ai-factory-austria)



Funded by



EuroHPC
Joint Undertaking



**Funded by
the European Union**

 **Federal Ministry
Innovation, Mobility
and Infrastructure
Republic of Austria**

under discussion with



AI Factory Austria AI:AT has received funding from the European High-Performance Computing Joint Undertaking (JU) under grant agreement No 101253078. The JU receives support from the Horizon Europe Programm of the European Union and Austria (BMIMI / FFG).