

# Bot-Camp

Wissensbasierte KI-Assistenten  
an Hochschulen einsetzen



## Einführung

Univ.-Prof. Dr. Malte Persike

 [persike@cls.rwth-aachen.de](mailto:persike@cls.rwth-aachen.de)

 +49 (170) 10 57 54 1

## Zum Auftakt werden wir:



Verstehen, wann und warum wir Custom Chatbots brauchen.



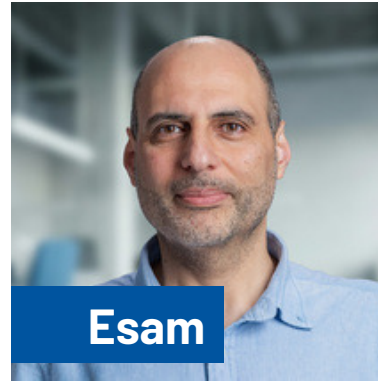
Die Komponenten eine Custom Chatbot kennenlernen.



Risiken für die Qualität der Antworten demonstrieren.

# KI für Heute

# Das Team





# Will Agentic AI Break Higher Education?

A new app offers to do students' work for them.  
Professors aren't ready.

ILLUSTRATION BY THE CHRONICLE; GETTY; WIKIMEDIA COMMONS

By *Jason Gulya*

March 3, 2026

**E**instein first came to my attention on February 25. It was promoted as an app that logged into a student's Canvas account and completed all the work on their behalf. As advertised on their website, "Einstein is an AI with a computer. He logs into Canvas every day, watches lectures, reads essays, writes papers, participates in discussions, and submits your homework — automatically."

Quelle: <https://www.chronicle.com/article/will-agentic-ai-break-higher-education>

## KI Agenten

auch Skills genannt

sind Custom Chatbots, die Schnittstellen mit anderen digitalen Systemen oder KI-Modellen haben.

Sie interagieren sowohl mit den Nutzenden als auch **mit anderen Softwarediensten** und können dort Aktionen auslösen.

## Custom Chatbots

auch CustomGPT, Gem oder Model genannt

sind KI-Modelle, die auf bestimmte Antworttypen optimiert sind.

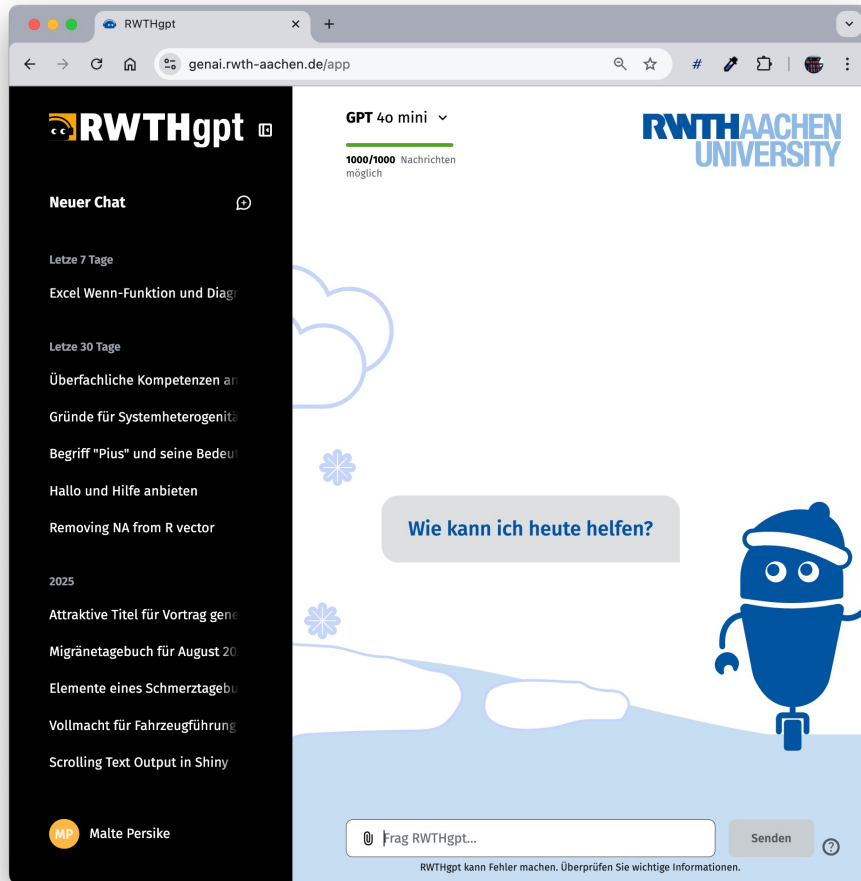
Sie interagieren nur mit den Nutzenden und können aus ihrem Chat-Interface **nicht „ausbrechen“**.

# 1

## Warum Custom Chatbots?

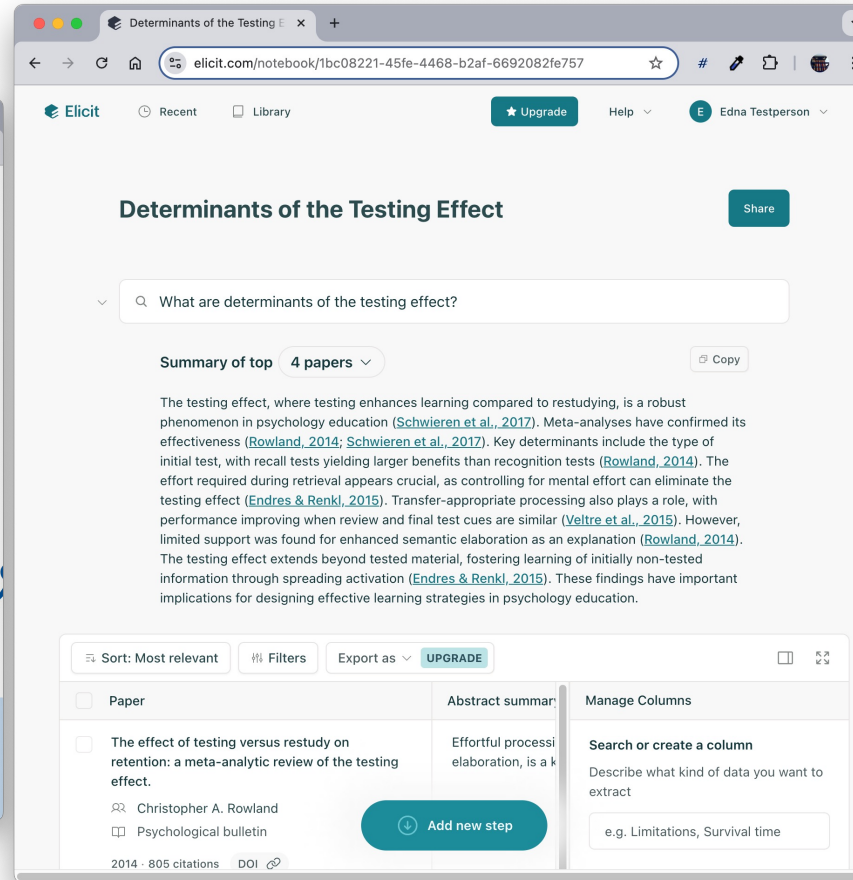
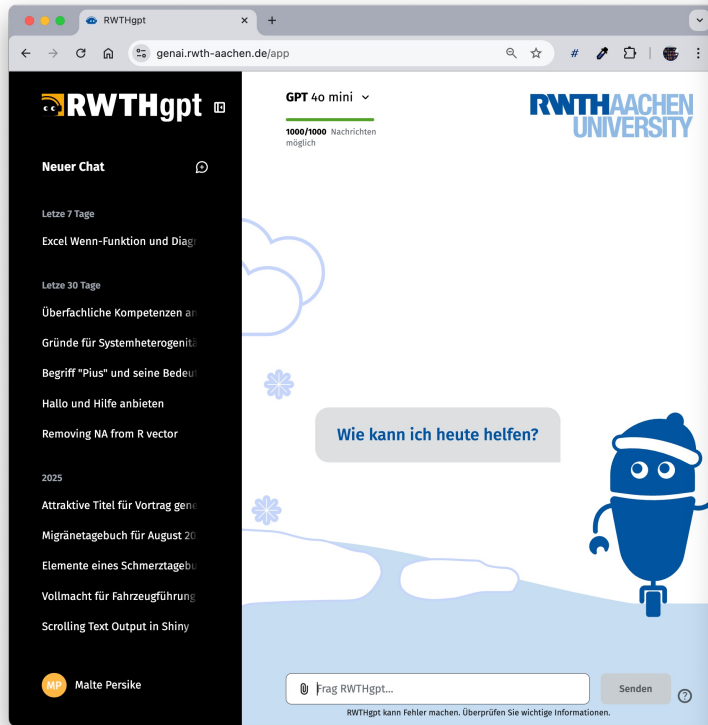


# Integrations Ebenen von KI-Systemen



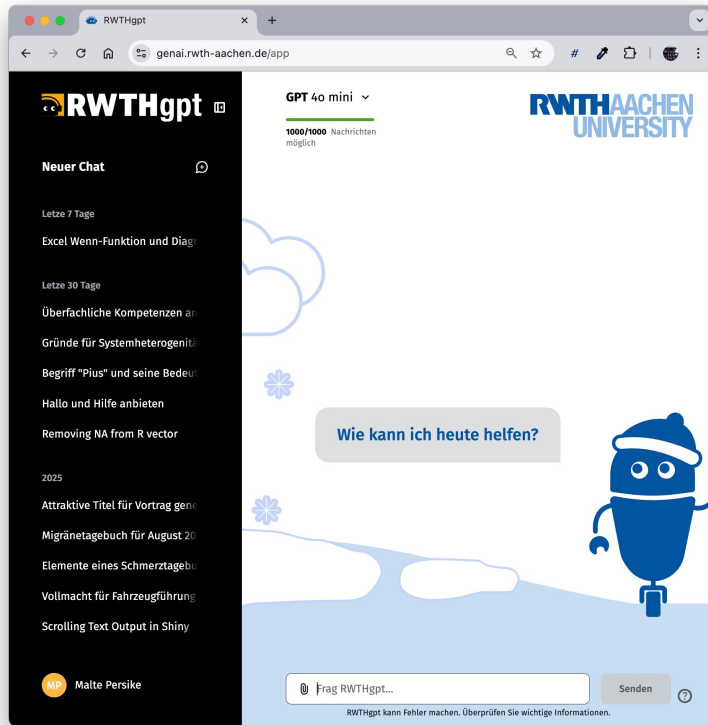
# Integrationsebenen von KI-Systemen

## General Purpose UI

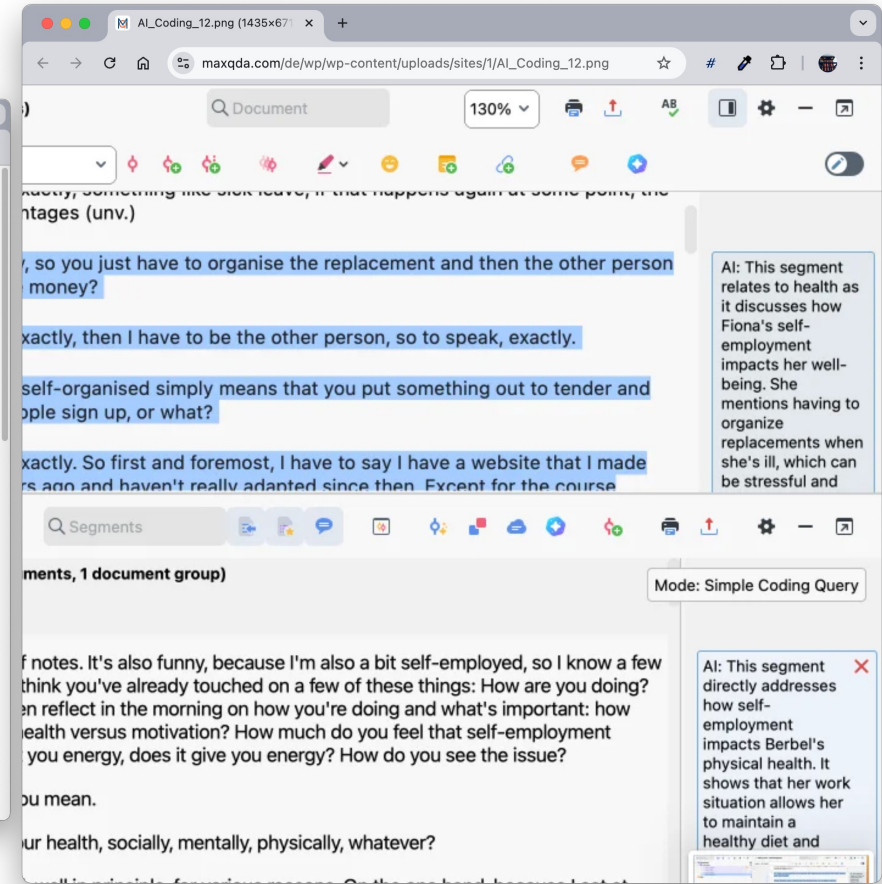
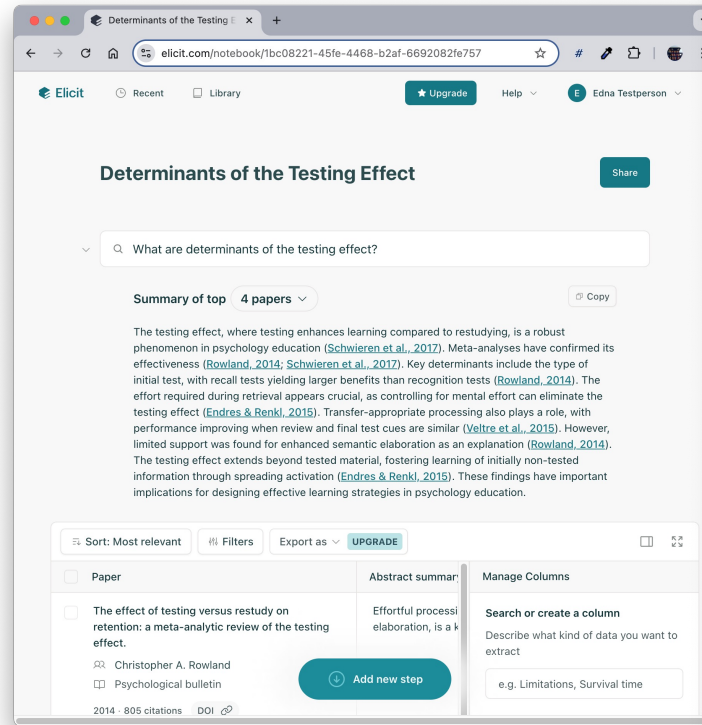


# Integrationsebenen von KI-Systemen

## General Purpose UI

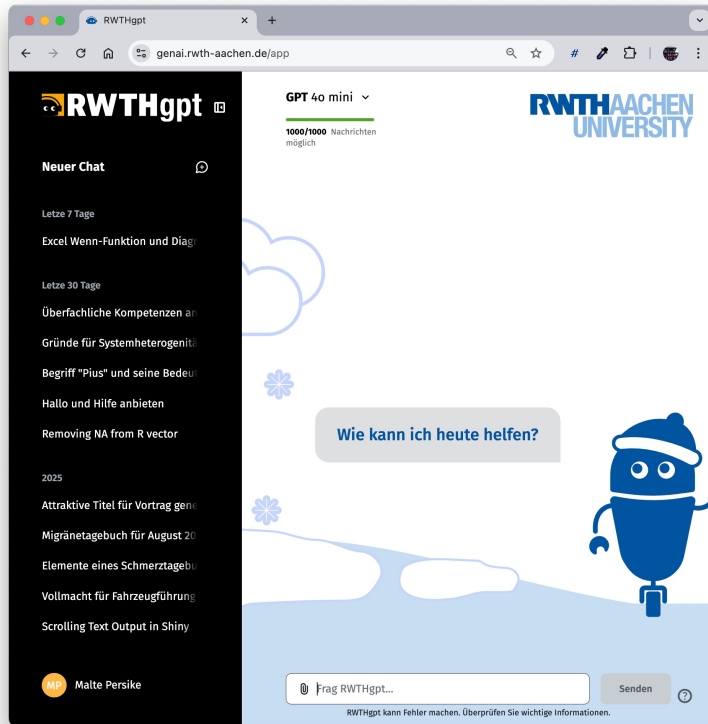


## Special Purpose UI / Custom Chatbot

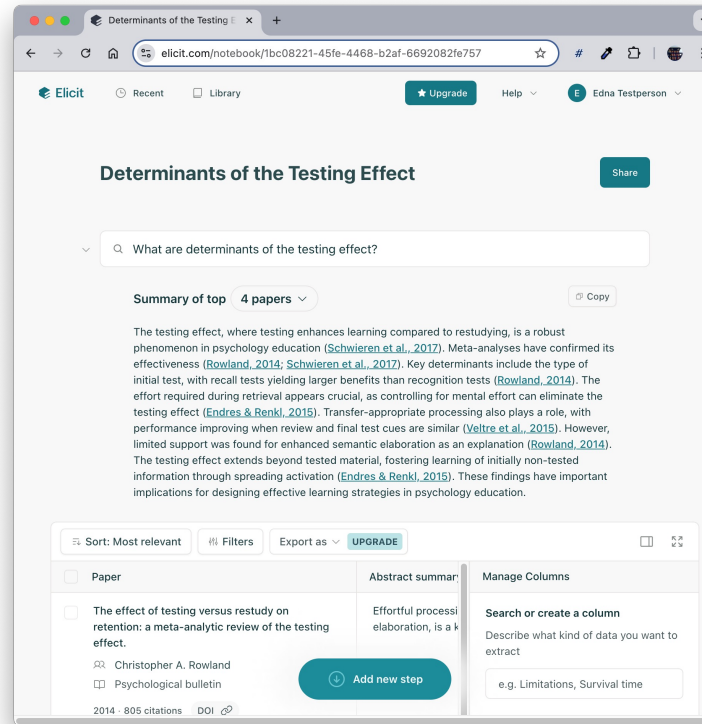


# Integrationsebenen von KI-Systemen

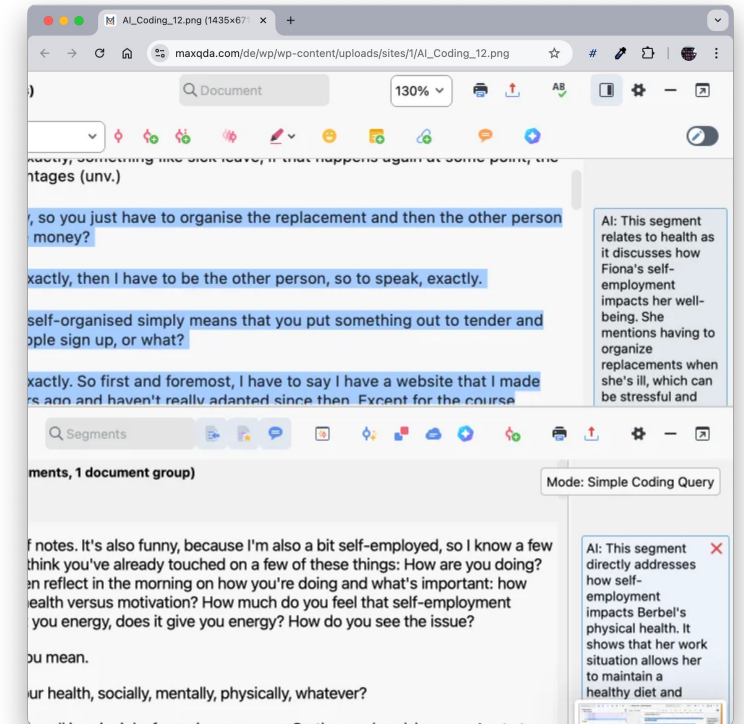
## General Purpose UI



## Special Purpose UI / Custom Chatbot



## In-App UI



Starke Kompetenzerfordernisse

Schwache Kompetenzerfordernisse

# KI-Systeme als Unterstützungswerkzeug für Alle

---

**Variante 1:** Selbständige Nutzung eines allgemeinen LLM wie z.B. ChatGPT



**Variante 2:** Nutzung eines allgemeinen LLM mit vorgefertigten Prompts und Dokumenten

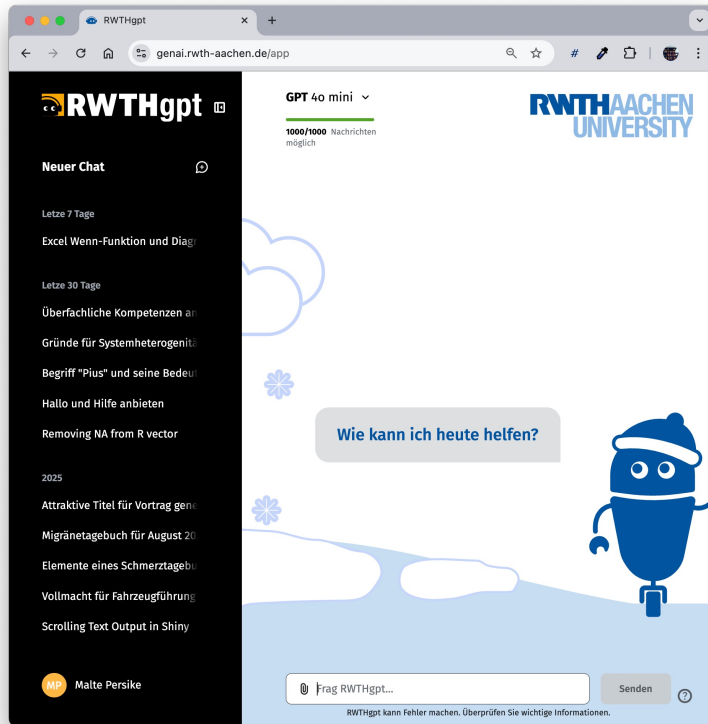


**Variante 3:** Nutzung von Custom Chatbots oder Spezialapplikationen mit eigens trainierten und instruierten LLMs

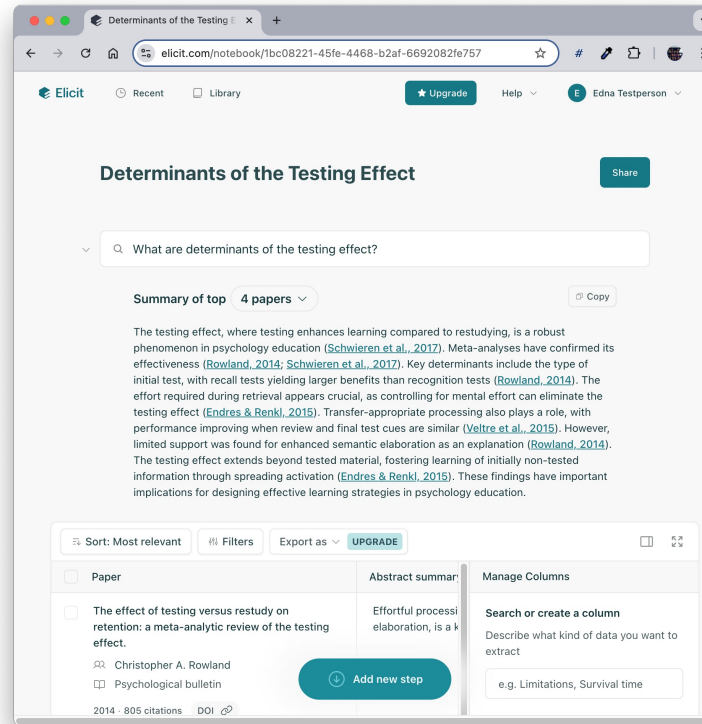


# Integrationsebenen von KI-Systemen

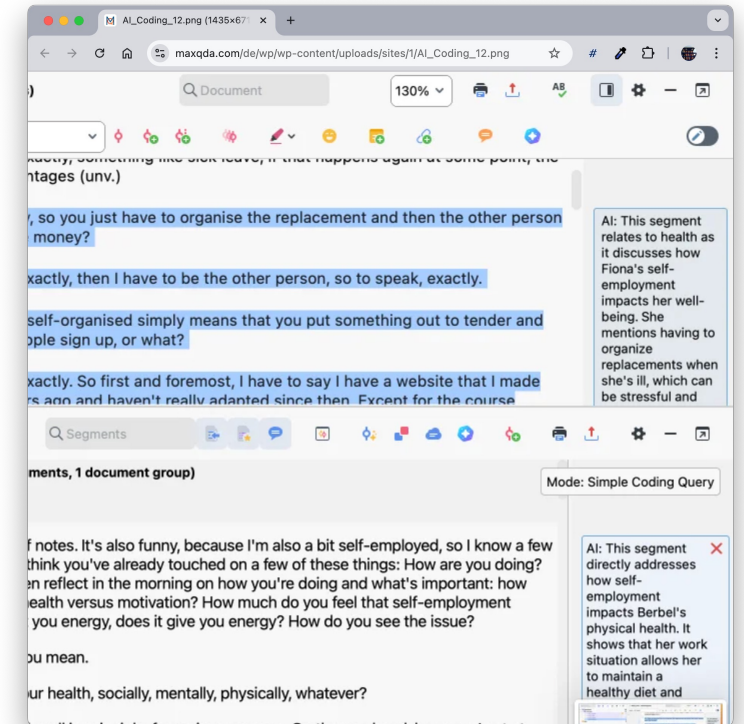
## General Purpose UI



## Special Purpose UI / Custom Chatbot



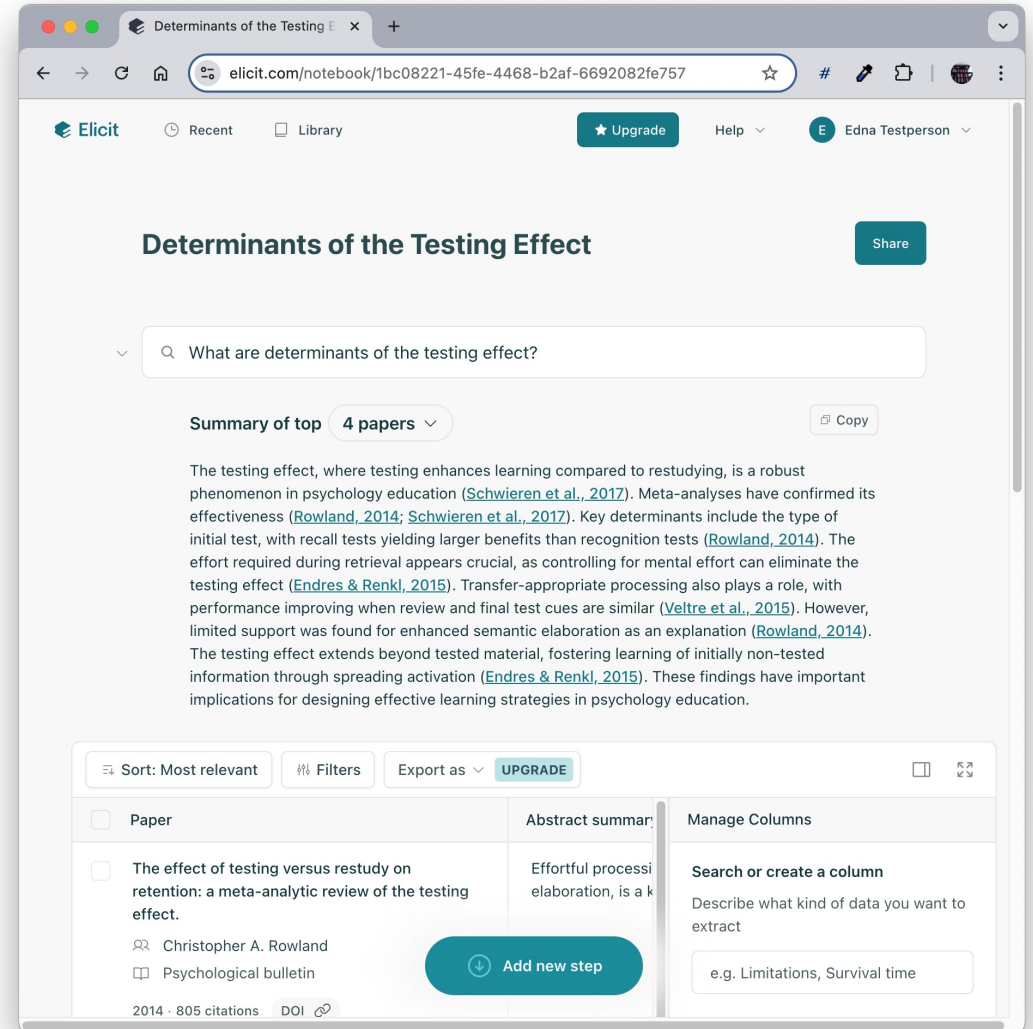
## In-App UI



# Spezialapplikation oder Custom Chatbot?

Gezeigt ist das Literaturrecherche-  
Werkzeug **Elicit** (<https://elicit.com>).

Im Kern ist auch eine solche Spezial-  
applikation kaum mehr als ein  
**Custom Chatbot** mit schicker  
Oberfläche.



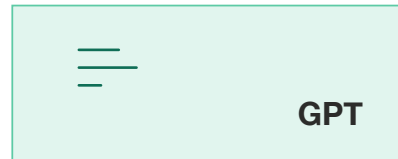
# 2

## **Komponenten eines Custom Chatbot**

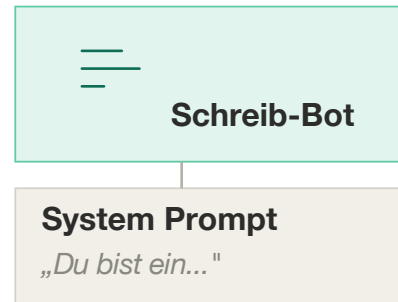


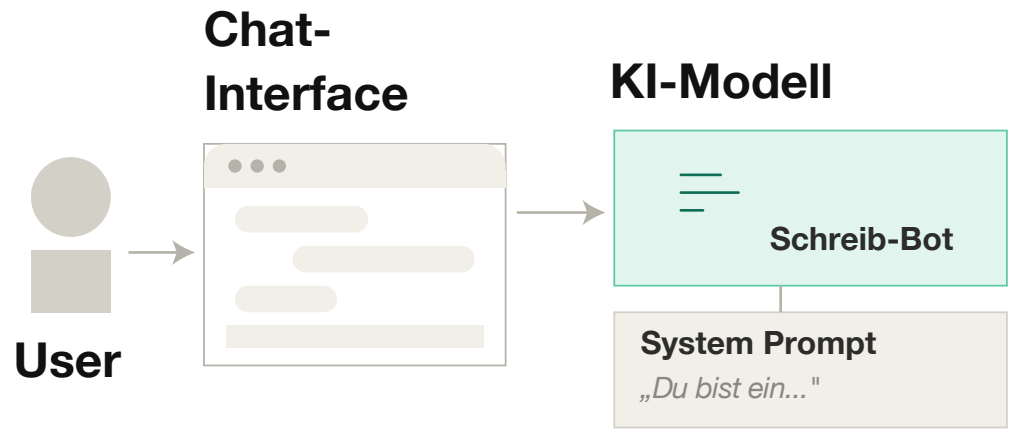
Bildquelle: <https://unsplash.com/photos/1K9T5YiZ2WU>

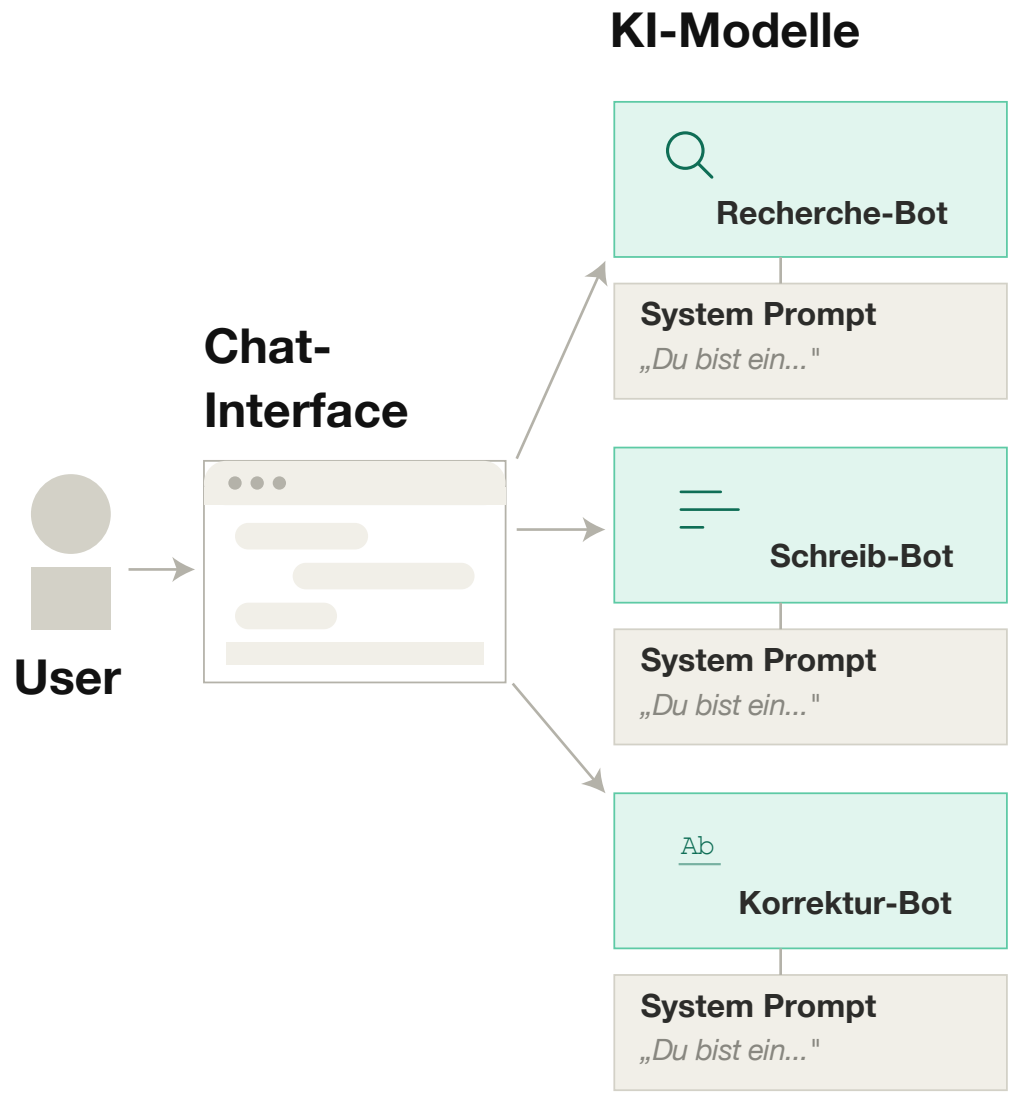
## KI-Modell



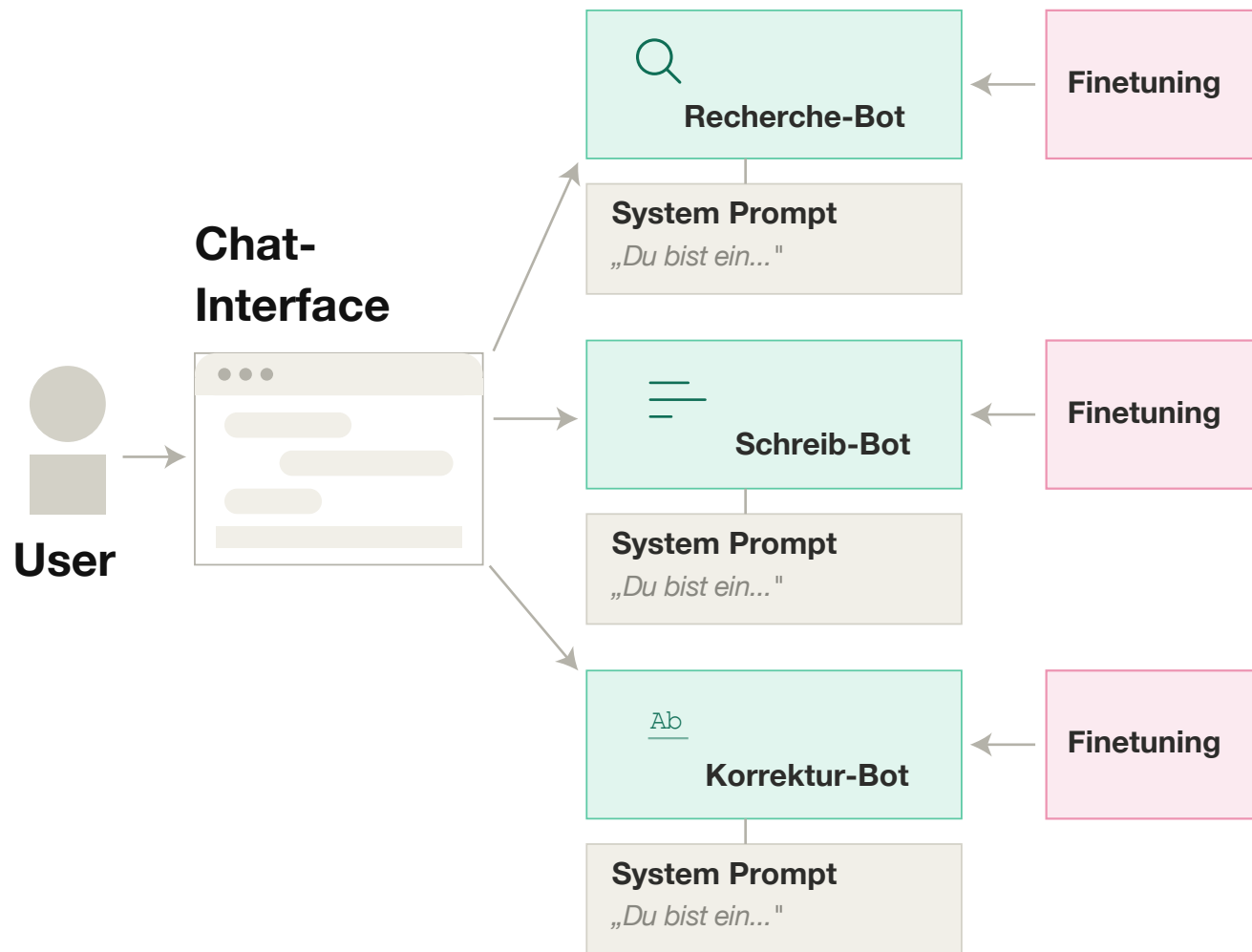
## KI-Modell

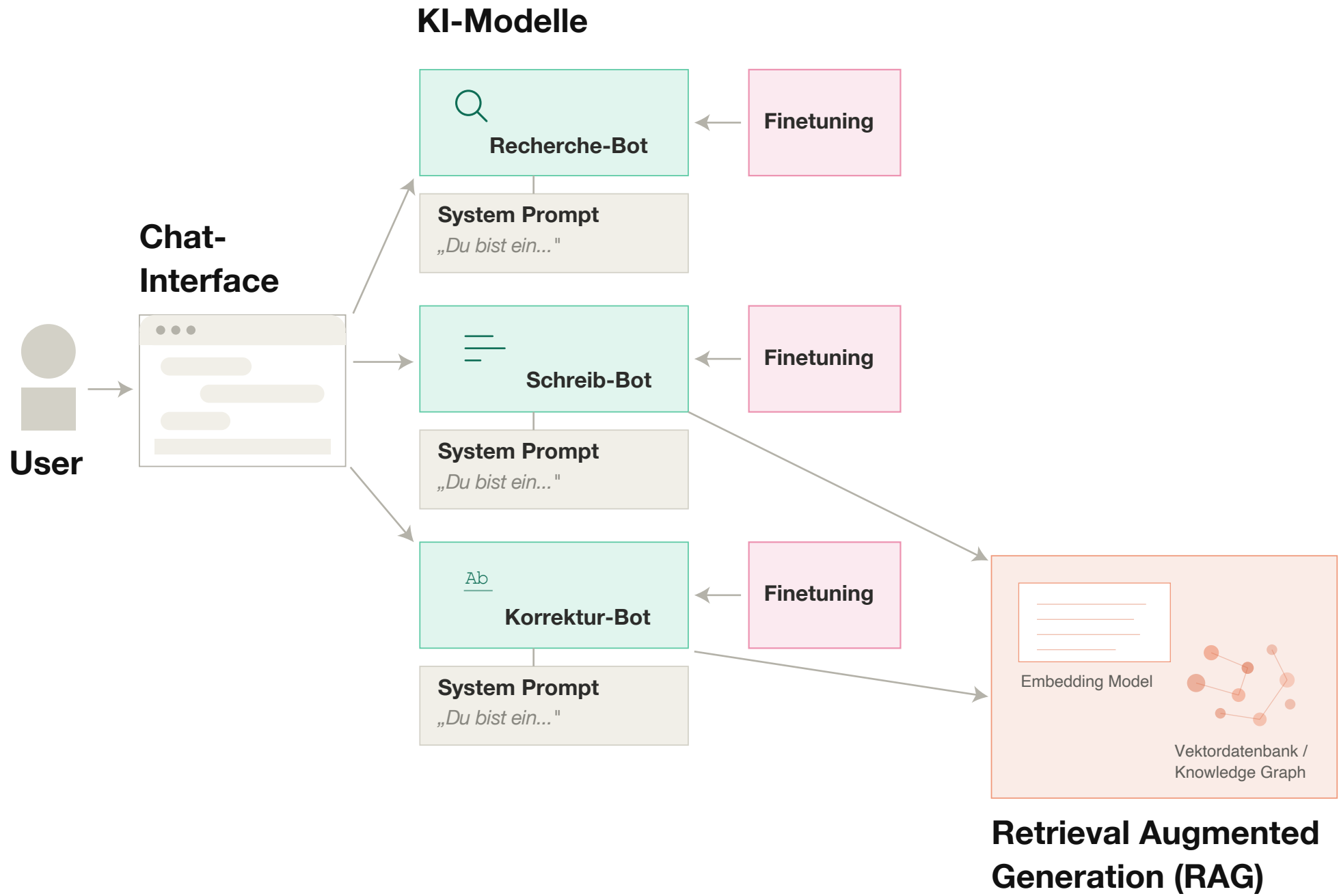


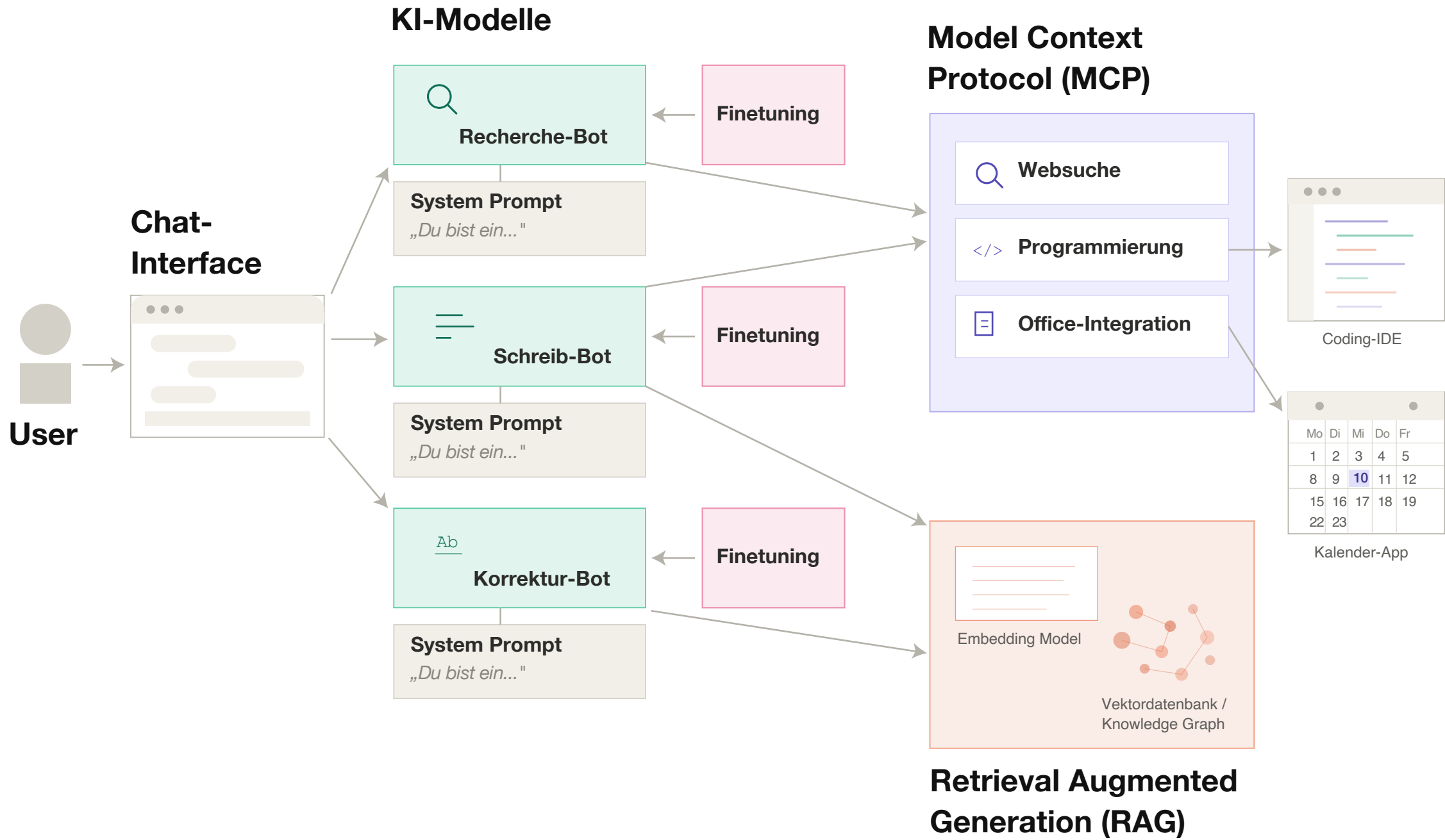




# KI-Modelle







## KI-Modelle

## Model Context Protocol (MCP)

## Retrieval Augmented Generation (RAG)

User

Chat-Interface

Recherche-Bot

Schreib-Bot

Korrektur-Bot

Finetuning

Finetuning

Finetuning

Websuche

Programmierung

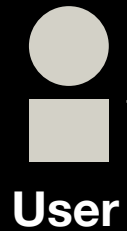
Office-Integration

Embedding Model

Vektordatenbank / Knowledge Graph

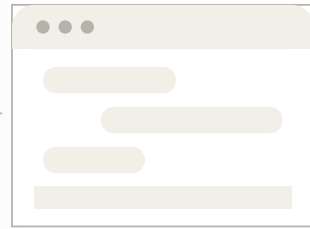
Coding-IDE

Kalender-App



User

### Chat-Interface



### KI-Modelle

Recherche-Bot

System Prompt  
„Du bist ein...“

Schreib-Bot

System Prompt  
„Du bist ein...“

Korrektur-Bot

System Prompt  
„Du bist ein...“

Finetuning

Finetuning

Finetuning

### Model Context Protocol (MCP)

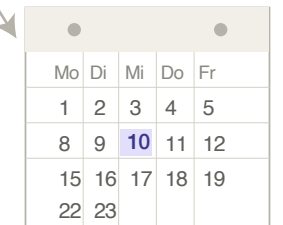
Web Suche

Programmierung

Office-Integration



Coding-IDE



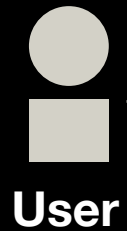
Kalender-App

Embedding Model

Vektordatenbank / Knowledge Graph

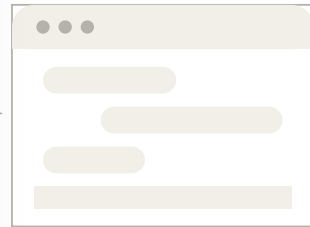
### Retrieval Augmented Generation (RAG)

# Custom Chatbot

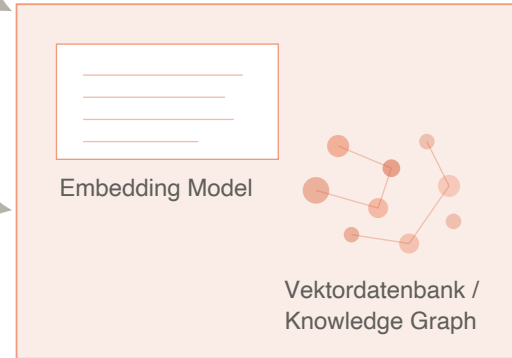
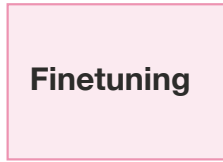
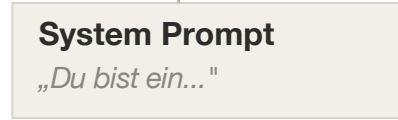
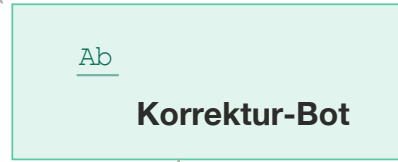
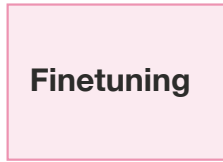
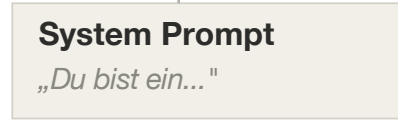
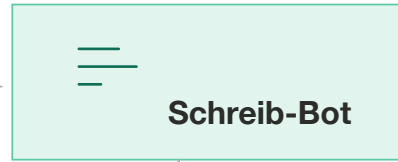
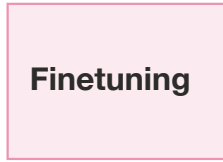
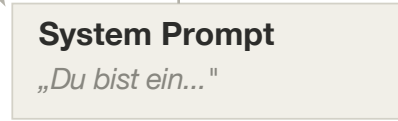
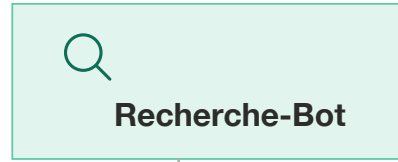


User

### Chat-Interface

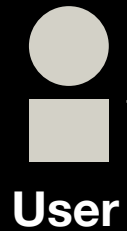


### KI-Modelle



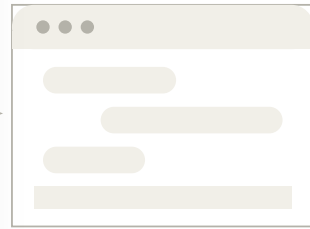
Retrieval Augmented Generation (RAG)

# Custom Chatbot

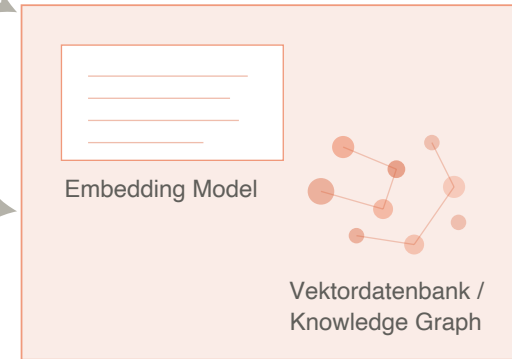
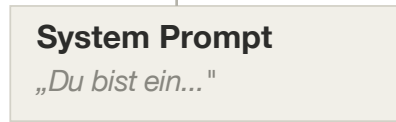
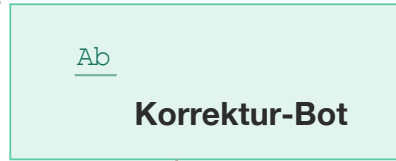
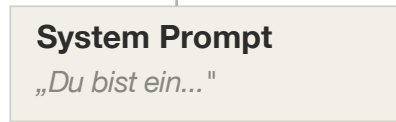
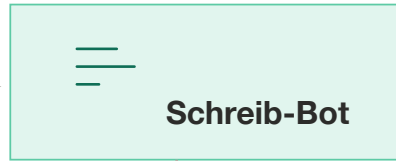
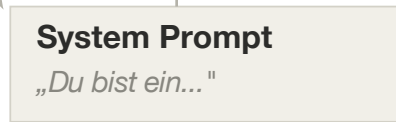
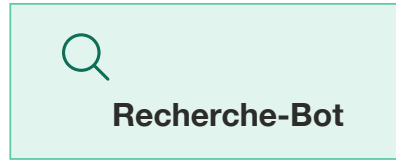


User

## Chat-Interface

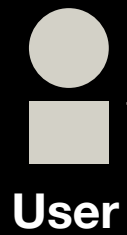


## KI-Modelle



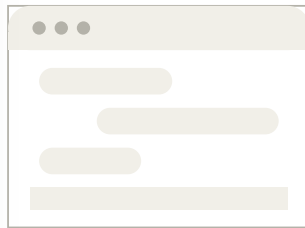
Retrieval Augmented Generation (RAG)

# Custom Chatbot

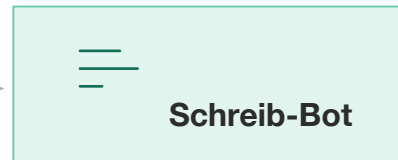


User

Chat-Interface



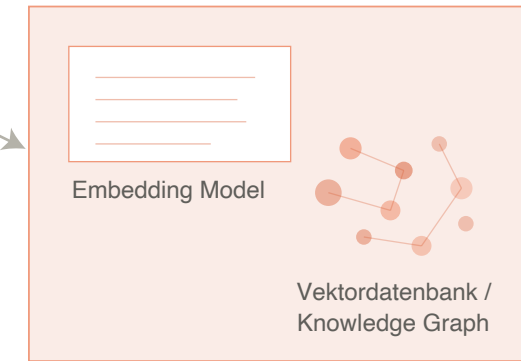
KI-Modell



Schreib-Bot

System Prompt

„Du bist ein...“



Retrieval Augmented Generation (RAG)

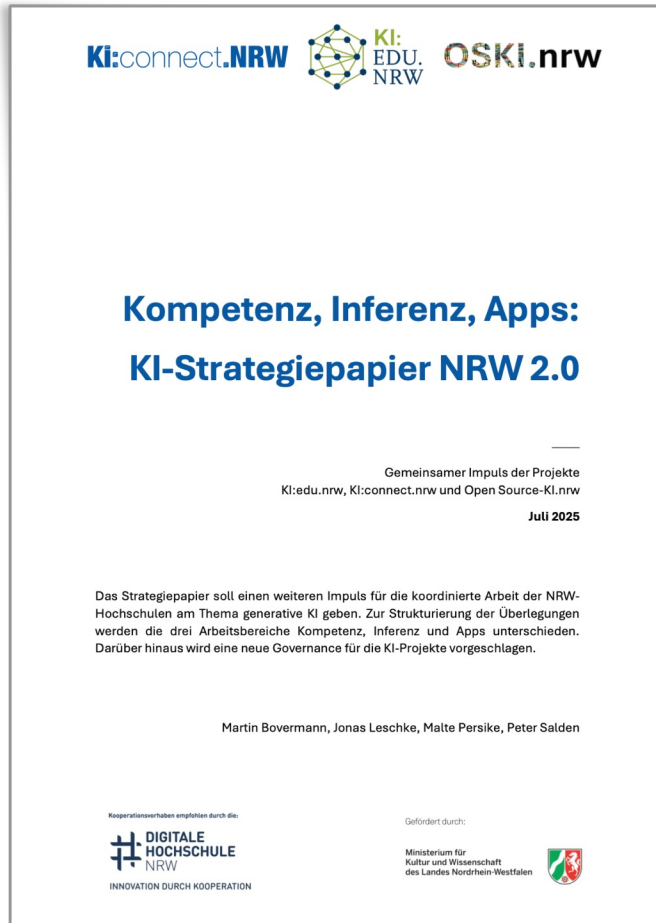
Bot-Camp



# 3 Risiken für die Antwortqualität



# 1 Wir laden PDFs hoch



# 2 Die PDFs werden in Sinneinheiten („Chunks“) gestückelt

Kompetenz, Inferenz, Apps:  
KI-Strategiepapier NRW 2.0

Gemeinsamer Impuls der  
Projekte KI:edu.nrw,  
KI:connect.nrw und Open  
Source-KI.nrw

Juli 2025

Das Strategiepapier soll einen  
weiteren Impuls für die  
koordinierte Arbeit der NRW-  
Hochschulen am Thema  
generative KI geben. Zur ...

# 3 Die Chunks werden in einer Datenbank gespeichert



Vector  
Datenbank

# 4 Relevante Chunks werden abgerufen und in die KI-Antworten eingebaut

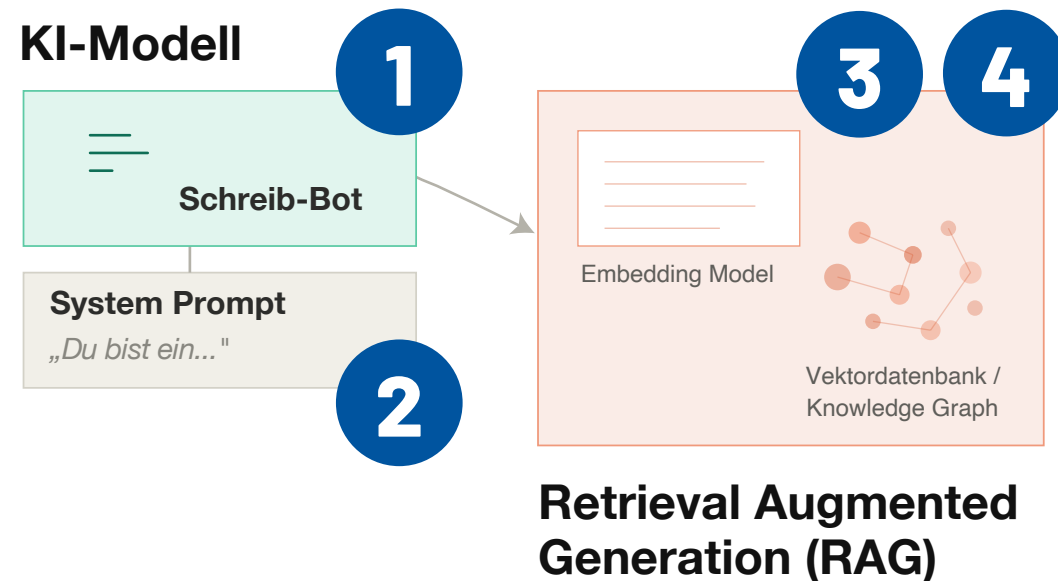


Im Juli 2025 wurde das Strategiepapier "Kompetenz, Inferenz, Apps: KI-Strategiepapier NRW 2.0" veröffentlicht. Es enthält einen Impuls für die koordinierte Arbeit der NRW-Hochschulen am Thema generative KI.



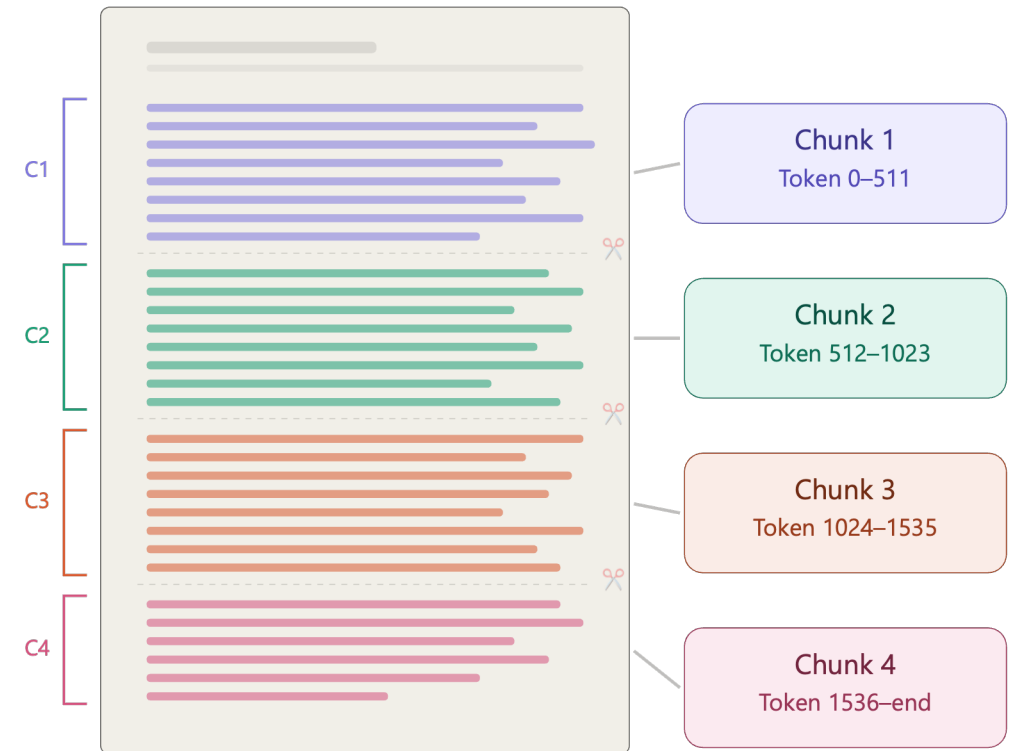
## Risiken für die Antwortqualität

- 1 Trivialer Fall:** KI-Halluzinationen
- 2 Vermeidbar:** Defizite im System Prompt
- 3 Problemfeld 1:** Chunking-Methoden beim RAG
- 4 Problemfeld 2:** Retrieval-Methoden beim RAG



# Einfluss der Chunking-Methode

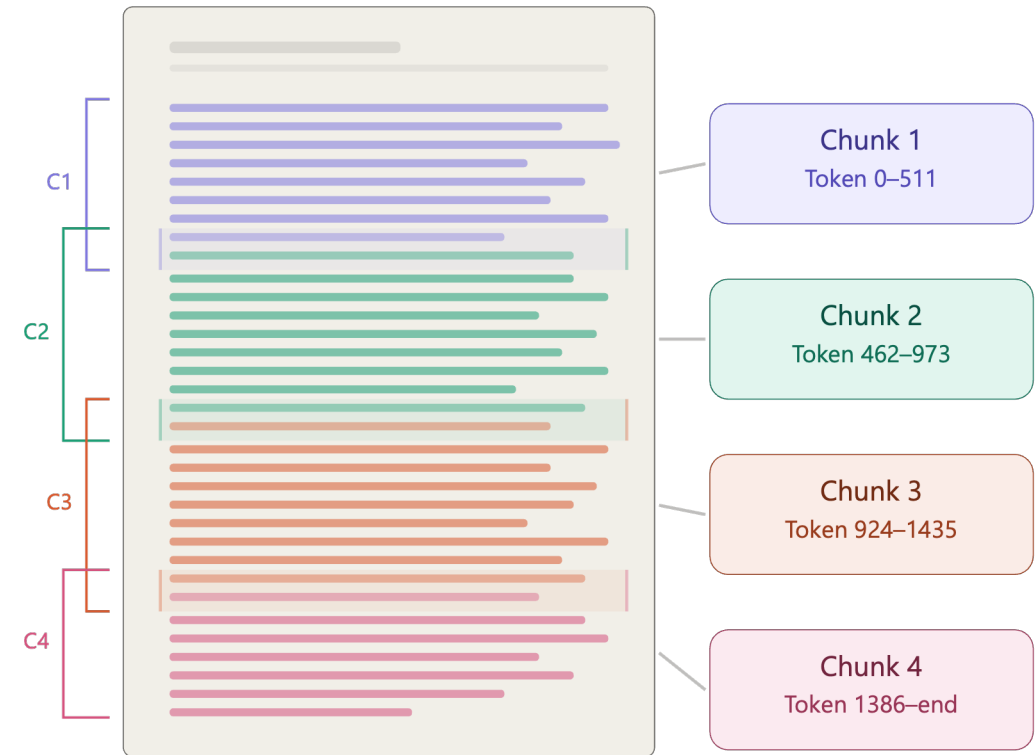
**Fixed Chunking:** Chunks sind immer eine feste Anzahl von Tokens.



# Einfluss der Chunking-Methode

**Fixed Chunking:** Chunks sind immer eine feste Anzahl von Tokens.

**Sliding Window:** Chunks überlappen sich, um den Sinngehalt der umgebenden Chunks zu erhalten

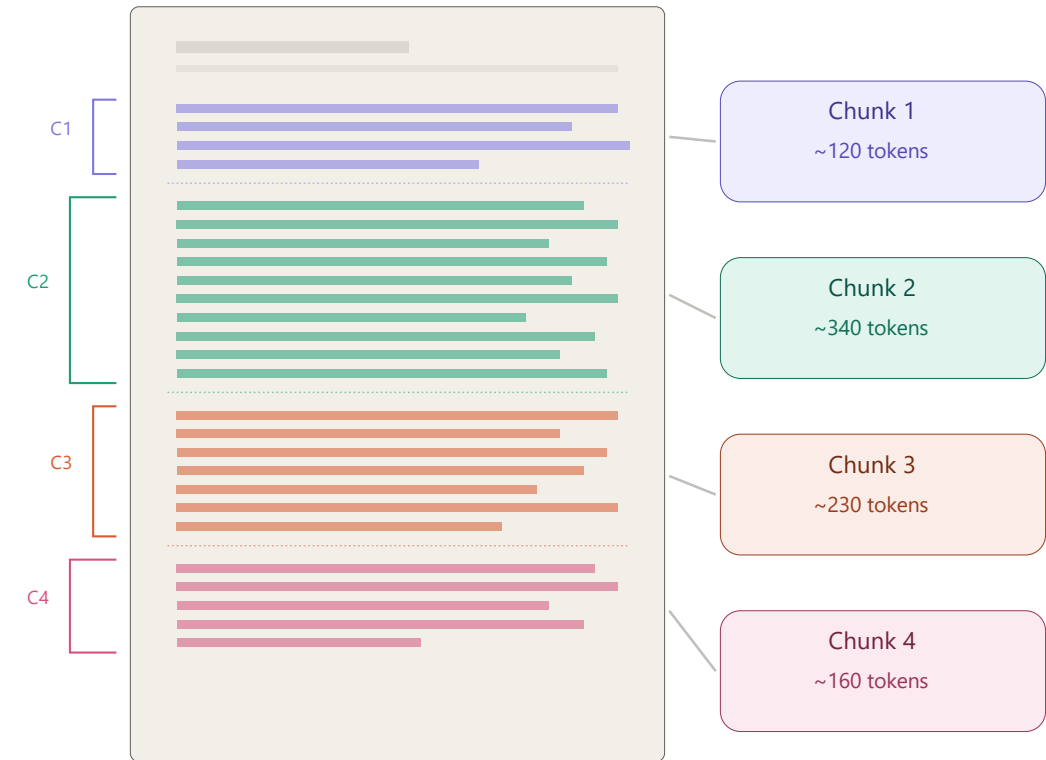


# Einfluss der Chunking-Methode

**Fixed Chunking:** Chunks sind immer eine feste Anzahl von Tokens.

**Sliding Window:** Chunks überlappen sich, um den Sinngehalt der umgebenden Chunks zu erhalten

**Semantic Chunking:** Sinneinheiten werden geschunkt.



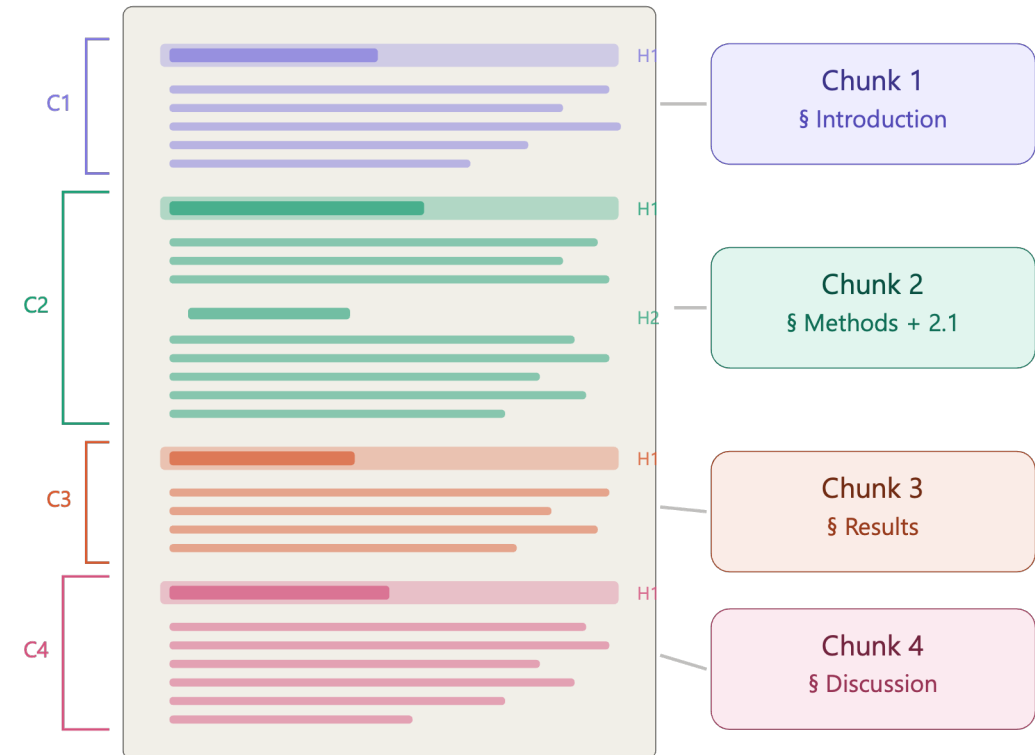
# Einfluss der Chunking-Methode

**Fixed Chunking:** Chunks sind immer eine feste Anzahl von Tokens.

**Sliding Window:** Chunks überlappen sich, um den Sinngehalt der umgebenden Chunks zu erhalten

**Semantic Chunking:** Sinneinheiten werden gechunkt.

**Document Aware Chunking:** Chunks folgen Strukturelementen des Dokuments (z.B. Überschriften).



## **Einfluss der Chunking-Methode**

**Fixed Chunking:** Chunks sind immer eine feste Anzahl von Tokens.

**Sliding Window:** Chunks überlappen sich, um den Sinngehalt der umgebenden Chunks zu erhalten

**Semantic Chunking:** Sinneinheiten werden geschunkt.

**Document Aware Chunking:** Chunks folgen Strukturelementen des Dokuments (z.B. Überschriften).

**Vision Based Chunking:** Bildelemente bleiben erhalten.

# Einfluss der Retrieval-Methode

## Sparse / Dense / Hybrid Retrieval

Suche mit einfachen Statistiken oder mit LLM?

## Reranking

Wie werden Fundstellen nach Wichtigkeit sortiert?

## Query Transformation

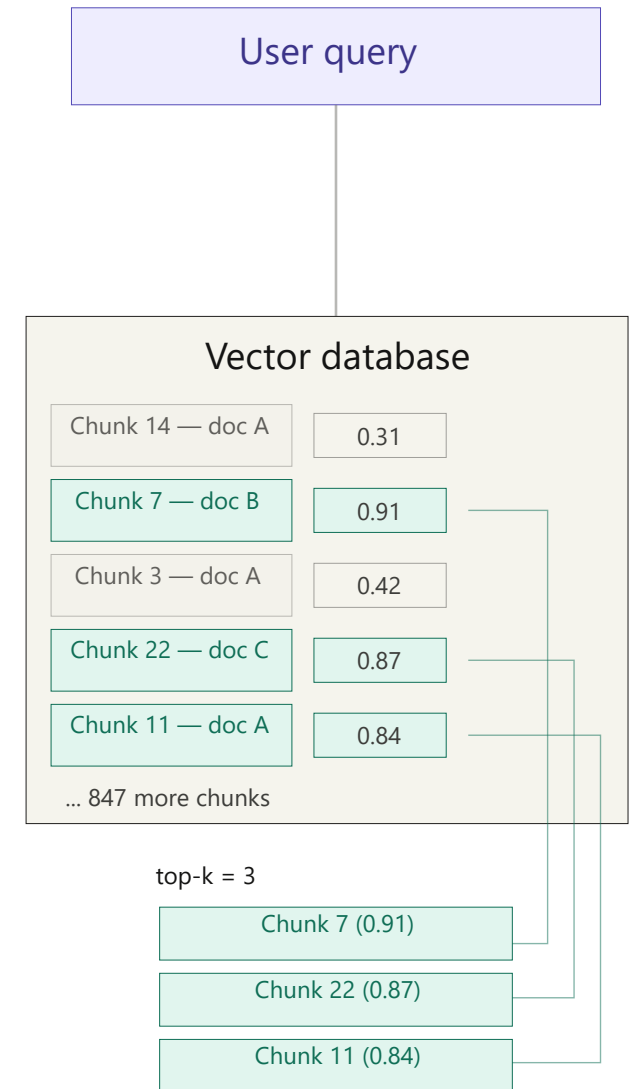
Wie werden Nutzendenanfragen möglichst günstig umformuliert?

## Fortgeschrittene Verfahren

Multi-Hop Retrieval, Contextual Retrieval etc.

## Art der Datenbank

Vektordatenbank oder Knowledge Graph  
Datenbank?



# Was nutzen wir im Bot-Camp?

**Erkenntnis:** Die eine optimale Chunking-Methode für alle Dokumentarten existiert nicht.

Ebenso verhält es sich mit den Retrieval-Methoden.

**Also:** Wir haben Methoden ausgewählt, die ein „möglichst wenig schädlicher“ Kompromiss aus Allgemeingültigkeit und Ressourcenaufwand sind.

**Folge:** Beim RAG für Eure Dokumente werden mit Sicherheit Fehler auftreten.

