



Timmy - Ein Chatbot für das Prüfungsamt

MSI – Teamprojekt – WS2023/24 – SS2024

Benjamin Brünau, Joel Merath, Jan Schmidt, Khadidja Kebaili

KADDI



Einführung

Gliederung

- Einführung
- Retrieval Augmented Generation (RAG)
- Evaluierung des RAG
- Finetuning
- Empfehlung
- Ausblick
- Demo

Einführung

- **Zielsetzung:**
Machbarkeitsprüfung für einen Chatbot für das Prüfungsamt.
- **Hintergrund:**
Ein Bot soll häufig gestellte Fragen (FAQ) beantworten, indem er mit Dokumenten trainiert wird, die vom Prüfungsamt zur Verfügung gestellt wurden.
- **Evaluierung:**
Leistungsfähigkeit des Chatbots, wie Korrektheit und Performance, werden anhand verschiedene Evaluierungsmetriken getestet.

Einführung

- **Verwendete Dokumente:**
 - Zulassungssatzung für die Masterstudiengänge (ZuSMa)
 - SPO Nr. 5 - Studiengang Informatik (MSI)
- **Machbarkeitsprüfung folgt primär für lokale LLMs (IOS Server)**
- **Methoden:**
 - Retrieval Augmented Generation (RAG) → unser Ansatz
 - Evaluierung
 - Finetuning
 - Prototyp

JAN

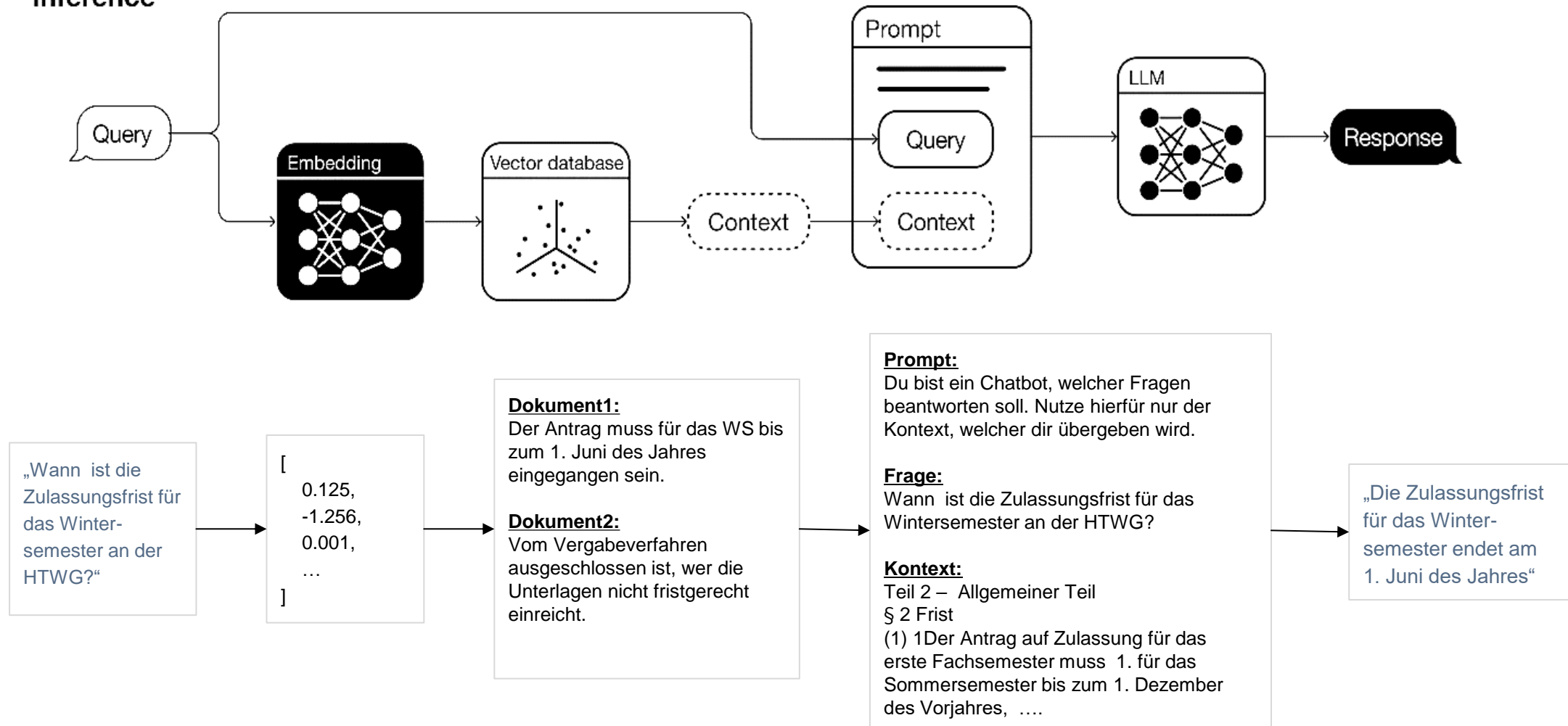


Retrieval Augmented Generation (RAG)



RAG - Pipeline

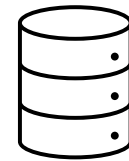
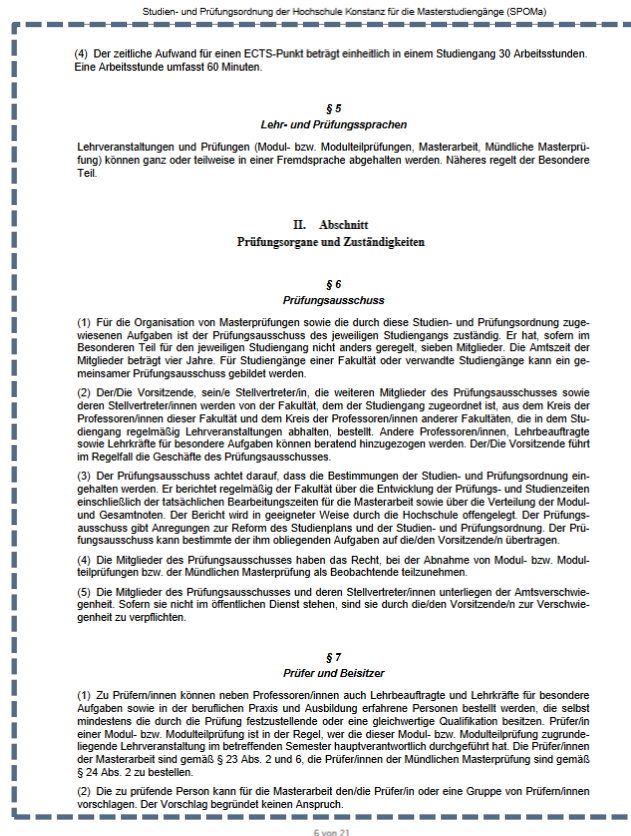
Inference



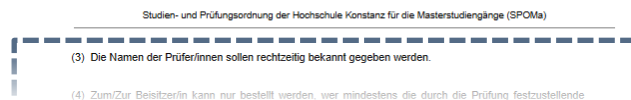
Retrieval

Erstellung der Datenbasis

- ➔ 1. Zuschneiden des Texts
2. Verbinden aller Seiten zu einem Text
3. Splitten nach X Zeichen an speziellen Zeichen, Bsp: (\n\n, \n\n\$)
4. Embedden & Persistieren



Vector DB



Retrieval

Vergleich verschiedener Modelle

Modell	Knn	Context size 1024	Context size 2048	Context size 4096
A (all-MiniLM-L6-v2)	1	0.24	0.28	0.21
B (all-mpnet-base-v2)	1	0.41	0.48	0.45
C (paraphrase-multilingual-MiniLM-L12-v2)	1	0.52	0.52	0.48

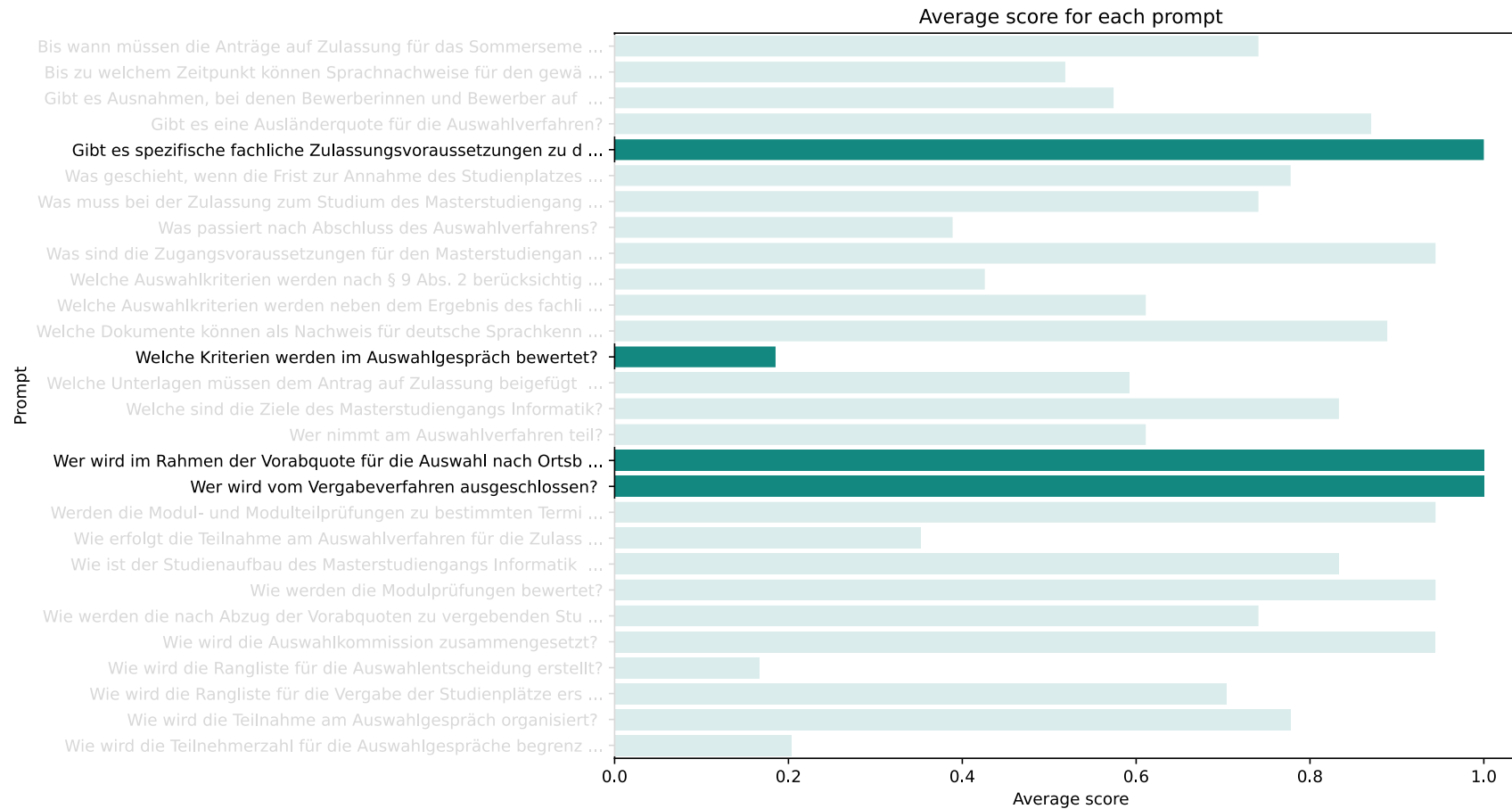
- Ergebnisse variieren stark je nach verwendetem Modell
- Die Context-Größe ändert die Qualität nur geringfügig

* Die Werte geben an, wie häufig (in Prozent) der richtige Absatz zu einer Frage gefunden wurde

* Grün hinterlegtes Modell ist das Beste in der Vergleichsgruppe

Retrieval

Analyse der Fragen



Retrieval Ergebnisse

Datenverarbeitung/ Aufbereitung:

- Datennahe Vorbereitung sinnvoll
 - erhöhte bei KNN (K-Nächste-Nachbarn) = 5 die Erfolgsrate um ca. 17%
- Tabellarische Daten schlecht extrahierbar

Settings:

- Embedding Model: paraphrase-multilingual-MiniLM-L12-v2
- Context size: 2048

Anwendung:

- Schrittweise Vergrößerung des KNN-Werts ist ratsam
(Bsp. LLM entscheiden lassen, ob die Infos in den Teilstücken enthalten sind)

JOEL

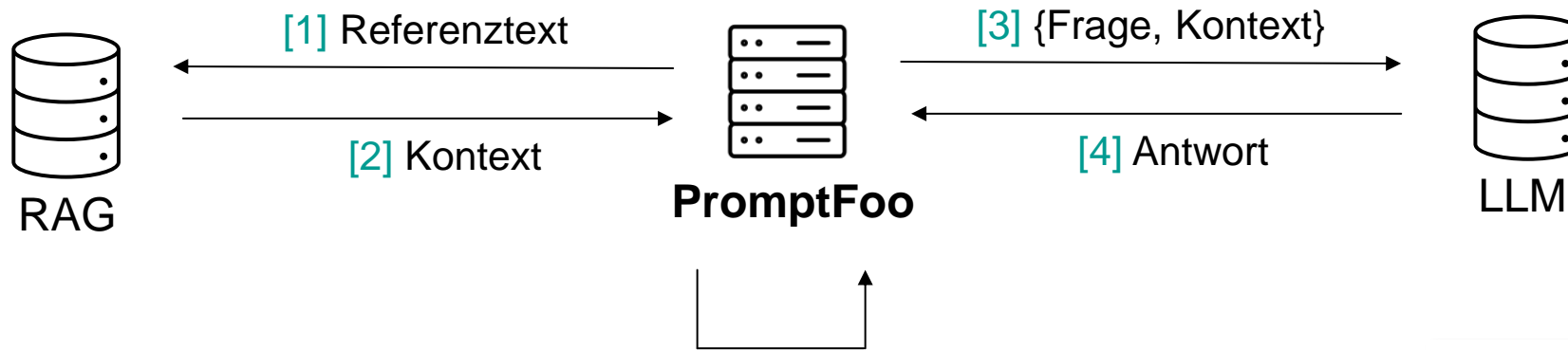


Evaluierung des RAG

PromptFoo

Testlaufautomatisierung

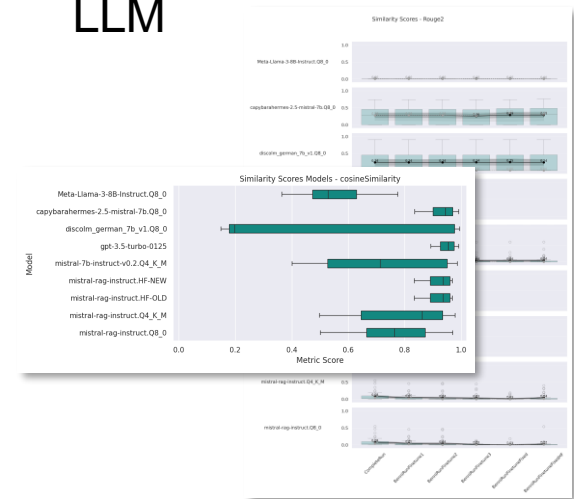
Loop: Für jeden Testfall



[5] Evaluation (Referenz vs. Antwort)

Metriken:
CosineSimilarity, KeywordScore,...

[6] Plotting



Testfälle

Basierend der SPO

- 15 Testfälle

- min. 1 Musterantwort

- Correctness (0.0 - 1.0)

- min. 3 Beispielfragen

- Specificity (0.0 - 1.0)
 - Relevance (0.0 - 1.0)

- 8 Metriken

- 9 LLMs

- Llama 3 (8b/Q8)
 - CapybaraHermesMistral (7b/Q8)
 - DiscoLM (7b/Q8)
 - Mistral (7b/Q4)

- ...

→ ~ 400 Testläufe = ~ 15.000 Datenpunkte

Fragebogen: Wie kann müssen die Anträge auf Zulassung für das Sommersemester und das Wintersemester 2023/2024 auf Zulassung für das Sommersemester 2024/2025 werden?

Antwort: Die Anträge auf Zulassung für das Sommersemester 2024/2025 werden...

cosineSimilarity	0.95
jaccardSimilarity	0.95
rouge1	0.95
rouge2	0.95
rougeL	0.95
editDistance	0.95
keywordScore	0.95
editDistance	0.95

Metriken

- **Semantik** (höher gewichtet):
 - CosineSimilarity
 - KeywordScore
 - ChatGPTScore
- **Wort- und Textebene** (niedriger gewichtet):
 - Edit Distance
 - Rouge - 1
 - Rouge - 2
 - Rouge - L
 - Jaccard Similarity

Metrik

CosineSimilarity

Definition

- Eine Metrik, die die Ähnlichkeit zwischen zwei Vektoren im Raum der Merkmale misst, indem der **Kosinus des Winkels** zwischen ihnen berechnet wird.

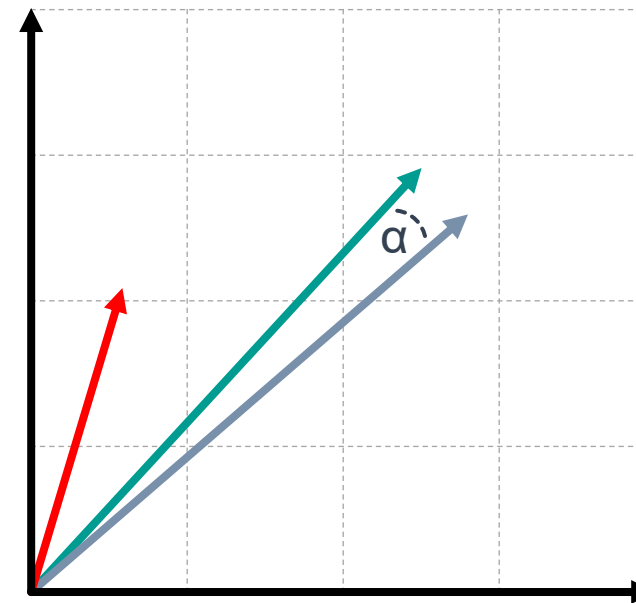
Die Relevanz für uns:

- Ermöglicht eine effektive Bewertung der semantischen Ähnlichkeit zwischen generierten Antworten und den erwarteten Antworten.
- Berücksichtigt auch, falls die Wortwahl in einer Antwort variiert, der Inhalt aber gleich bleibt.

A: "gemäß dem Terminplan der Hochschule Konstanz"

B: "laut dem Zeitplan der Hochschule Konstanz"

C: "Nachweis über einen Hochschulabschluss"



Beispielhafte Darstellung

$$\alpha = 6.0^\circ$$

CosineSimilarity
Zwischen A u. B

→ 0.96
(oder 96.0%)

Metrik

KeywordScore

Definition

- Ermittelt, ob eine Anzahl an definierten Keywords in der Antwort vorhanden sind:

Die Relevanz für uns:

- Prüfen, ob die wichtigsten Kerninformationen als Keyword vorhanden ist (wichtig bei genauer Zitierung)

Keywords:

['Bewerbungsschluss', 'Antrag auf Zulassung', '§ 6 Abs. 5']



Antwort:

„Sprachnachweise für den gewählten Studiengang, die bis zum **Bewerbungsschluss** nicht vorgelegt werden können, [...], für das der **Antrag auf Zulassung** gestellt wurde, nachgereicht werden. Die Zulassung erfolgt in diesem Fall gemäß **§ 6 Abs. 5** unter Vorbehalt.“

→ **KeywordScore:** 3/3 = 1.00 (oder 100%)

Metrik

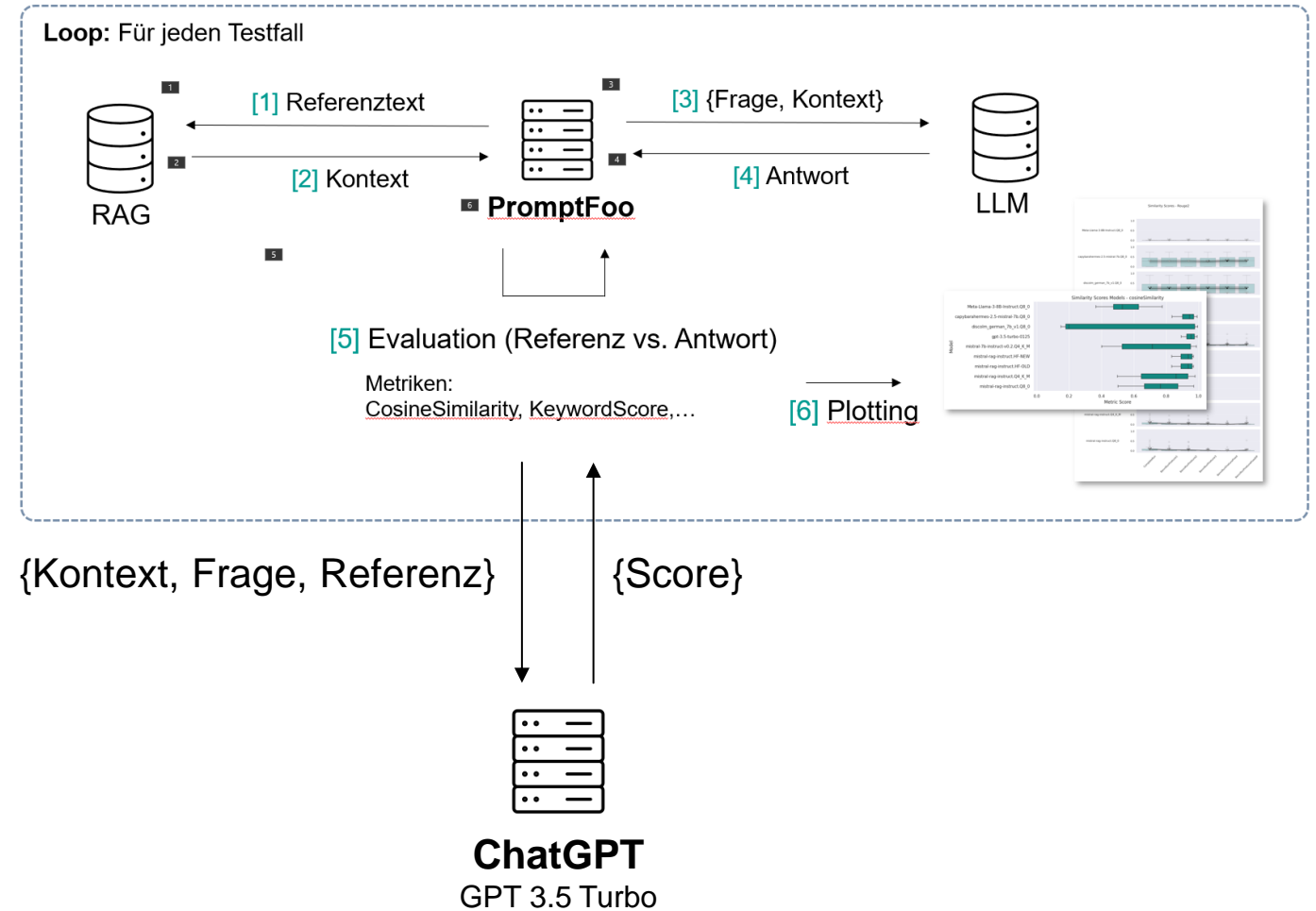
ChatGPT Score

Definition

- Nutzung der OpenAI-Schnittstelle, um die Ähnlichkeit zwischen der Antwort des LLMs und dem Kontext zu bewerten.

Die Relevanz für uns:

- Diese externe Validierung stärkt die Auswahl des LLMs und untermauert die Ergebnisse weiterer Metriken.

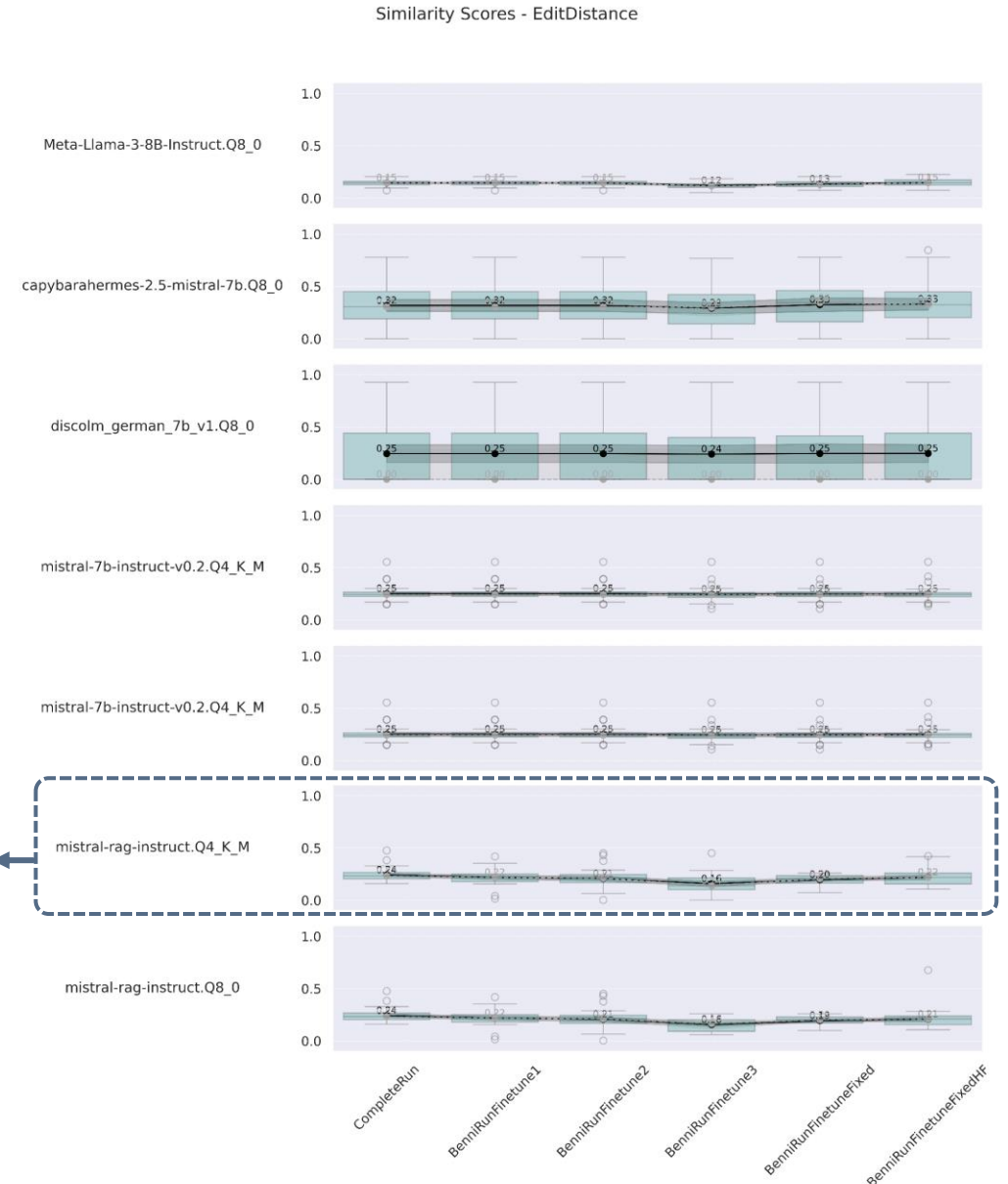
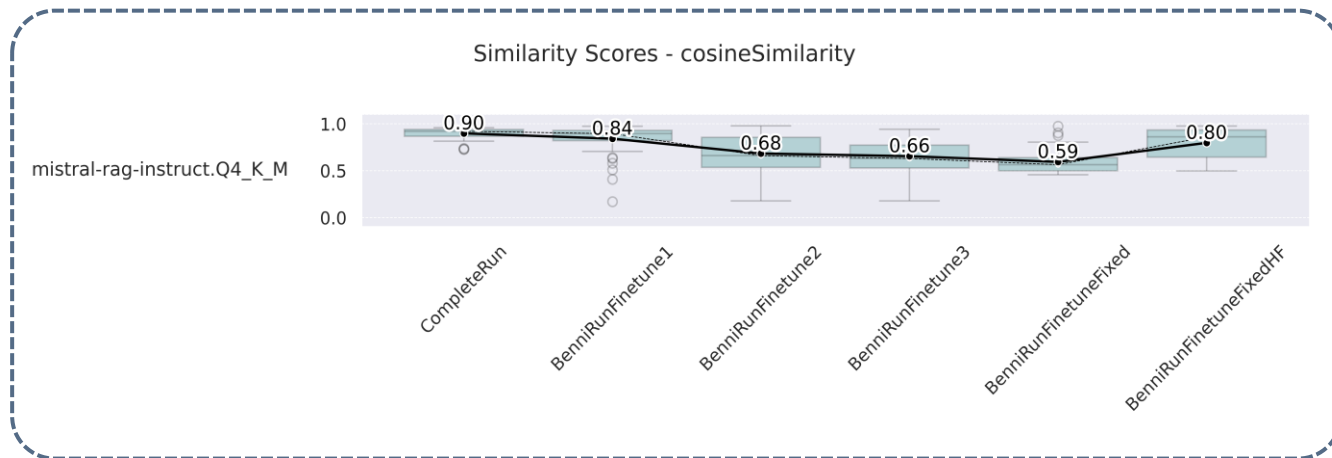


Weitere Metriken

Metrik	Definition	Vorteil
Edit Distance	Minimale Anzahl von Bearbeitungsschritten (Einfügungen, Löschungen, Ersetzungen), von eine Zeichenkette Antwort → Referenz	Berücksichtigt Ähnlichkeit auf Textebene
Rouge - 1	Einzelwort Übereinstimmungen [Terminplan, <u>der</u> , Hochschule, Konstanz] [Zeitplan, <u>der</u> , Hochschule, Konstanz]	Vergleicht nicht nur einzelne Wörter, sondern auch Bigramme und längere Sequenzen. → Feingranulare Bewertung der Textähnlichkeit
Rouge - 2	Bigram-Übereinstimmungen [Terminplan <u>der</u> , <u>der Hochschule</u>] [Zeitplan <u>der</u> , <u>der Hochschule</u>]	
Rouge - L	Die längste gemeinsame Teilfolge (LCS) <u>gemäß dem Terminplan der Hochschule</u> laut <u>dem Zeitplan der Hochschule</u>	
Jaccard Similarity	Misst die Ähnlichkeit zwischen zwei Mengen von Elementen [gemäß, <u>dem</u> , Terminplan, <u>der</u> , Hochschule, Konstanz] [laut, <u>dem</u> , Zeitplan, <u>der</u> , Hochschule, Konstanz] [gemäß, laut, dem, Terminplan, Zeitplan, ..., Konstanz]	Bewertung von Textähnlichkeiten, bei denen die <u>Präsenz und Verteilung</u> von Wörtern von Bedeutung sind.

Plotting Test-Run

- Auswertung mehrerer Test-Runs
 - Ermöglicht direkte Analyse von Veränderungen
 - Bsp.: Finetuning
 - Ermöglicht einen Gesamteindruck der Prozesse

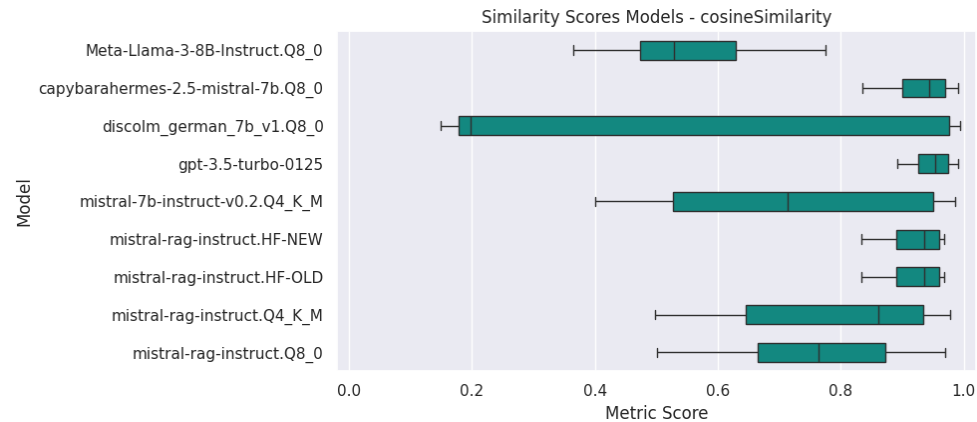


Plotting

Einzelauswertungen

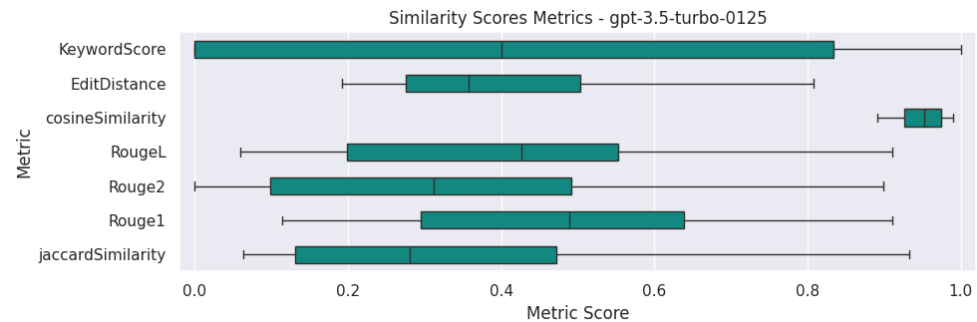
- **Spezifisches Metrik Ergebnisse** bezogen auf einen spezifischen Test-Run für jedes Modell

- Bsp.: **CosineSimilarity**



- **Alle Metrik Ergebnisse** bezogen auf einen spezifischen Test-Run für ein bestimmtes Modell

- Bsp.: **gpt-3.5-turbo-0125**

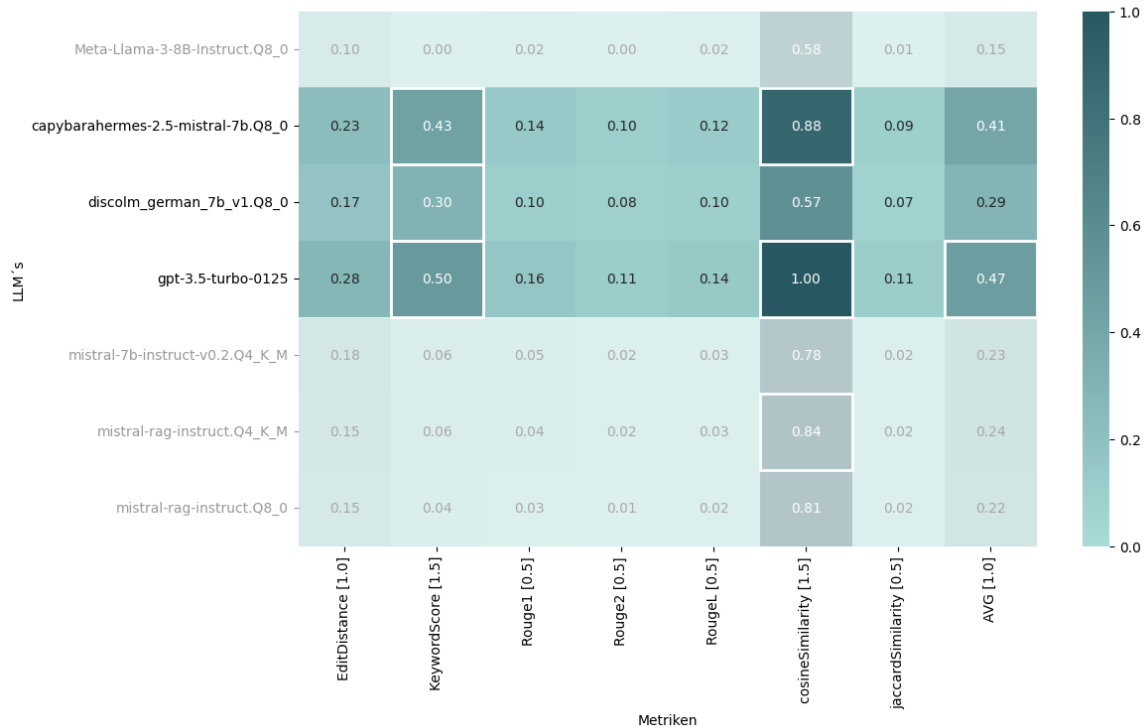


20240805_231203
conclusions
models
combined
timeline_models_capybarahermes-2.5-mistral-7b...
timeline_models_capybarahermes-2.5-mistral-7b...
timeline_models_discolm_german_7b_v1.Q8_0_ci.p
timeline_models_discolm_german_7b_v1.Q8_0.png
timeline_models_gpt-3.5-turbo-0125_ci.png
timeline_models_gpt-3.5-turbo-0125.png
timeline_models_Meta-Llama-3-8B-Instruct.Q8_0_.
timeline_models_Meta-Llama-3-8B-Instruct.Q8_0...
timeline_models_mistral-7b-instruct-v0.2.Q4_K_M.
timeline_models_mistral-7b-instruct-v0.2.Q4_K_M..
timeline_models_mistral-rag-instruct.HF-NEW_ci.p.
timeline_models_mistral-rag-instruct.HF-NEW-OL..
timeline_models_mistral-rag-instruct.HF-NEW.png
timeline_models_mistral-rag-instruct.HF-OLD_ci.pr
timeline_models_mistral-rag-instruct.HF-OLD.png
timeline_models_mistral-rag-instruct.Q4_K_M_ci.p.
timeline_models_mistral-rag-instruct.Q4_K_M.png
timeline_models_mistral-rag-instruct.Q8_0_ci.png
timeline_models_mistral-rag-instruct.Q8_0.png
singles
singles
metrics
models
singles_models_capybarahermes-2.5-mistral-7b.Q8..
singles_models_capybarahermes-2.5-mistral-7b.Q8..
singles_models_capybarahermes-2.5-mistral-7b.Q8..
singles_models_capybarahermes-2.5-mistral-7b.Q8..
singles_models_discolm_german_7b_v1.Q8_0_Relev.
singles_models_discolm_german_7b_v1.Q8_0_simila.
singles_models_discolm_german_7b_v1.Q8_0_Specif
singles_models_discolm_german_7b_v1.Q8_0_Specif
singles_models_gpt-3.5-turbo-0125_Relevance.png
singles_models_gpt-3.5-turbo-0125_similarity.png
singles_models_gpt-3.5-turbo-0125_Specificity_Rele
singles_models_gpt-3.5-turbo-0125_Specificity.png
singles_models_Meta-Llama-3-8B-Instruct.Q8_0_Rel.
singles_models_Meta-Llama-3-8B-Instruct.Q8_0_si...

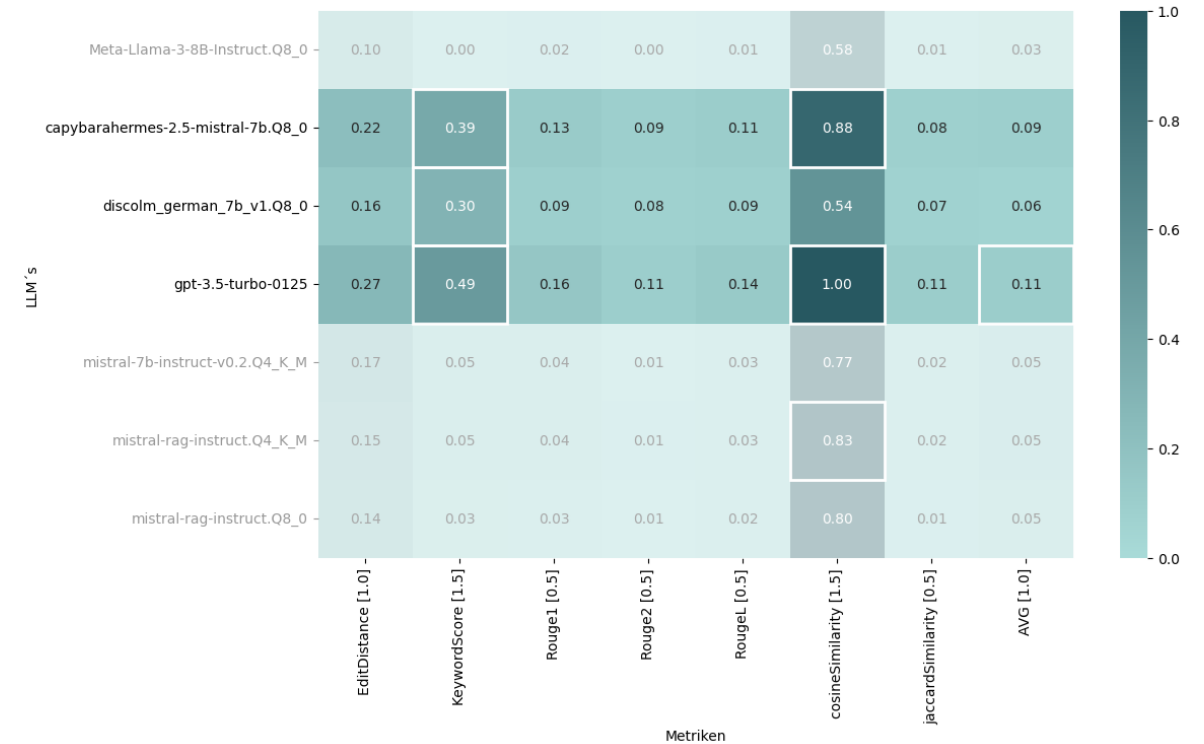
Auswertung

Gewichtet und Ungewichtet

Gewichtete Similarity der LLM's mit Metriken



Gewichtete Similarity der LLM's mit Metriken (unter Berücksichtigung von Specificity und Relevance)

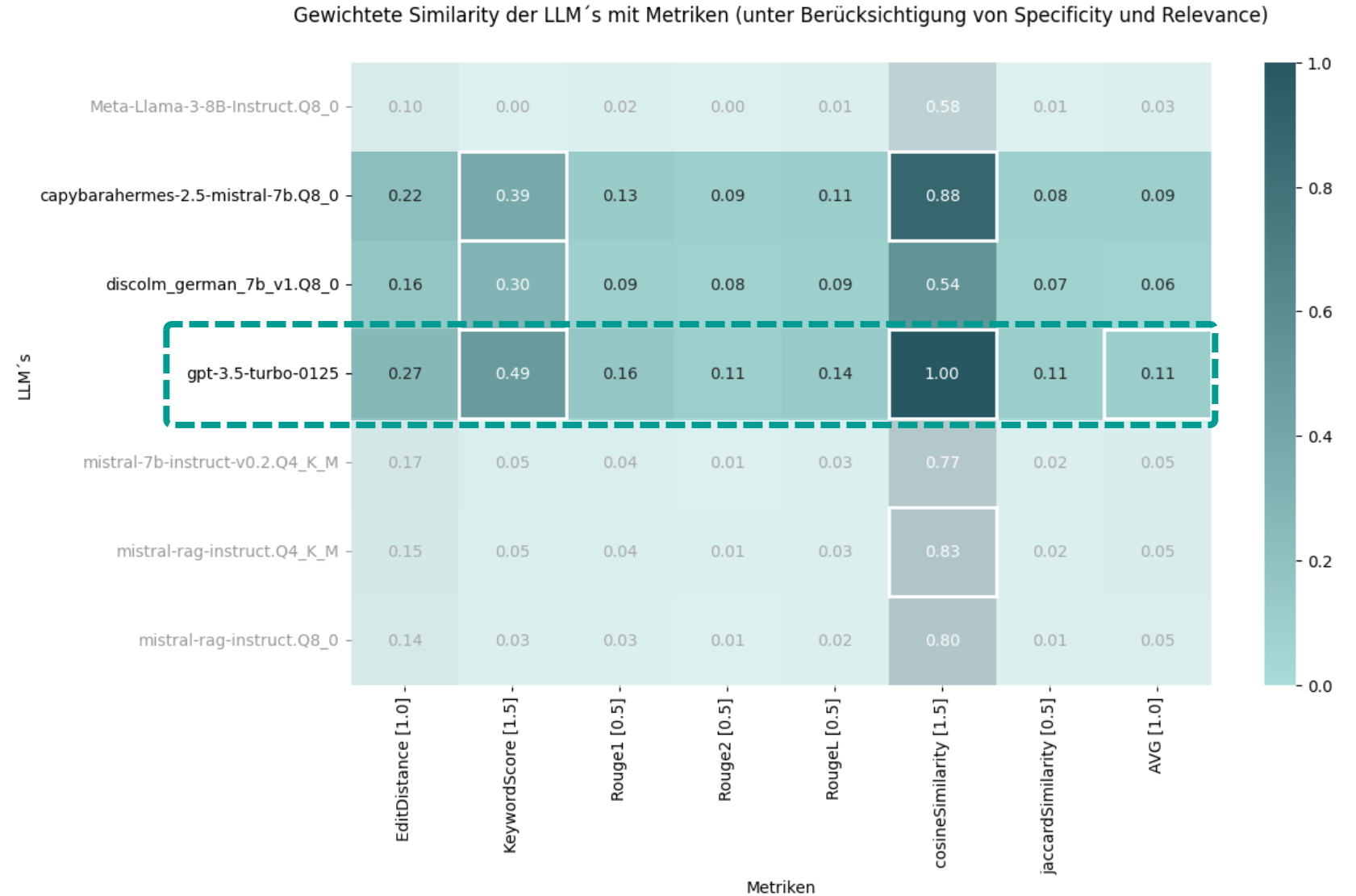


→ **gpt-3.5-turbo-0125** stellt sich als das bestes LLM heraus

Auswertung

Ungewichtet und Gewichtet

→ **gpt-3.5-turbo-0125**
stellt sich als das
bestes LLM heraus



Metrik

EditDistance

Definition

- Minimale Anzahl von Bearbeitungsschritten (Einfügungen, Löschungen, Ersetzungen), von eine Zeichenkette $A \rightarrow B$

Die Relevanz für uns:

- Quantifiziert die Ähnlichkeit zwischen:

Generierten Antworten und den erwarteten Antworten
- Berücksichtigt Ähnlichkeit auf Textebene.

Metrik

Rouge (1,2,L)

Definition

- Eine Metrik, die die Übereinstimmung zwischen zwei Texten bewertet, basierend auf:
 - **Rouge - 1**
Einzelwort Übereinstimmungen
 - **Rouge - 2**
Bigram-Übereinstimmungen
 - **Rouge - L**
Die längste gemeinsame Teilfolge (LCS)

Referenz:

„gemäß dem Terminplan der Hochschule Konstanz“

LLM-Antwort:

„laut dem Zeitplan der Hochschule Konstanz“

1-Gram:

[gemäß, dem, Terminplan, der, Hochschule, Konstanz]

[laut, dem, Zeitplan, der, Hochschule, Konstanz]

→ **Rouge-1:** $4/6 = 0.66$ (oder 66%)

2-Gram:

["gemäß dem", "dem Terminplan", "Terminplan der", "der Hochschule", "Hochschule Konstanz"]

["laut dem", "dem Zeitplan", "Zeitplan der", "der Hochschule ", "Hochschule Konstanz"]

→ **Rouge-2:** $1/5 = 0.2$ (oder 20%)

Metrik

Rouge (1,2,L)

Definition

- Eine Metrik, die die Übereinstimmung zwischen zwei Texten bewertet, basierend auf:
 - **Rouge - 1**
Einzelwort Übereinstimmungen
 - **Rouge - 2**
Bigram-Übereinstimmungen
 - **Rouge - L**
Die längste gemeinsame Teilfolge (LCS)

Referenz:

„gemäß dem Terminplan der Hochschule Konstanz“

LLM-Antwort:

„laut dem Zeitplan der Hochschule Konstanz“

Vergleich:

["gemäß dem", "dem Terminplan", "Terminplan der", "der Hochschule", "Hochschule Konstanz"]

["laut dem", "dem Zeitplan", "Zeitplan der", "der Hochschule ", "Hochschule Konstanz"]

→ **Rouge-2:** $1/5 = 0.2$ (oder 20%)

Metrik

Rouge (1,2,L)

Definition

- Eine Metrik, die die Übereinstimmung zwischen zwei Texten bewertet, basierend auf:
 - **Rouge - 1**
Einzelwort Übereinstimmungen
 - **Rouge - 2**
Bigram-Übereinstimmungen
 - **Rouge - L**
Die längste gemeinsame Teilfolge (LCS)

Referenz:

„gemäß dem Terminplan der Hochschule Konstanz“

LLM-Antwort:

„laut dem Zeitplan der Hochschule Konstanz“

Vergleich:

„gemäß dem Terminplan der Hochschule Konstanz“

„laut dem Zeitplan der Hochschule Konstanz“

- **Länge der LCS:**
4 Wörter [„dem“, „der“, „Hochschule“, „Konstanz“]
 - **Gesamtlänge der Referenz:** 6 Wörter
 - **Gesamtlänge der LLM-Antwort:** 6 Wörter
- **Rouge-L:** $4/6 = 0.67$ (oder 67%)

Metrik

Rouge (1,2,L)

Die Relevanz für uns:

- Rouge (1,2,L) ermöglicht eine feine Bewertung der Textähnlichkeit, indem sie nicht nur einzelne Wörter, sondern auch Bigramme und längere Sequenzen berücksichtigt.

Metrik

Jaccard Similarity

Definition

- Misst die Ähnlichkeit zwischen zwei Mengen von Elementen
- Anzahl der gemeinsamen Elemente zur Gesamtanzahl der Elemente in den Mengen.

Die Relevanz für uns:

- Besonders nützlich für die Bewertung von Textähnlichkeiten, bei denen die Präsenz und Verteilung von Wörtern von Bedeutung sind.

Referenz:

„gemäß dem Terminplan der Hochschule Konstanz“

LLM-Antwort:

„laut dem Zeitplan der Hochschule Konstanz“

Vergleich:

[gemäß, dem, Terminplan, der, Hochschule, Konstanz]

[laut, dem, Zeitplan, der, Hochschule, Konstanz]

Schnittmenge: 4

[gemäß, laut, dem, Terminplan, Zeitplan, der, Hochschule, Konstanz]

→ Jaccard-Similarität

$$= \frac{\text{Anzahl der gemeinsamen Wörter}}{\text{Anzahl der einzigartigen Wörter}} = \frac{4}{8} = 0.5 \text{ (oder 50\%)}$$

BENNI



Finetuning

Ziele / Ansätze

- Dokumente (Kontext) dem LLM beibringen
- RAG verbessern (Darauf finetunen Fragen mithilfe eines Kontexts zu beantworten)

Finetuning Methoden

- **Full Fine Tuning**
 - Alle Gewichte des Models werden aktualisiert
- **Parameter Efficient Fine Tuning (PEFT)**
 - Es werden nur Teile der Gewichte aktualisiert (der Rest bleibt wie zuvor)

Full Fine Tuning

Eine kleine Rechnung...

- Ungefähre Speicheranforderungen (GPU)
 - LLM laden: $parameter_precision_as_bytes \times number_parameters$
 - Finetuning:
 - Optimizer states: 2 states per parameter
 - Gradients: $4\ bytes \times number_parameters$

- Mistral7B (16Bit Float Precision)

LLM laden	$2\ bytes \times 7B\ parameters$	$= 14B\ bytes = 14GB$
Optimizer states	$2 \times 2\ bytes \times 7B\ Parameters$	$= 28GB$
Gradients	$4\ bytes * 7B\ parameters$	$= 28GB$
Total		70 GB

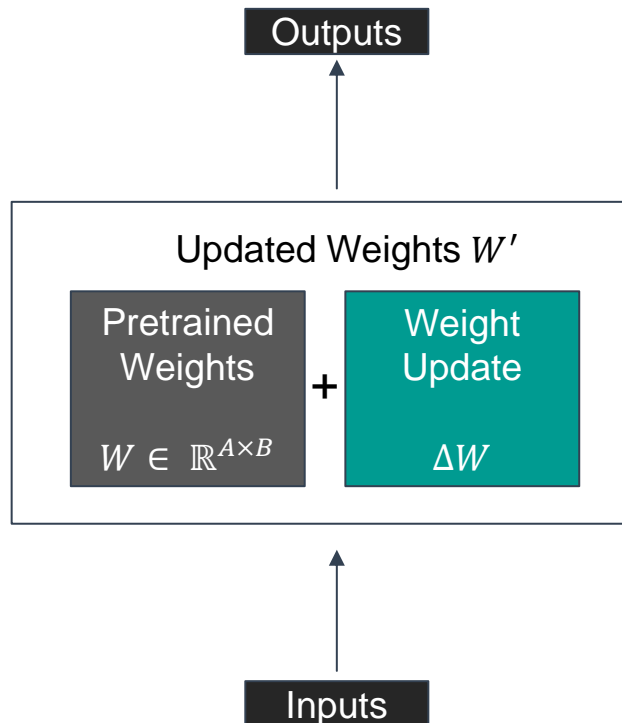
- NVIDIA RTX Titan – 24 GB VRAM

Parameter Efficient Finetuning

LoRA

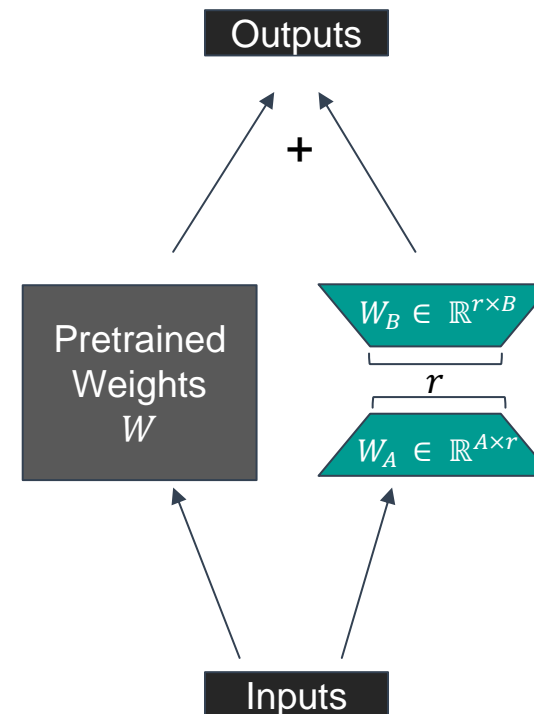
- Low Rank Matrizen werden zum Base Model hinzugefügt
- Funktioniert wie eine Art Adapter (das Base Model bleibt unverändert)

Weight update in **Full Fine Tuning**



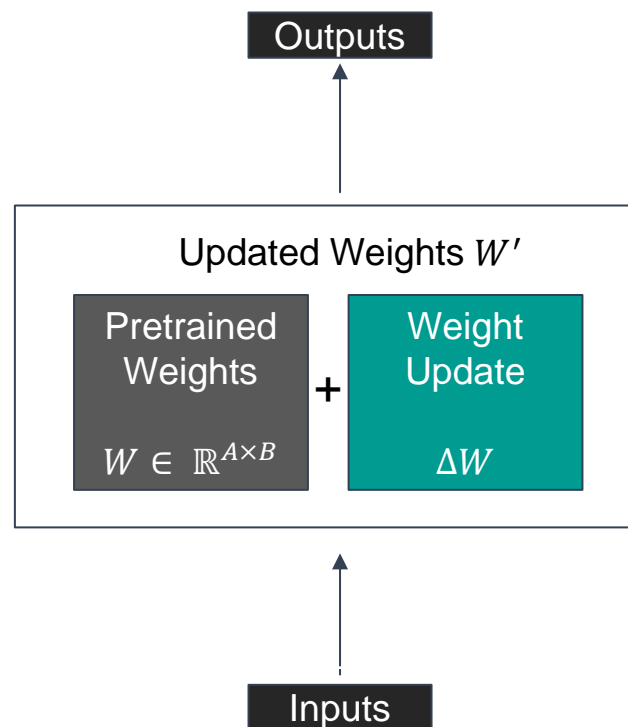
Approximate ΔW
with $\Delta W = W_A \times W_B$

Weight update with **LoRA**

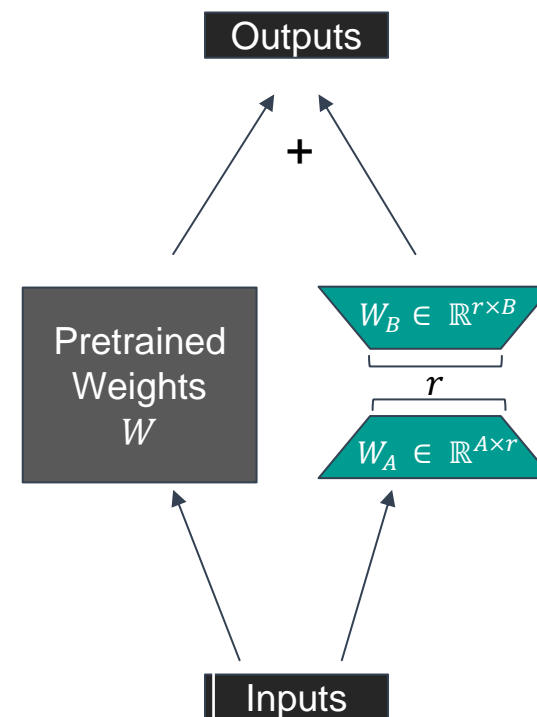


Parameter Efficient Finetuning

LoRA



Approximate ΔW
with $\Delta W = W_A W_B$



Number Parameters to train	$83666 * 32 * 2$	$= 5.354.624$
Optimizer states	$2 \times 2 \text{ bytes} \times 5M \text{ Parameters}$	$= 0.0214 \text{ GB}$
Gradients	$4 \text{ bytes} * 5M \text{ parameters}$	$= 0.0214 \text{ GB}$
Total		$14,54 \text{ GB}$

Finetuning Ablauf



- Mit Huggingface Transformers (baut auf pytorch auf)
- Objective: next token prediction
- Modell während des Trainings überprüfen

```
runs.summary["sample_predictions"]
```

	prompt	generation	max_length
1	You are a smart helpful	Step 1:	20
2	You are a smart helpful	Step 1:	20
3	You are a smart helpful	Step 1: Identify the key	20
4	You are a smart helpful	Step-by-	20
5	You are a smart helpful	Step 1:	20
6	You are a smart helpful	Step 1:	20
7	You are a smart helpful	Step 1:	20



Finetuning

Auf Dokumente

- Dokumente (SPO's, Modulhandbücher) in 140 Paragraphen aufgeteilt
- LLM übernimmt Schreibstil der Dokumente aber lernt kein neues Wissen

```
[
  {
    "paragraph": "Hochschule Konstanz Technik, Wirtschaft und Gestaltung Seite 1 von
                  43 Zulassungssatzung...",
  },
  ...,
  {
    "paragraph": "Hochschule Konstanz Technik, Wirtschaft und Gestaltung Seite 42 von
                  43 Zulassungssatzung... ",
  },
  {
    "paragraph": "Hochschule Konstanz Modulhandbuch des Studiengangs Fakultät Informatik ...",
  },
  ...,
]
```

Finetuning

Auf Dokumente

- LLM übernimmt Schreibstil der Dokumente aber lernt kein neues Wissen

Frage:

„Wer ist der Koordinator des Data-Science Modules an der Hochschule Konstanz?“

Referenz-Antwort:

„Prof. Dr. O. Dürr“

LLM-Antwort:

„Prof. Dr. M. Krause“

Finetuning

Instruction Finetuning auf RAG

- Trainieren auf 192 QAC-Triplets
(Question, Answer, Context)

```
[
{
  "question": "What are the goals of the 'Natural Language Processing' module?",
  "answer": "Students will gain a fundamental understanding of NLP, apply NLP methods on practical use cases...",
  "context": "Lernziele des Moduls\nFachliche Kompetenzen\nStudents will\n- gain a fundamental understanding of NLP and text processing,\n- apply NLP methods on practical use cases,\n- understand the mathematical concepts..."
},
...
{
  "question": "Was ist die Deadline für eine Bewerbung auf den Master Informatik im Sommersemester?",
  "answer": "1. Dezember des Jahres",
  "context": "Der Antrag auf Zulassung f\u00fcr das erste Fachsemester muss 2. f\u00fcr das Wintersemester bis zum 1. Dezember des Jahres"
},
...
]
```

[INST]You are a smart helpful assistant for the HTWG Konstanz.

Answer the following question based only on the provided context.

It is mandatory to answer in GERMAN:

Context: {context}

Question: {question}

[/INST]

{answer}

Finetuning

Instruction Finetuning auf RAG mit COT Antwort

- COT = Chain of Thought
- LLM geht die Aufgabe Schritt für Schritt an:
 1. Kontext betrachten und relevanten Abschnitt raussuchen + zitieren
 2. Wichtigen Informationen aus Kontextabschnitten extrahieren
 3. Antwort formulieren

```
[
  {
    "question": "Wie wird die Berufserfahrung im Auswahlprozess der Hochschule Konstanz bewertet?",
    "context": "\\n\\nFachspezifische Berufst\u00e4tigkeit nach Abschluss des Bachelorstudiums bis zum Zeitpunkt der Bewerbung: \\n- F\u00fcr mindestens neun Monate T\u00e4tigkeit bis zu einem Jahr: 0,1 Notenpunkte...",
    "cot_answer": "...",
  },
  ...
]
```

To answer the question, we need to look at the information provided in the context about the grading points that can be earned for professional activity after the completion of Bachelor's studies.

The context states:

##begin_quote##

\"Fachspezifische Berufst\u00e4tigkeit nach Abschluss des Bachelorstudiums bis zum Zeitpunkt der Bewerbung:

- F\u00fcr mindestens neun Monate T\u00e4tigkeit bis zu einem Jahr: 0,1 Notenpunkte

- F\u00fcr mindestens ein Jahr T\u00e4tigkeit bis zu zwei Jahren: 0,2 Notenpunkte

- F\u00fcr mindestens zwei Jahre T\u00e4tigkeit: 0,4 Notenpunkte. - Einer Berufst\u00e4tigkeit gleichwertige Fortbildungen k\u00f6nnen einmalig mit 0,1 Notenpunkte aufgewertet werden. Somit sind insgesamt maximal 0,5

Notenpunkte m\u00f6glich.\"

##end_quote##

This translates to:

- For at least nine months of activity up to one year: 0.1 grading points
- For at least one year of activity up to two years: 0.2 grading points
- For at least two years of activity: 0.4 grading points

Equivalent professional training can be upgraded once with 0.1 grading points. Therefore, a maximum of 0.5 grading points is possible in total.

<ANSWER>: A maximum of 0.5 university grading points can be earned for specific professional activity after the completion of the Bachelor's studies.",

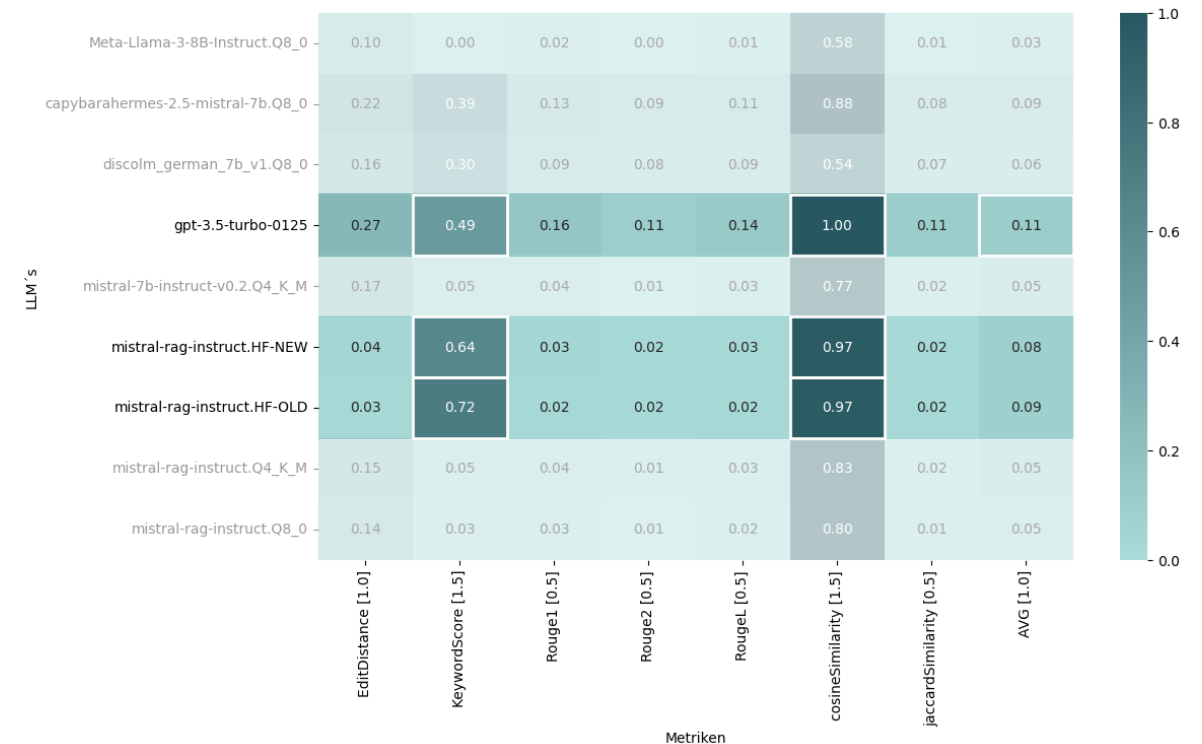
Finetuning

Auswertung - Gewichtet und Ungewichtet

Gewichtete Similarity der LLM's mit Metriken



Gewichtete Similarity der LLM's mit Metriken (unter Berücksichtigung von Specificity und Relevance)

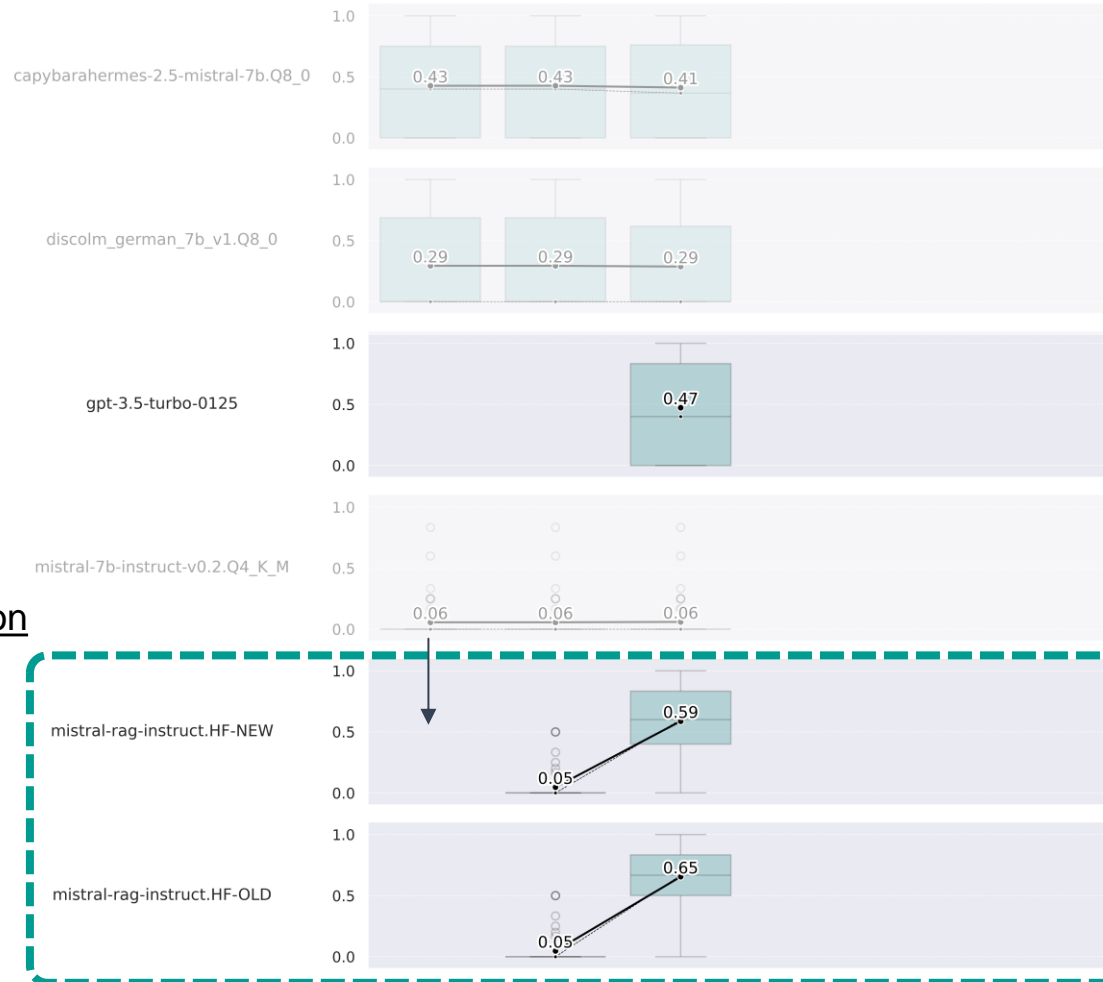


- **gefinetunete Modelle** kommen an GPT3 ran (HF-NEW & HF-OLD)
- sind in Metriken über Textähnlichkeit deutlich schlechter durch Chain of Thought

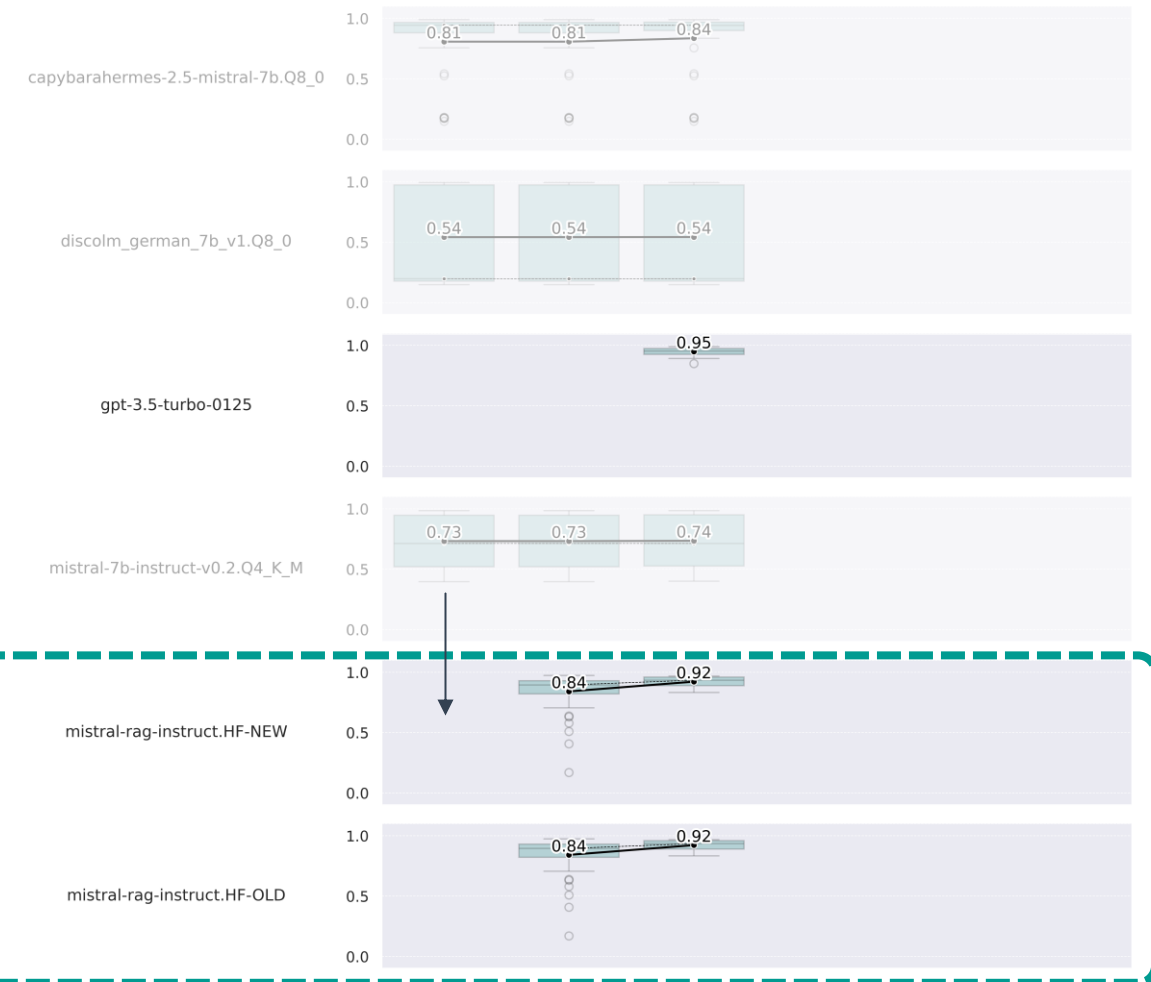
Finetuning

Auswertung – QAC Finetuning vs QAC + COT

Keyword Score



Cosine Similarity

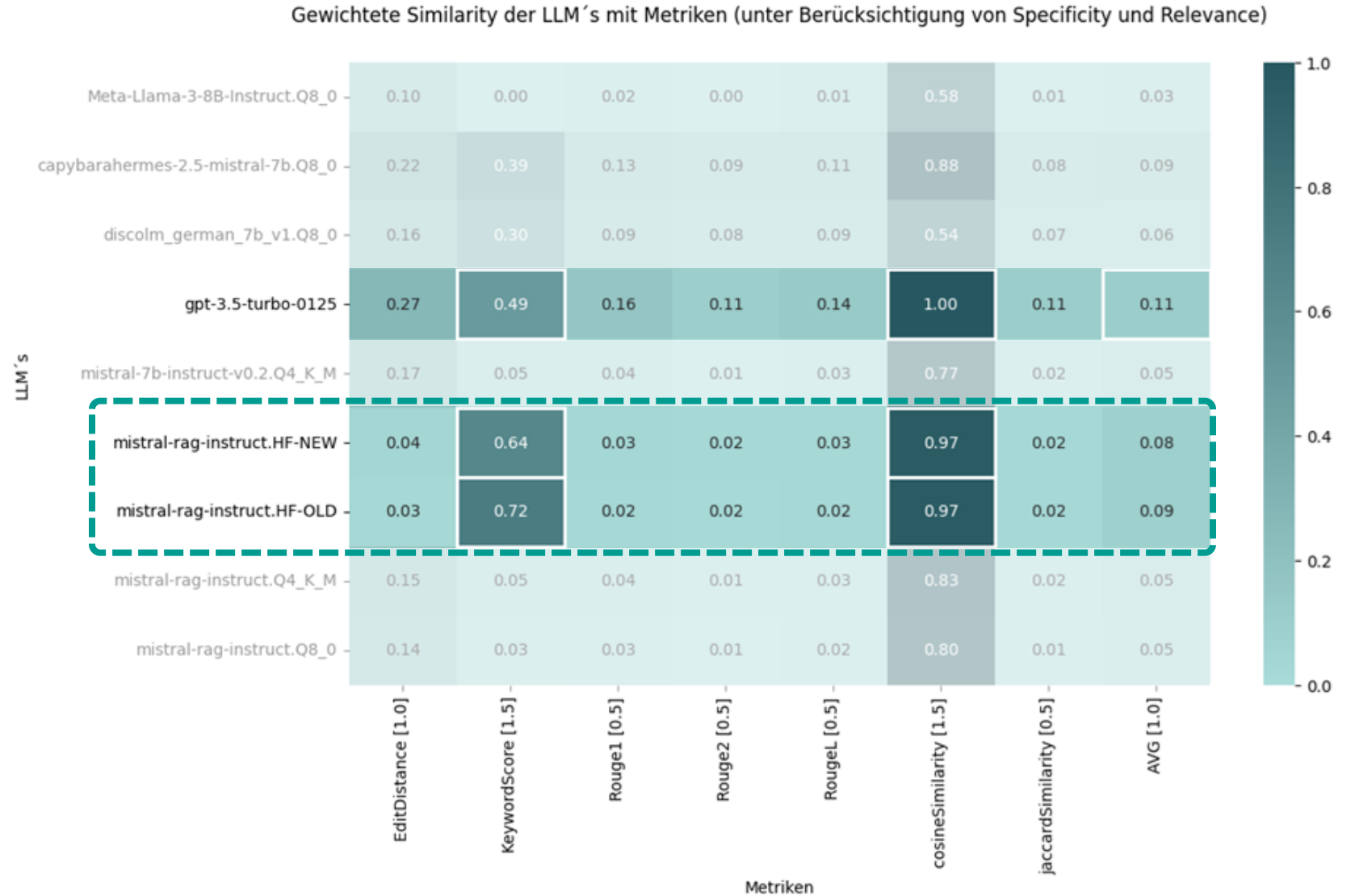


Finetuning

Auswertung – Ungewichtet und Gewichtet

→ **gefinetunete Modelle**
kommen an GPT3 ran
(HF-NEW & HF-OLD)

→ sind in Metriken über
Textähnlichkeit deutlich
schlechter durch Chain of
Thought



KADDI

Relevanz Prompt-Engineering

- **Prompt:** "Bitte alle Fragen in Deutsch beantworten und nur mithilfe des gegebenen Kontexts beantworten. {context} Frage: {question}„
 - **Resultat:** Einige Dutzend Antworten werden dennoch auf Englisch beantwortet und teilweise werden sogar Fragen erfunden, die nicht gestellt wurden.
- **Prompt:** "Please answer only in german and only with the knowledge of the given context, else answer 'Ich weiß es nicht'. Do not make up any information or question. {context} Frage: {question}„
 - **Resultat:** Eine Handvoll Antworten werden dennoch auf Englisch beantwortet und es werden weiterhin zusätzliche Fragen erfunden.
- **Prompt:** "ANSWER ONLY IN GERMAN AND ONLY TO ONE QUESTION, THAT YOU RECEIVE WITH THE KNOWLEDGE OF THE GIVEN CONTEXT. IF YOU CANT ANSWER BASED ON THE GIVEN CONTEXT WRITE 'Ich weiß es nicht'. BITTE NUR DIE ERSTE FRAGE BEANTWORTEN UND KEINE WEITEREN. {context} Frage: {question}„
 - **Resultat:** Alle Fragen auf Deutsch beantwortet, nur die gestellte Frage wird beantwortet, jedoch wird diese wenn zu unpräzise gestellt, dennoch zu Ende erfunden und beantwortet. Bisher immer noch nicht die Antwort "Ich weiß es nicht" zurückgegeben.

Empfehlung und Demo

- Die Entwicklung des Chatbots für das Prüfungsamt ist machbar ist, allerdings bleibt ein erhebliches Maß an Unsicherheit bestehen, da große Sprachmodelle (LLMs) dazu neigen, falsche oder erfundene Informationen zu generieren, wenn sie nicht über ausreichende Kenntnisse verfügen.
- Dieses Problem besteht sogar bei fortschrittlichen LLMs wie ChatGPT, trotz der erheblichen finanziellen Investitionen in ihre Entwicklung, so dass unser lokales Modell wahrscheinlich nicht mit ihrer Leistung mithalten kann.
- In Anbetracht der Wichtigkeit genauer Informationen im Kontext eines Prüfungsamtes empfehlen wir, den Chatbot mit Vorsicht einzusetzen und sicherzustellen, dass er immer **Verweise auf die Quelldokumente und spezifische Seitenzahlen liefert.**



Empfehlung

Kostenanalyse

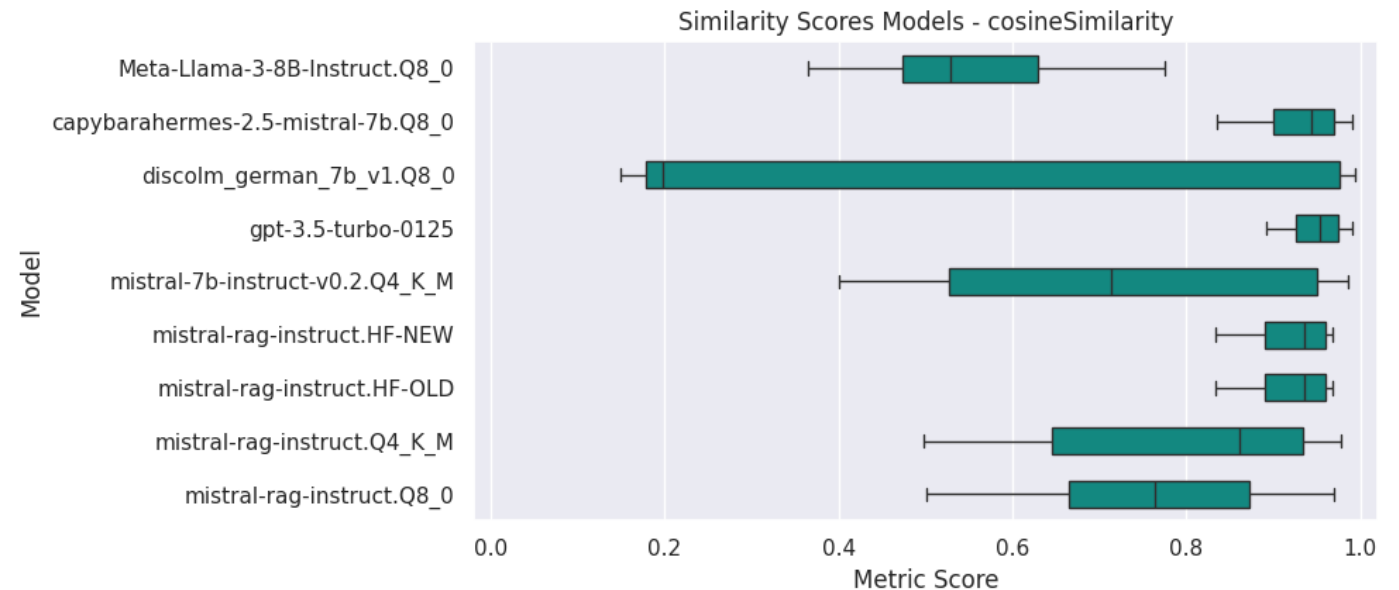
- DB müsste lokal laufen (keine großen Anforderungen)

- Annahme Jeder Student (800) nutzt die Funktion einmal pro Semester

Mistral (Small)	7.79\$
GPT 3.5 turbo 0125	4.03\$
Llama 3	2.62\$
GPT4o-latest	40.34\$

- LLM würden wir „GPT 3.5 Turbo“ empfehlen (Kosten / Nutzen)

- Lokales Ausführen erfordert teure Hardware im Bereich von mehreren tausend Euro
- Lokale Skalierung ist nicht trivial



Empfehlung

- Machbarkeit gegeben, aber hohe Unsicherheit
- LLMs neigen zu fehlerhaften/erfundenen Infos
- Verwaltung lokaler Modelle ist deutlich aufwendiger als Cloud-Dienste (wie z.B. Azure/OpenAI)
- **Wichtige Maßnahmen:**
 - Vorsichtiger Einsatz des Chatbots
 - Verweise auf Quelldokumente und spezifische Seitenzahlen sicherstellen



Ausblick

Relevanz Prompt-Engineering

Prompt	Resultat
„Bitte alle Fragen in Deutsch beantworten und nur mithilfe des gegebenen Kontexts beantworten. {context} Frage: {question}“	<ul style="list-style-type: none">• Einige Dutzend Antworten auf Englisch• Teilweise werden Fragen erfunden
„Please answer only in german and only with the knowledge of the given context, else answer "Ich weiß es nicht". Do not make up any information or question. {context} Frage: {question}“	<ul style="list-style-type: none">• Einige Antworten auf Englisch• Zusätzliche Fragen werden erfunden
„ANSWER ONLY IN GERMAN AND ONLY TO ONE QUESTION, THAT YOU RECEIVE WITH THE KNOWLEDGE OF THE GIVEN CONTEXT. IF YOU CANT ANSWER BASED ON THE GIVEN CONTEXT WRITE 'Ich weiß es nicht'. BITTE NUR DIE ERSTE FRAGE BEANTWORTEN UND KEINE WEITEREN. {context} Frage: {question}“	<ul style="list-style-type: none">• Alle Fragen auf Deutsch beantwortet• Nur die gestellte Frage wird beantwortet• Unpräzise Fragen werden gegebenenfalls ergänzt und beantwortet• Keine Antwort mit "Ich weiß es nicht" gegeben

Conversational Prompt Injection



Fortführende Themen

- Performance Optimierung Lokale LLM's
 - Momentan kann nur eine Anfrage gleichzeitig beantwortet werden
- Kontextuelle Konversation
 - Frühere Fragen und Antworten mit einbeziehen (momentan wird nur auf die aktuelle Frage reagiert)
- Vergleich Full Fine Tuning mit LoRA
- „Guardrails“ für LLM's -> Input und Output Validierung

Demo

JAN

