

## KI-Systeme für das Wissensmanagement aufbauen: CustomGPTs, RAG und Infrastrukturentscheidungen

ITSM Summit, Referent: Alexander Poschke Hamburg, 24.09.2025





# Anwendungsfälle Wissensmanagement





# Wissensmanagement Was sind Large Language Modelle?

LLM sind Lernmodelle, die auf riesigen Datenmengen trainiert werden, um Sprache zu verstehen und zu generieren.

#### Grenzen:

- LLM haben eine Wissensgrenze (ab dem Zeitpunkt des Trainings) können keine Echtzeitdaten abgerufen werden.
- LLM verstehen nicht wirklich Sprache, sondern arbeiten auf Basis von Mustererkennung und Wahrscheinlichkeiten.
- Die Modelle k\u00f6nnen Verzerrungen aus den Trainingsdaten \u00fcbernehmen.
- Es werden teilweise fehlerhafte oder widersprüchliche
   Informationen generiert, da keine Logiken hinterlegt sind.



Quelle: Pexels, https://www.pexels.com/



# CustomGPT Was ist das?

- Ein CustomGPT ist ein speziell auf die Bedürfnisse eines Unternehmens zugeschnittene Version eines GPT-Modells.
- Häufig spricht man auch von "Finetuning eines Sprachmodells"
- Es wird auf spezifische Datenquellen oder branchenspezifische Inhalte trainiert, um spezialisiertes Wissen und maßgeschneiderte Antworten zu liefern.

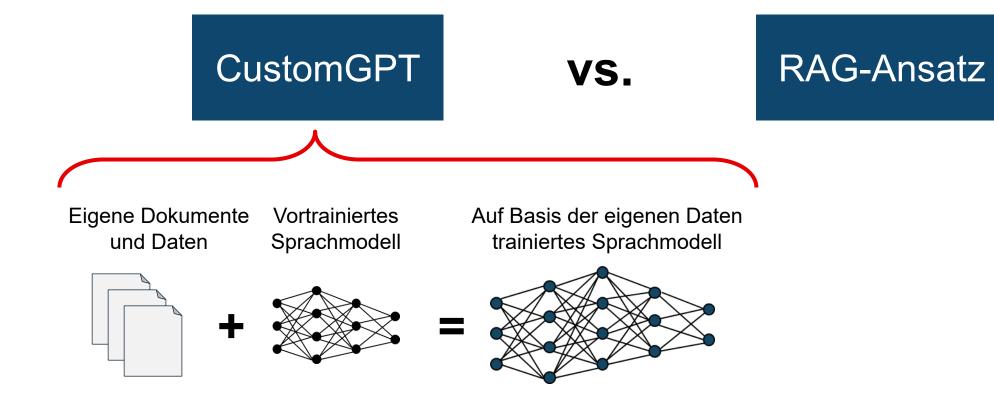




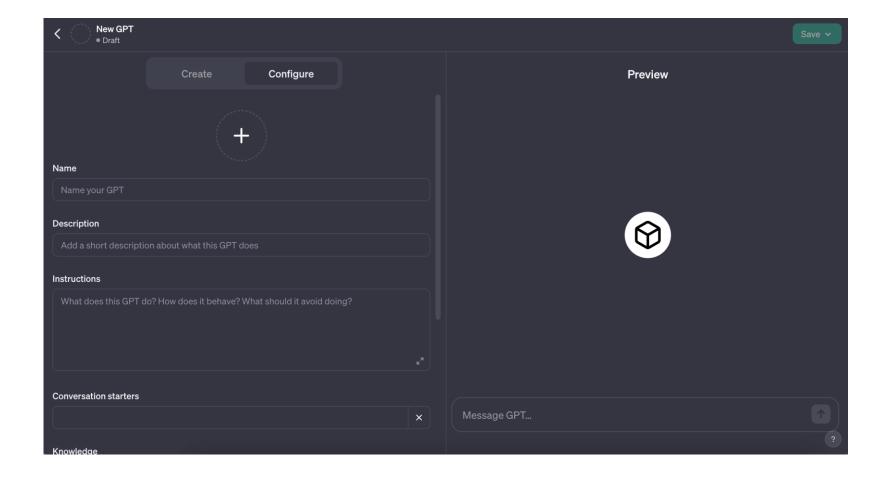
VS.

RAG-Ansatz







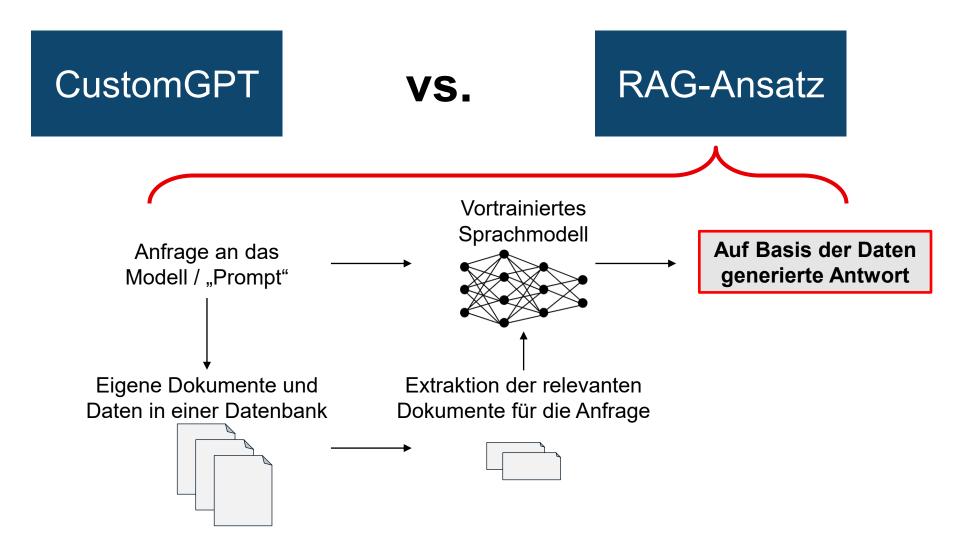




# RAG-Ansatz Was ist das?

- RAG steht für Retrieval Augmented Generation und kombiniert die Fähigkeit eines Modells, relevante Informationen aus einer externen Wissensbasis abzurufen (Retrieval), mit der Fähigkeit eines generativen Modells (wie GPT), auf Grundlage dieser Informationen neue Texte zu erstellen.
- Durch den Zugriff auf eine externe Wissensquelle kann RAG spezifischere und genauere Antworten liefern, besonders bei Anfragen, die detailliertes oder aktualisiertes Wissen erfordern.
- RAG ist besonders n\u00fctzlich in Bereichen, die hochgradig spezialisiertes Wissen erfordern, indem es das generative
   Modell mit den relevanten Informationen unterst\u00fctzt, die nicht im Modell selbst vortrainiert sind.
- Der RAG-Ansatz ermöglicht es, das generative Modell dynamisch mit aktuellen oder spezifischen Daten zu versorgen, wodurch es anpassungsfähiger und vielseitiger wird, ohne dass das Modell selbst umfangreich nachtrainiert werden muss.







### Wichtige Entscheidungen

## **CustomGPT**

Es ist besser kontrollierbar, was in das Modell fließt.

Wenn Daten hinzukommen, muss das Modell komplett neutrainiert werden.

Keine Nachverfolgung, woher die Informationen herkommen

## **RAG-Ansatz**

Es besteht hier die Gefahr, dass bei Übernahme kompletter bestehender Datenbanken viele schlechte Daten in der Datenbank landen.

Wenn Daten hinzukommen, muss das Modell nicht neutrainiert werden. Die Datenbank wird um die neuen Daten ergänzt.

Gute Nachverfolgung durch den Zugriff auf entsprechende Dokumente.



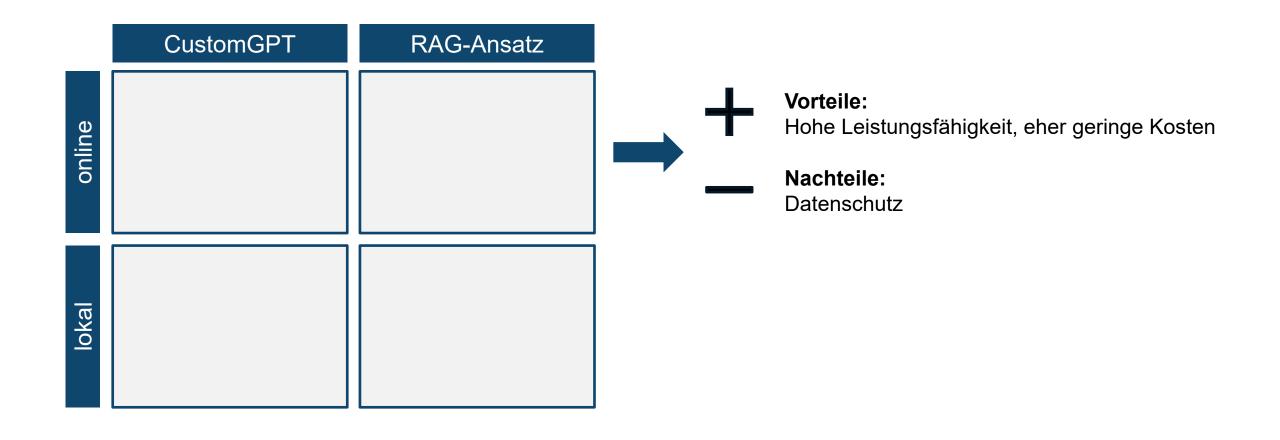
## Wichtige Entscheidungen

CustomGPT VS. RAG-Ansatz

lokal VS. online

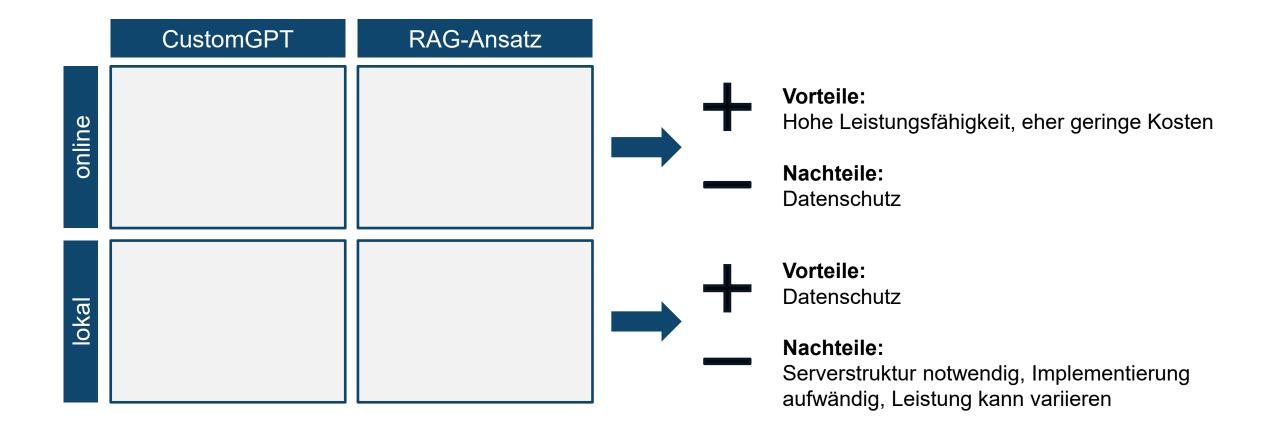


# Entscheidungsmatrix





## Entscheidungsmatrix





## Exkurs Gesetzeslage

DSGVO - Datenschutzgrundverordnung

TDDDG - Telekommunikations-Telemetrie-Datenschutz-Gesetz

KI-VO

**UrhG - Urheberrecht** 

AGG – Allgemeines Gleichbehandlungsgesetz



#### KI und Datenschutz

- KI-Systeme nutzen eine große Menge an sensiblen
   Informationen und Daten, um Muster zu erkennen und Prognosen zu treffen.
- Trainingsdaten können personenbezogene Daten beinhalten.
- Regeln der Datenschutzgrundverordnung (DSGVO) anwendbar

Vorsicht: Werden die Regeln des DSGVO nicht eingehalten, drohen **Geldstrafen von bis zu 10 Millionen Euro** oder 4 % des weltweiten Jahresumsatzes DSGVO -Datenschutzgrundverordnung TDDDG – Telekommunikations-Telemetrie-Datenschutz-Gesetz

KI-VO

UrhG -Urheberrecht

AGG – Allgemeines Gleichbehandlungsgesetz



#### TDDDG - Telekommunikations-Telemetrie-Datenschutz-Gesetz

- Regelung von Datenschutz bei digitalen Diensten und **Telekommunikation**
- setzt die ePrivacy-Richtlinie um
- Schutz von Grundrechten und Privatsphäre

Beispiel: Ein KI-Chatbot erhebt Nutzungsdaten zur Echtzeit-Personalisierung.

DSGVO -Datenschutzgrundverordnung

TDDDG -Telekommunikations-Telemetrie-Datenschutz-Gesetz

KI-VO

UrhG -Urheberrecht

AGG – Allgemeines Gleichbehandlungsgesetz



## Was sind personenbezogene Daten?

"Im Sinne dieser Verordnung bezeichnet der Ausdruck:

- "personenbezogene Daten" alle Informationen, die sich auf eine identifizierte oder identifizierbare natürliche **Person** (im Folgenden "betroffene Person") beziehen; als identifizierbar wird eine natürliche Person angesehen, die direkt oder indirekt, insbesondere mittels **Zuordnung zu einer Kennung** wie einem **Namen**, zu einer **Kennnummer**, zu **Standortdaten**, zu einer **Online-Kennung** oder zu einem oder mehreren besonderen **Merkmalen**, die Ausdruck der physischen, physiologischen, genetischen, psychischen, wirtschaftlichen, kulturellen oder sozialen **Identität** dieser natürlichen Person sind, identifiziert werden kann:
- [...]"

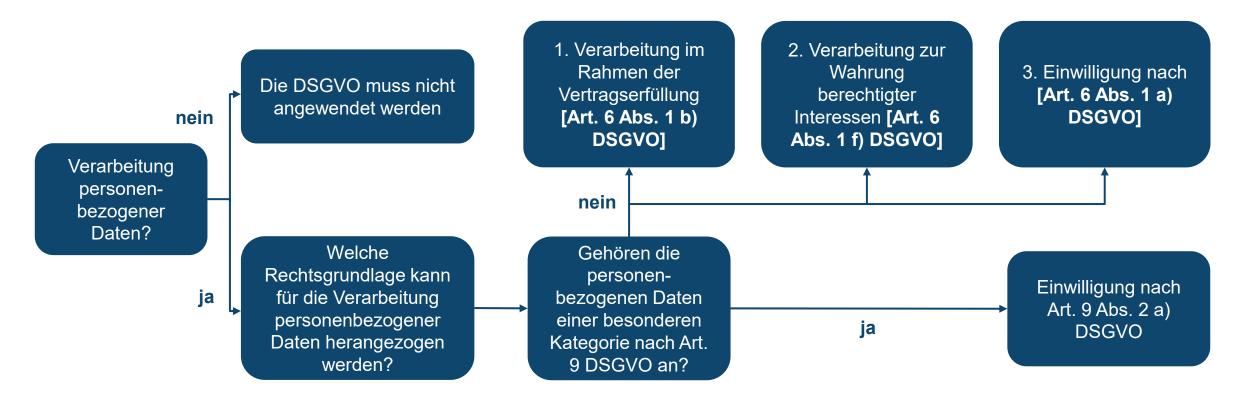


## Verarbeitung besonderer Kategorien personenbezogener Daten (Art. 9 DSGVO)

(1) Die Verarbeitung personenbezogener Daten, aus denen die rassische und ethnische Herkunft, politische Meinungen, religiöse oder weltanschauliche Überzeugungen oder die Gewerkschaftszugehörigkeit hervorgehen, sowie die Verarbeitung von genetischen Daten, biometrischen Daten zur eindeutigen Identifizierung einer natürlichen Person, Gesundheitsdaten oder Daten zum Sexualleben oder der sexuellen Orientierung einer natürlichen Person ist untersagt.

# *IPH*

#### **DSGVO**



Quelle: <a href="https://digitalzentrum-chemnitz.de/wissen/ki-anwendungen-und-dsgvo/">https://digitalzentrum-chemnitz.de/wissen/ki-anwendungen-und-dsgvo/</a>



## Grundsätze für die Verarbeitung personenbezogener Daten (Art. 5 DSGVO)

### Rechtmäßigkeit, Verarbeitung nach Treu und Glauben, Transparenz

Daten müssen rechtmäßig, fair und transparent für die betroffene Person verarbeitet werden

#### Zweckbindung

• Daten nur für spezifische, legitime und ausdrücklich angegeben Zwecke erheben und verwenden

#### **Datenminimierung**

Nur absolut erforderliche Daten erheben

#### Richtigkeit

Daten korrekt und aktuell

#### Speicherbegrenzung

Aufbewahrung von Daten nicht länger als erforderlich

#### Integrität und Vertraulichkeit

 Personenbezogene Daten müssen sicher verarbeitet werden, um unbefugten oder unrechtmäßigen Zugriff, Verlust oder Beschädigung der Daten zu verhindern



## Datenschutzrechtliche Herausforderungen im Unternehmen

- vor der Nutzung von personenbezogenen Daten: Datenschutz-Folgeabschätzung (Art. 35 DSGVO) anfertigen und datenschutzbeauftragte Person kontaktieren
- schriftliche Einwilligung vor der Nutzung von personenbezogenen Daten einholen

Beispiel: Ein KI-Tool wertet Metadaten aus E-Mails, Kalendern und Microsoft Teams aus, um Stressindikatoren zu prognostizieren.



## Datenschutz online vs. lokal

	Lokale Nutzung	Online-Nutzung
Datenfluss	Eingaben und Ausgaben bleiben auf Gerät oder internen Server	<ul> <li>Eingaben werden an Server des KI-Anbieters geschickt, dort verarbeitet und Antwort zurückgesendet</li> </ul>
Risiken	<ul> <li>keine Übertragung an Dritte → gering</li> <li>Verantwortung für Datensicherheit verbleibt im Unternehmen (z. B. Verschlüsselung, Zugriffsrechte)</li> </ul>	<ul> <li>Daten verlassen das System</li> <li>Risiko von Zugriffen durch Dritte oder unklaren Speicher-/Verarbeitungsorten</li> </ul>
DSGVO-Relevanz	<ul> <li>Kein externer Auftragsverarbeiter im Spiel →         es findet in der Regel keine "Übermittlung an         Dritte" statt</li> <li>Unternehmen allein Verantwortlicher im Sinne         der DSGVO</li> <li>Wenn sensible oder personenbezogene Daten         verarbeitet werden, muss das Unternehmen         selbst technische und organisatorische         Maßnahmen nach Art. 32 DSGVO         sicherstellen</li> </ul>	<ul> <li>Auftragsverarbeitungsvertrag (AVV) nach Art. 28 DSGVO erforderlich</li> <li>Zusätzliche Anforderungen, falls Daten außerhalb der EU verarbeitet werden</li> <li>Transparenzpflichten: Du musst deine Nutzer/Kunden darüber informieren, dass und wie ihre Daten an einen KI-Dienst übermittelt werden.</li> </ul>



## Buchempfehlung

#### Rechtsleitfaden KI im Unternehmen

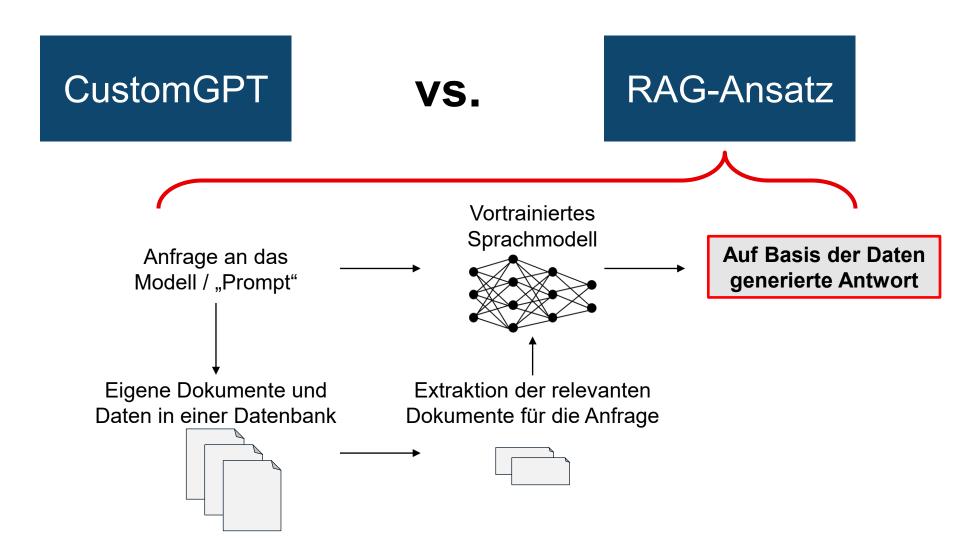
von Niklas Mühleis, Nick Akinci

ISBN: 978-3-367-10098-9

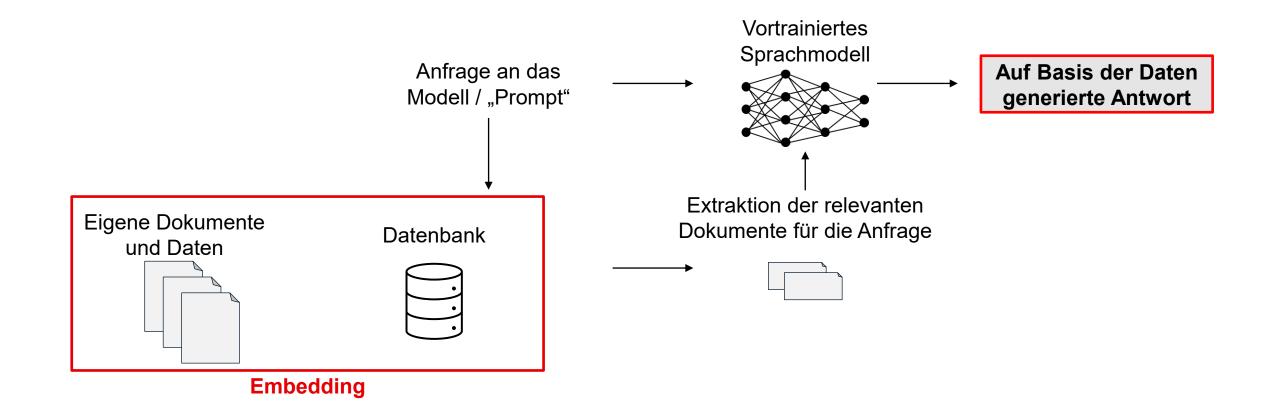
https://www.rheinwerk-verlag.de/rechtsleitfaden-ki-im-unternehmen/?srsltid=AfmBOoptxmu5XE3t9tiAeOFjGE-vgzKKQKtVM8vJfD-fh0x-7H1iK-mE



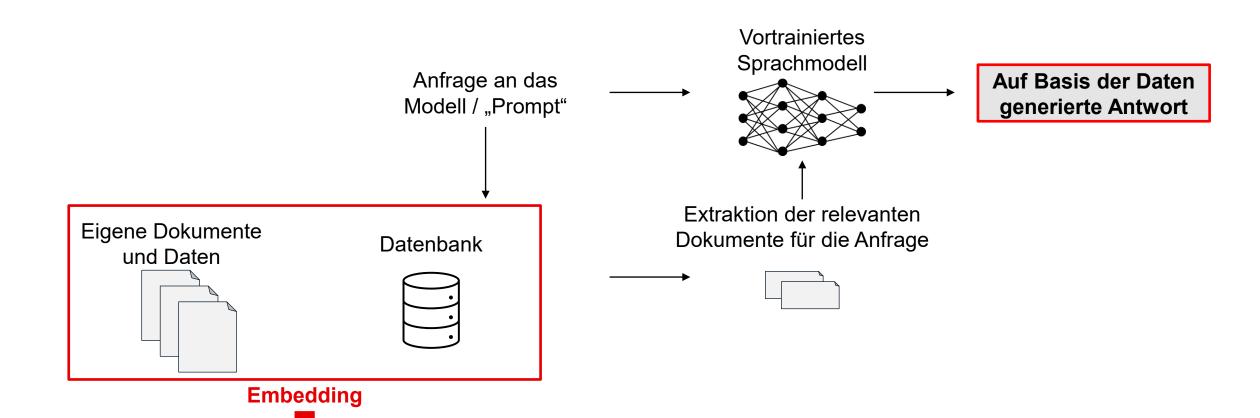












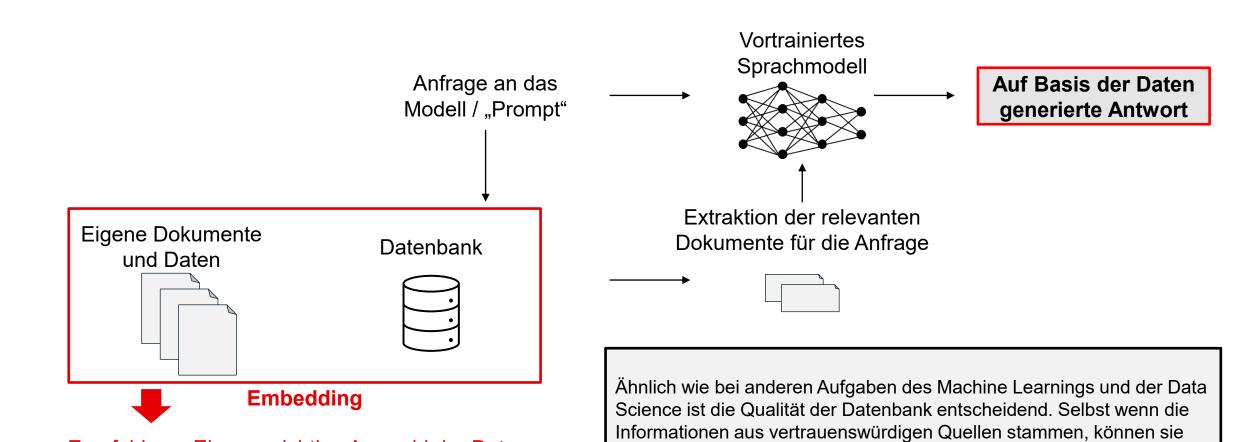
© IPH | Paulina Merkel, Alexander Poschke

Auswahl eines geeigneten Embedding-Algorithmus (z. B. Mistral, Sauerkraut NeMo, Jina)



Empfehlung: Eine vorsichtige Auswahl der Daten

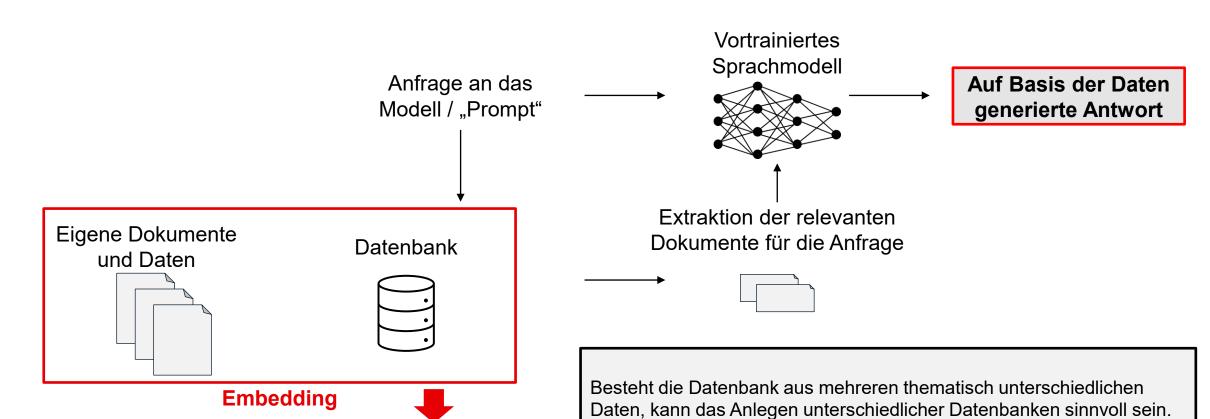
vornehmen, nicht einfach alle "reinschmeißen".



© IPH | Paulina Merkel, Alexander Poschke | Produktion erforschen und entwickeln | 27

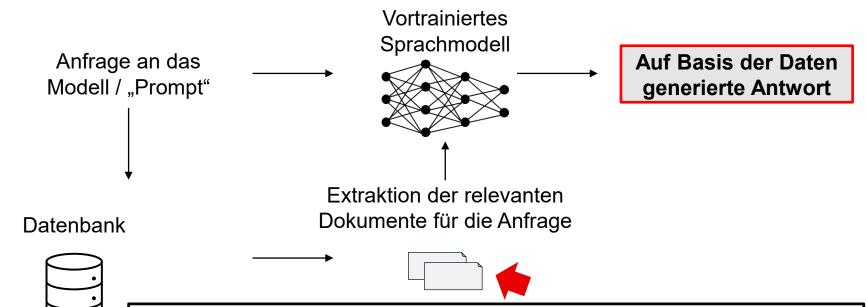
irrelevanten und irreführenden Text enthalten.



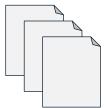


Empfehlung: Entwicklung einer Vektordatenbank (z. B. Milvus, FAISS, Annoy)

Es könnte auch eine bei der Vorverarbeitung eine Zusammenfassung von Dokumenten erstellt und gespeichert werden.

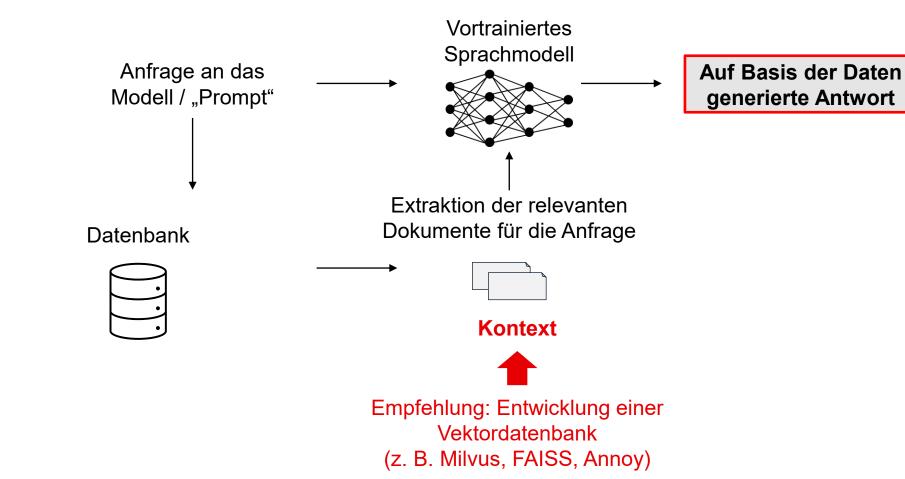


**Eigene Dokumente** und Daten



Abhängig von der Art der Datenbank und der Größe des LLM kann der Prozess der Datenübergabe langsam sein. Bevor alle gefundenen Daten an das LLM weitergegeben werden, könnte das LLM zur Verbesserung der Geschwindigkeit z. B. eine Zusammenfassung der Daten vornehmen oder die gefundenen Daten können neu geordnet werden. Eine Möglichkeit wäre auch ein Metadatenfilter (z. B. eine Priorisierung).



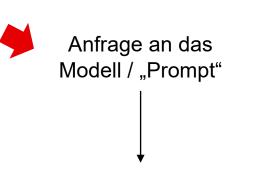


Eigene Dokumente

und Daten

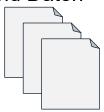


Empfehlung: Vorverarbeitung der Anfrage, indem sie in thematisch getrennte Teilanfragen zerlegt wird



Auf Basis der Daten generierte Antwort

Eigene Dokumente und Daten



Datenbank

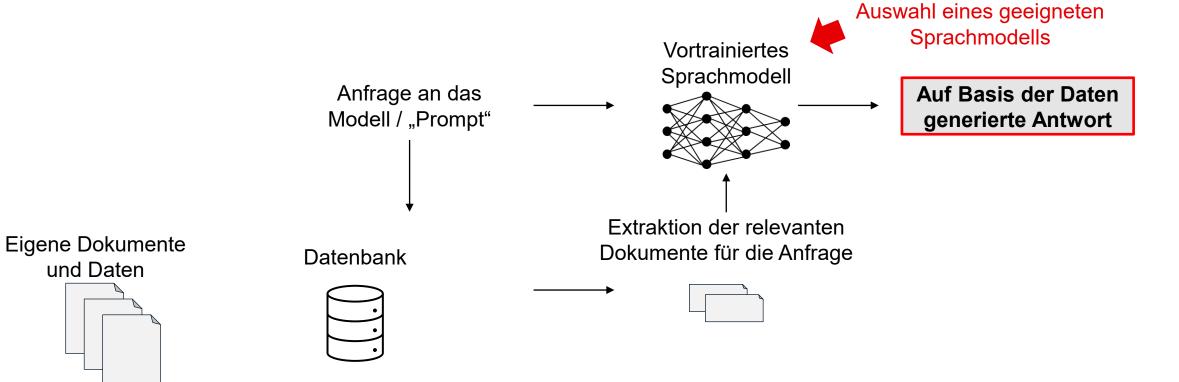


Extraktion der relevanten Dokumente für die Anfrage

Vortrainiertes

Sprachmodell





© IPH | Paulina Merkel, Alexander Poschke

und Daten



## Wichtige Entscheidungen



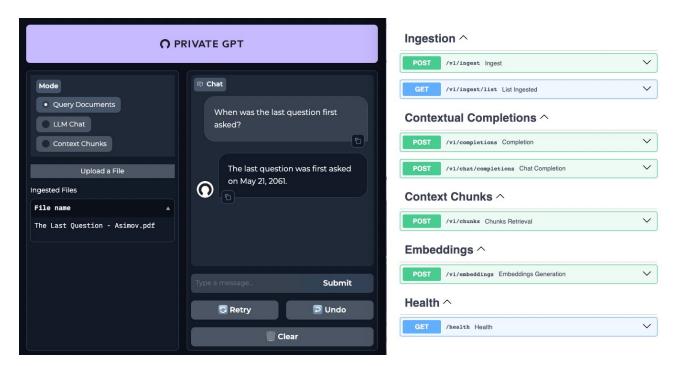
VS.

RAG-Ansatz



## Lokale Systeme aufbauen

- Programmvorschläge:
  - LM-Studio
  - PrivateGPT (Wechsel auf die kommerzielle Software Zylon.ai möglich)
  - pdfGPT
  - GPT4AII
  - Axolotl
  - AnythingLLM





# **RAG** Einkaufen fertiger Lösungen

Berücksichtigung der DSGVO besonders einfach, wenn RAG-Anbieter deutsch ist.

Deutsche Anbieter für RAG-Systeme:

https://www.it-p.de/leistungen/wisbee/

https://www.kamium.de

https://www.kencube.com/de/loesungen/produktion.htm

→ Es gibt noch viele weitere!





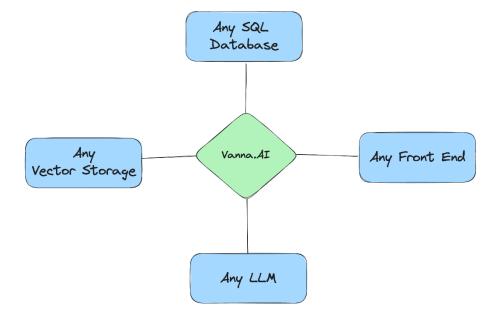




#### **RAG**

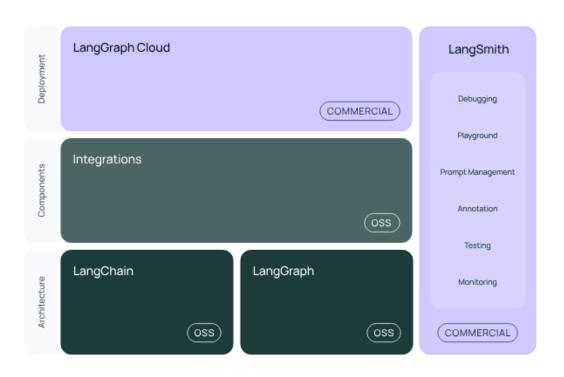
## Aufbauen eigener Lösungen → Vanna.ai und Langchain

Use AI to Interact With Your Database



#### Link zum Ausprobieren:

https://github.com/vanna-ai/vanna

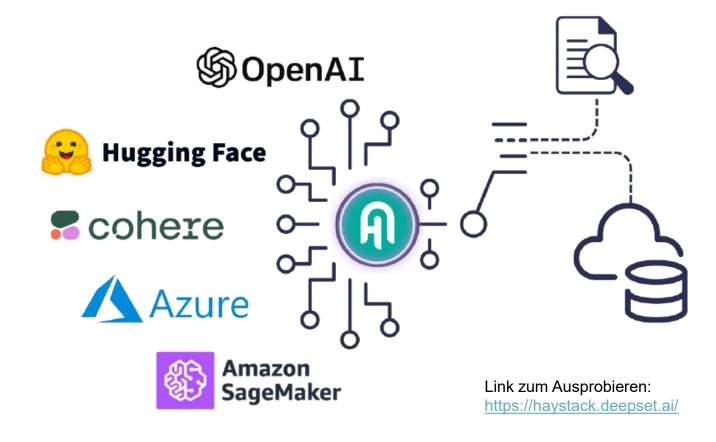


#### Link zum Ausprobieren:

https://js.langchain.com/docs/introduction/



# **RAG** Haystack





# Sprachmodelle was gibt es open source vs proprietär (Vor- und Nachteile)

Kriterium	Open Source	Proprietär
Leistung	Gut	Sehr gut
Kostenmodell	Einmalige Infrastrukturkosten	Variable API-Kosten
Anpassbarkeit	Vollständig (bei Lizenz erlaubt)	Eingeschränkt, meist API-basiert
Datenschutz/Kontrolle	Sehr hoch (on-premise möglich)	Geringer, abhängig vom Anbieter
Support	Community-getrieben, eigene Verantwortung	Professioneller Support, SLAs etc.
Lizenz	Offen (teilweise mit Einschränkungen)	Nutzungsbedingungen strikt



### Sprachmodelle aktuelle Benchmarks

#### Verschiedene Benchmarks (Auszug):

- Scale com SFAL Leaderboard
- LLM-Stats.com
- Vellum AI LLM Leaderboard
- LM Arena Leaderboard
- → Bestenlisten abhängig von Art der Tests
- → Viele Änderungen bei Anpassung der Modelle
- → Austauschbarkeit kann von Vorteil sein!

#### Hinweis:

Unter <a href="https://lmarena.ai/">https://lmarena.ai/</a> können Sie verschiedene Modelle gegeneinander antreten lassen und anschließend selber bewerten!





## KI-Training und das Urheberrecht



#### Der Hamburgische Beauftragte für Datenschutz und Informationsfreiheit

13.11.2023

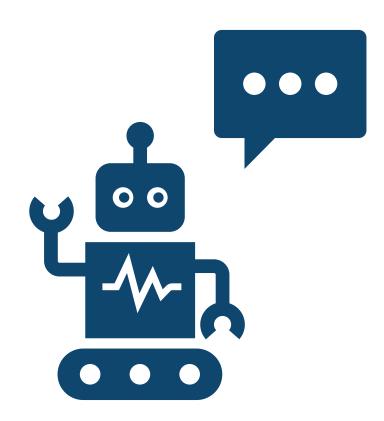
#### **Checkliste zum Einsatz LLM-basierter Chatbots**

Generative KI in Form von Chatbots bietet die Möglichkeit, schnell und unkompliziert Inhalte zu erstellen. Bekannte Large Language Models (LLM) sind ChatGPT, Luminous oder Bard. In vielen Einrichtungen sind die Tools mittlerweile Teil des Arbeitsalltags geworden, oft jedoch ohne verbindliche Vorgaben zur Nutzung. Dass die Sprachmodelle üblicherweise in einer Cloud betrieben werden, birgt verschiedene Datenschutzrisiken. Zum einen ist der Schutz vertraulicher Daten ge-



#### Blick in die Zukunft

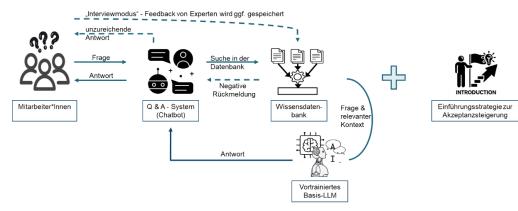
- Für Anwendungen, in denen die Qualität der Ausgabe sehr wichtig ist, können KI-basierte Agenten genutzt werden.
- Sog. Agentic RAG haben trainierte KI-Agenten, die die Aussagen auf Plausibilität und Richtigkeit prüfen.
- Anwendungsfälle wären z. B., wenn Kunden im Kontakt mit der KI stehen und die Ausgabe mit hoher Sicherheit richtig sein muss oder wenn es um sicherheitsrelevante Themen geht.





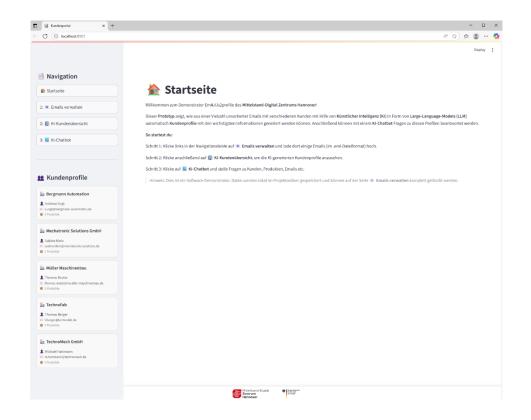
#### Zusammenarbeit mit dem IPH

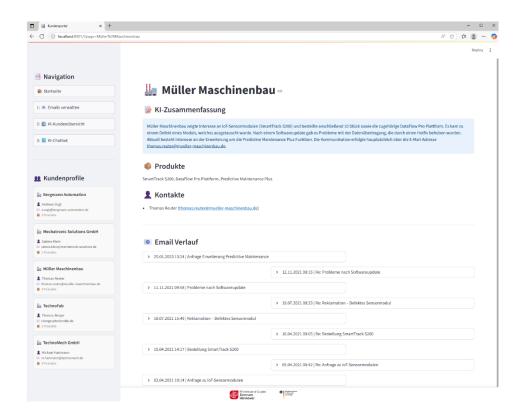
- Inanspruchnahme von innovativen Dienstleistungen rund um CustomGTPs und RAG-Modelle.
  - Mögliche Förderprogramme sind z. B. Innovationsgutschein, Digitalbonus, usw.
- Partizipation an Forschungsprojekten, die sich mit Sprachmodellen und der Wissensspeicherung beschäftigen, z. B.:
  - Entwicklung eines KI-basierten, unternehmensindividuellen Chatbots zur Aufnahme und Nutzung von explizitem und implizitem Expertenwissen im Einkauf (EnKi-Bot)
  - Befähigung von produzierenden KMUs zum selbstständigen und sicheren Umgang mit dem Al Act und weiteren KI-Regularien (Al Act-Ready)
  - Befähigung von KMU zur Implementierung nachhaltiger Automatisierungslösungen in der Produktion und Entwicklung einer Einführungsstrategie (AutoSus)
  - Und viele mehr: https://www.iph-hannover.de/de/forschung/forschungsprojekte/





## Praxisbeispiel MDZ/IPH







— <a href="https://github.com/iph-hannover/MDZ-Demonstrator">https://github.com/iph-hannover/MDZ-Demonstrator</a>



### Kontakt



M. Sc. Alexander Poschke

IPH – Institut für Integrierte Produktion Hannover gGmbH Hollerithallee 6 30419 Hannover



+49 (0)511 27976-229



poschke@iph-hannover.de



www.iph-hannover.de



M. Sc. Paulina Merkel

IPH – Institut für Integrierte Produktion Hannover gGmbH Hollerithallee 6 30419 Hannover



+49 (0)511 27976-331



merkel@iph-hannover.de



www.iph-hannover.de