

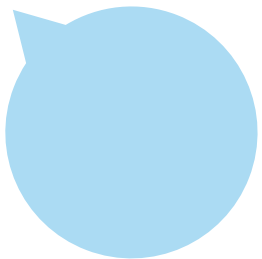
Methodische Evaluierung von Sprachmodellen

Modelle für die Öffentliche Verwaltung Evaluieren

René Walter, Dr. Thilo Michael

Agenda

1. Übersicht KI-Kompetenzcenter
2. Was ist MÖVE?
3. Was wird evaluiert?
4. Ergebnisse
5. Ausblick



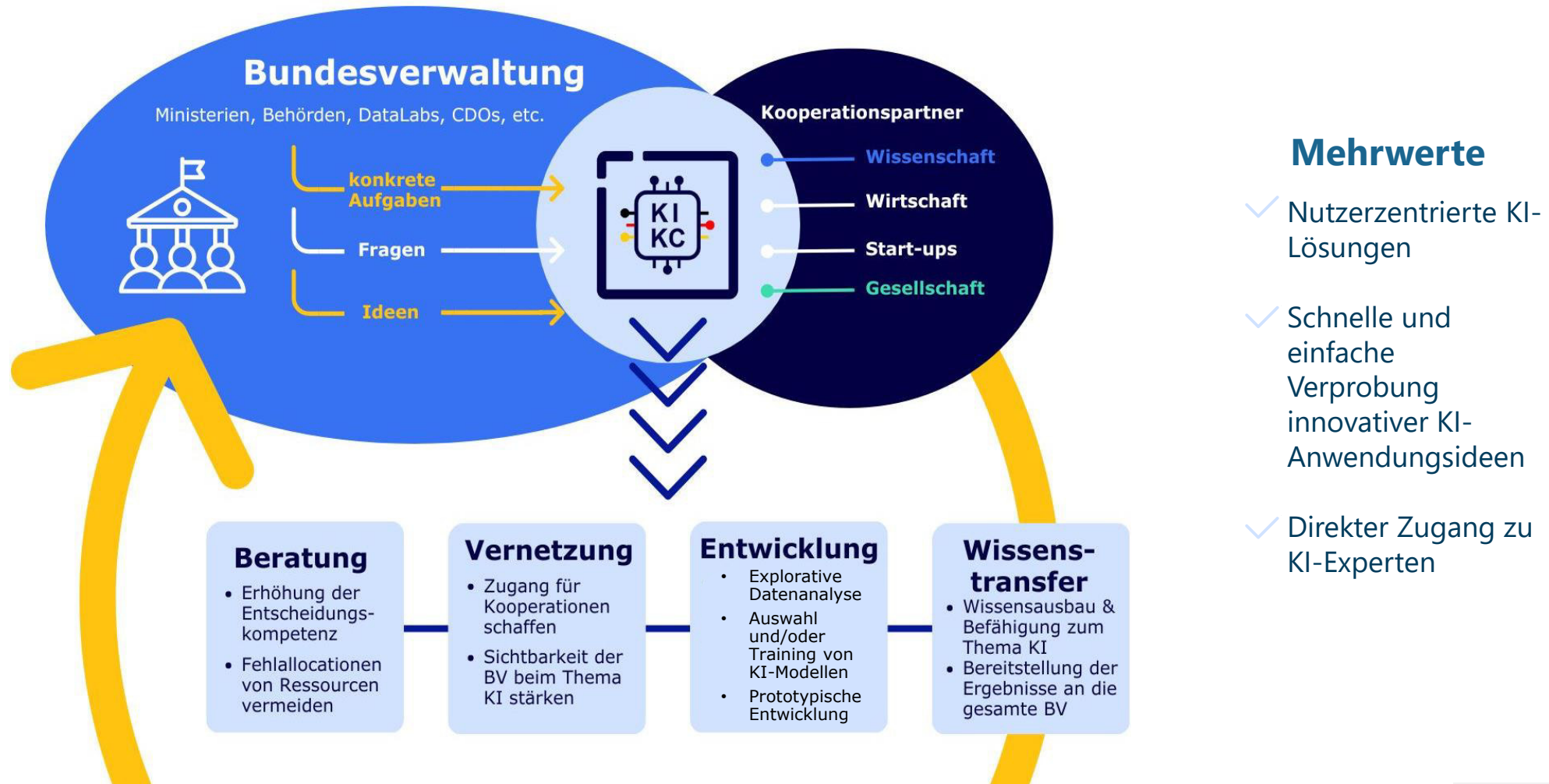


1. KI-Kompetenzcenter

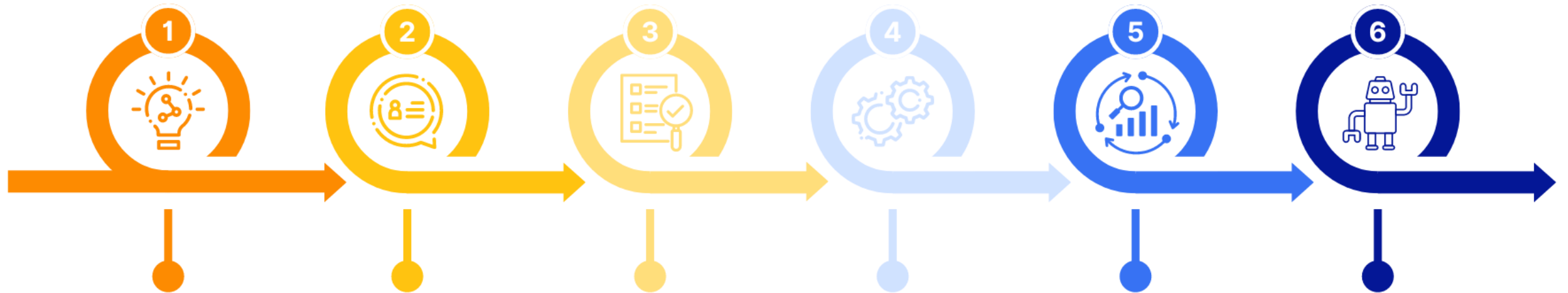
Übersicht und Vorgehen



Befähigung der Bundesverwaltung zum souveränen Einsatz von KI



Iterativ von der Idee zum KI-Prototyp



Use Cases identifizieren & Machbarkeit bewerten

- Definition der Aufgabe, die mit KI gelöst werden soll
- Bewertung der Datenlage

Mitarbeitende der Verwaltung einbinden

- Analyse von Arbeitskontext und Rahmenbedingungen für einen KI-Einsatz
- Verstehen der Anforderungen der Nutzenden

Analyse der verfügbaren Daten

- Bereinigen der Daten
- Explorative Datenanalyse

Auswahl und/oder Entwicklung eines KI-Modells

- Auswahl geeigneter Modelle
- Trainieren von Modellen

Evaluation von Modellen und Mehrwerten

- Bewerten und auswählen des KI-Modells mit Zielnutzenden
- Iterative Entwicklung einer prototypischen Anwendung

Projekt- und Prototyp-Übergabe

- KI-Modell für Nutzenden-gruppen implementieren
- Prototyp und fachliche Erkenntnisse bereitstellen

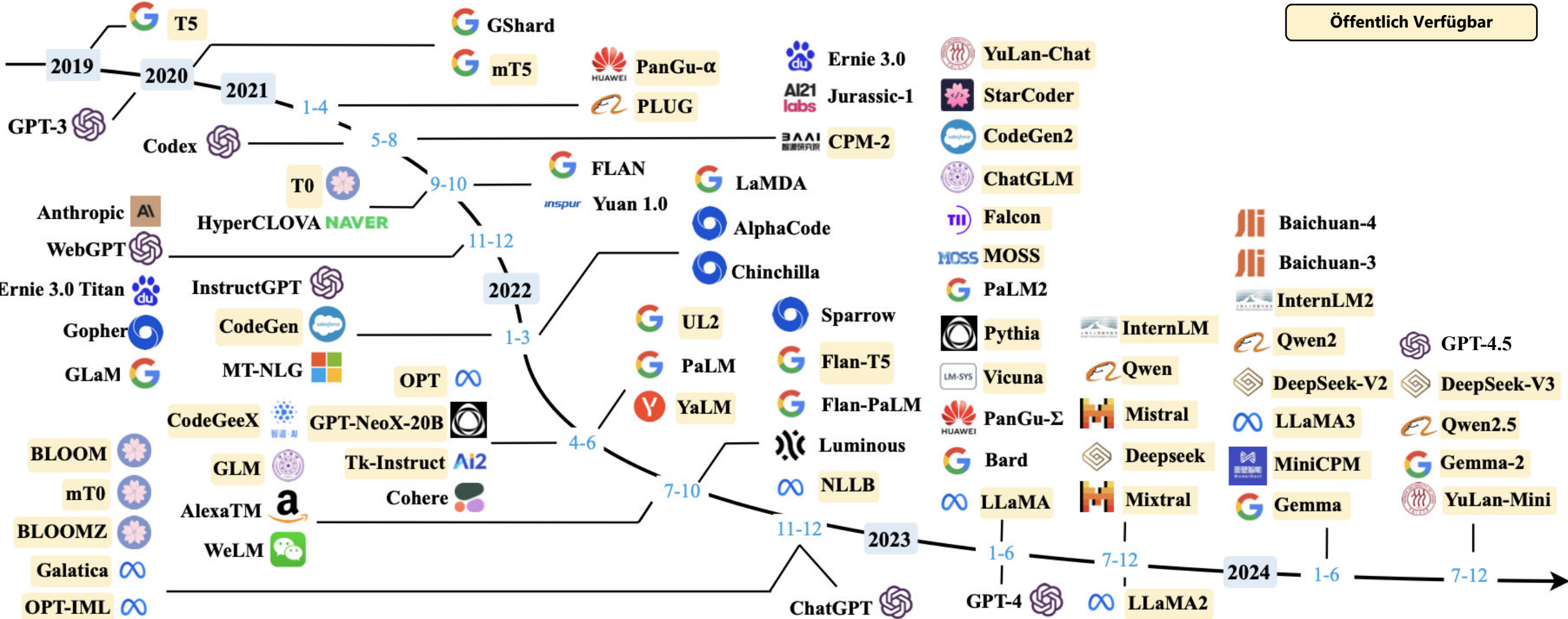


2. Was ist MÖVE?

Übersicht des Forschungsvorhabens



Sprachmodelle: Nicht nur ChatGPT



Was ist MÖVE

- Ein Sprachmodell-Benchmark für die öffentliche Verwaltung
- Forschungsprojekt im Rahmen des KI-Kompetenzcenters
- Aktive Entwicklung seit ca. 6 Monaten

Ziel

Der Benchmark soll ein frei verfügbares, nützliches Werkzeug für den Einsatz von Sprachmodellen in der öffentlichen Verwaltung sein

Es gibt schon viele LLM-Benchmarks!

Wieso noch einer?

- **Fokus auf relevanten Aufgaben für die Bundesverwaltung**

Andere: Puzzle, Generelles Wissen, "AGI"

MÖVE: Relevante Aufgaben, die in der öffentlichen Verwaltung Anwendung finden

- **Performance und Governance Kriterien**

Andere: Fokus auf Performance in einer spezifischen Aufgabe

MÖVE: beinhaltet "Governance Kriterien" wie Bias, Sicherheit, Halluzinationen

- **Deutsche Sprache**

Andere: Häufig in englische oder automatisch übersetzte Datensätze

MÖVE: speziell ausgewählte, teils händisch erstellte Datensätze, meistens direkt auf Deutsch

- **Validierte Datensätze**

Andere: Sehr große Datensätze mit teils schlechter Qualität

MÖVE: Goldstandard Datensätze mit Domänenwissen direkt aus der Verwaltung

Übersicht der Performance und Governance Kriterien

PERFORMANCE Evaluierung



GOVERNANCE Evaluierung



Übersicht der Performance und Governance Kriterien

Eine Performance oder Governance Evaluierung besteht aus mehreren "**Szenarien**":

Ein Szenario wird definiert durch:

- Ein LLM (z.B. GPT-4o, mistral-large-2411, ...)
- Eine Aufgabe (z.B. Zusammenfassen, Question Answering, ...)
- Einen Datensatz
- Einen Prompt (z.B. "*Fasse das folgende Gesetz zusammen: {Gesetzestext}*")
- Mehrere Metriken ("Klassische" Metriken, ML-basierte Metriken, LLM-as-a-judge Metriken)



3. Was wird evaluiert?

Modelle, Datensätze, Metriken



LLMs: Welche Modelle werden evaluiert

Wir evaluieren über 32 Modelle (hauptsächlich **Open-Weights Modelle** und OpenAI's **GPT Modelle**)

llama-3.3-70b-instruct eurollm-1.7b-instruct-q8
phi3.5-3.8b openelm-3b-instruct-q4_0
luminous-base-control-20240215 mistral-large-2411
mistral-nemo-12b llama-3.1-70b-instruct
llama3.2-1b gpt-35-turbo-16k gpt-4o-mini
smollm-1.7b gemma2-2b
gpt-4-128k mixtral-8x22b gpt-4o
mistral-7b sauerkrautlm-nemo-12b-instruct-q4_0
llama3.2-3b smollm2-1.7b
mixtral-8x7b sauerkrautlm-phi-3-medium-q4_0
gemma2-9b llama3.1-8b
mistral-large-2407 luminous-extended-control-20240215

Performance-Evaluierung

Aufgaben und Datensätze

Aktuell evaluieren wir 2 Aufgaben mit insgesamt 6 Datensätzen:

ZUSAMMENFASSEN

- **Eur Lex Sum**
Zusammenfassen von europäischen Gesetzen
- **Swiss Leading Decisions**
Zusammenfassen von Schweizer Gerichtsentscheidungen
- **KIKC-Summary**
Händisch erstellte Zusammenfassungen von Dokumenten aus verschiedenen Ministerien

BEANTWORTUNG VON FRAGEN

- **German QuAD**
Fragen aus der deutschen Wikipedia
- **KIKC-QA**
Händisch erstellte Fragen und Antworten aus Dokumenten aus verschiedenen Ministerien
- **FAQ Law**
FAQ zu mehreren deutschen Gesetzen

Woran wir gerade arbeiten...

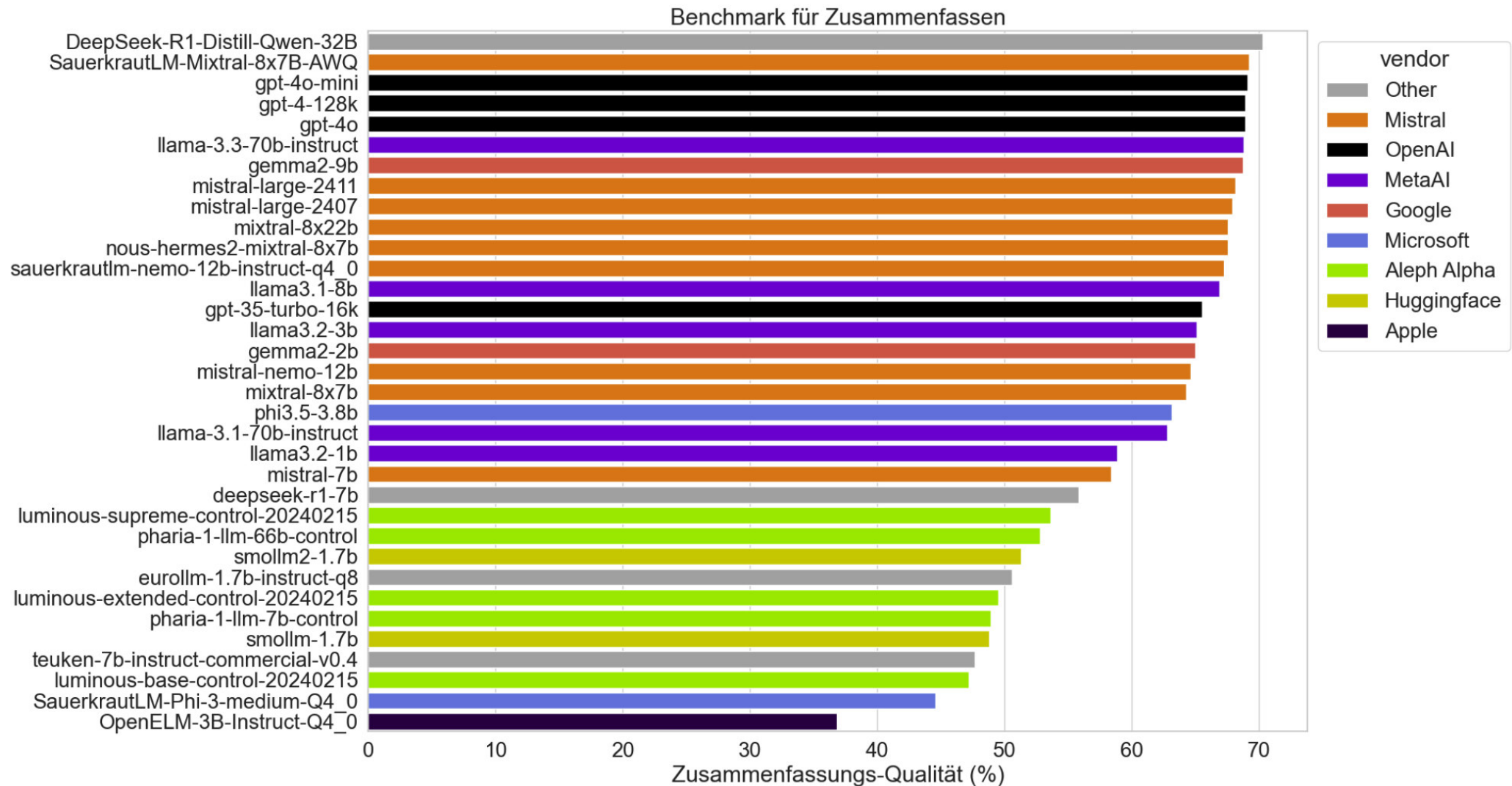
KIKC-Keywords (händisch erstellt), **SUM_Q1** (Veröffentlichungen des Umweltbundesamt), **SUM_Q4** (DIP Veröffentlichungen), **SUM_Q10** (social science open access repository), Wahlomat (Datensatz zur deutschen Politik)

Unsere Metriken

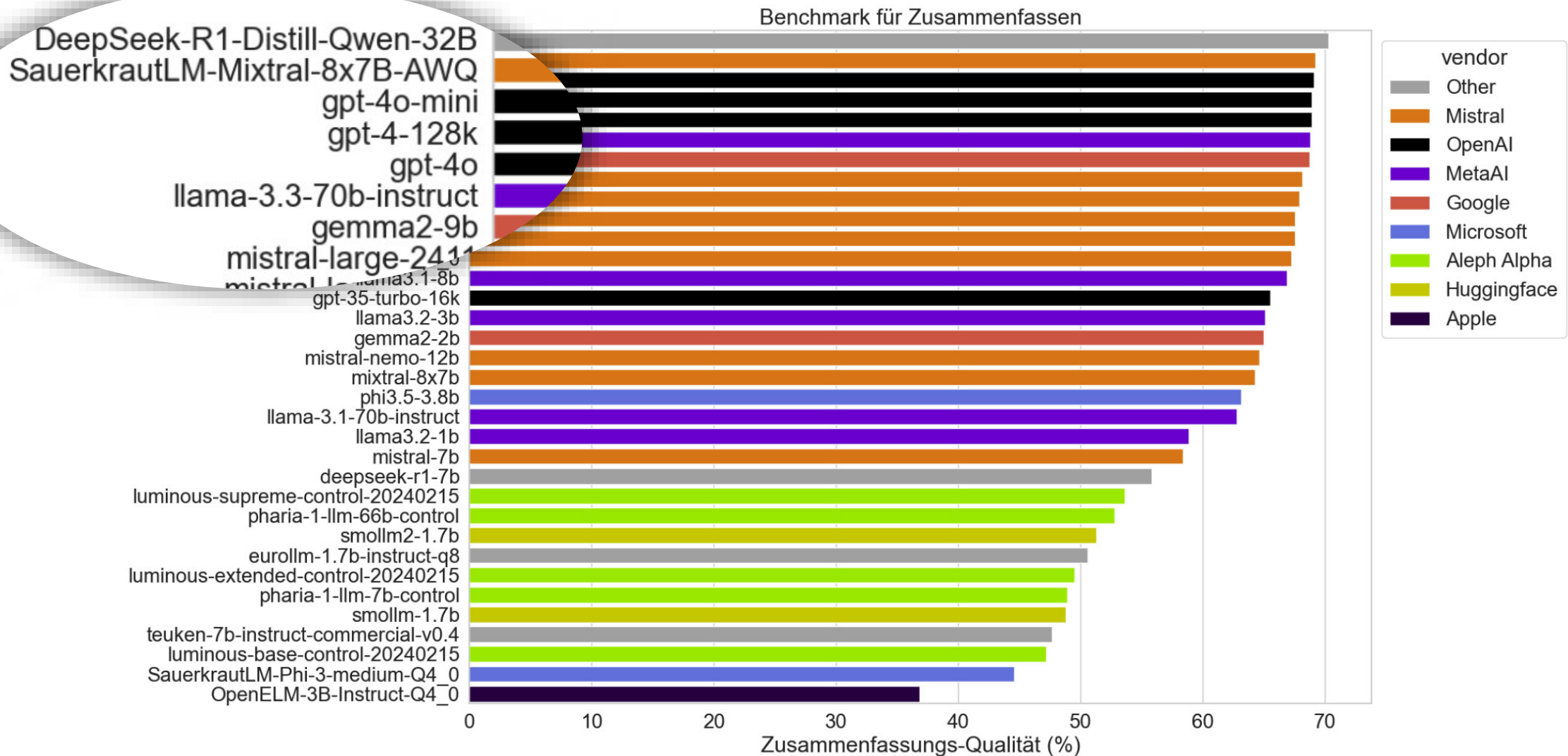
- Die **gewählten Metriken** messen **unterschiedliche Facetten** der Modelle
- Wir untersuchen wie sich die **Metriken untereinander unterscheiden**
- Wir wählen aus, **welche Metriken repräsentativ für eine Aufgabe** sind

exactmatch:case_insensitive_match:median
bleu:precision:mean bertscore:recall:median
bertscore:recall:mean rouge:rouge2:median
bertscore:precision:median
bleu:bleu:mean rouge:rouge1:mean
bertscore:precision:mean
bertscore:f1:median
exactmatch:exact_match:mean rouge:rouge2:mean
exactmatch:exact_match:median bertscore:f1:mean
exactmatch:fuzzy_match:mean bleu:bleu:median
rouge:rouge1:median bleu:precision:median
exactmatch:case_insensitive_match:mean

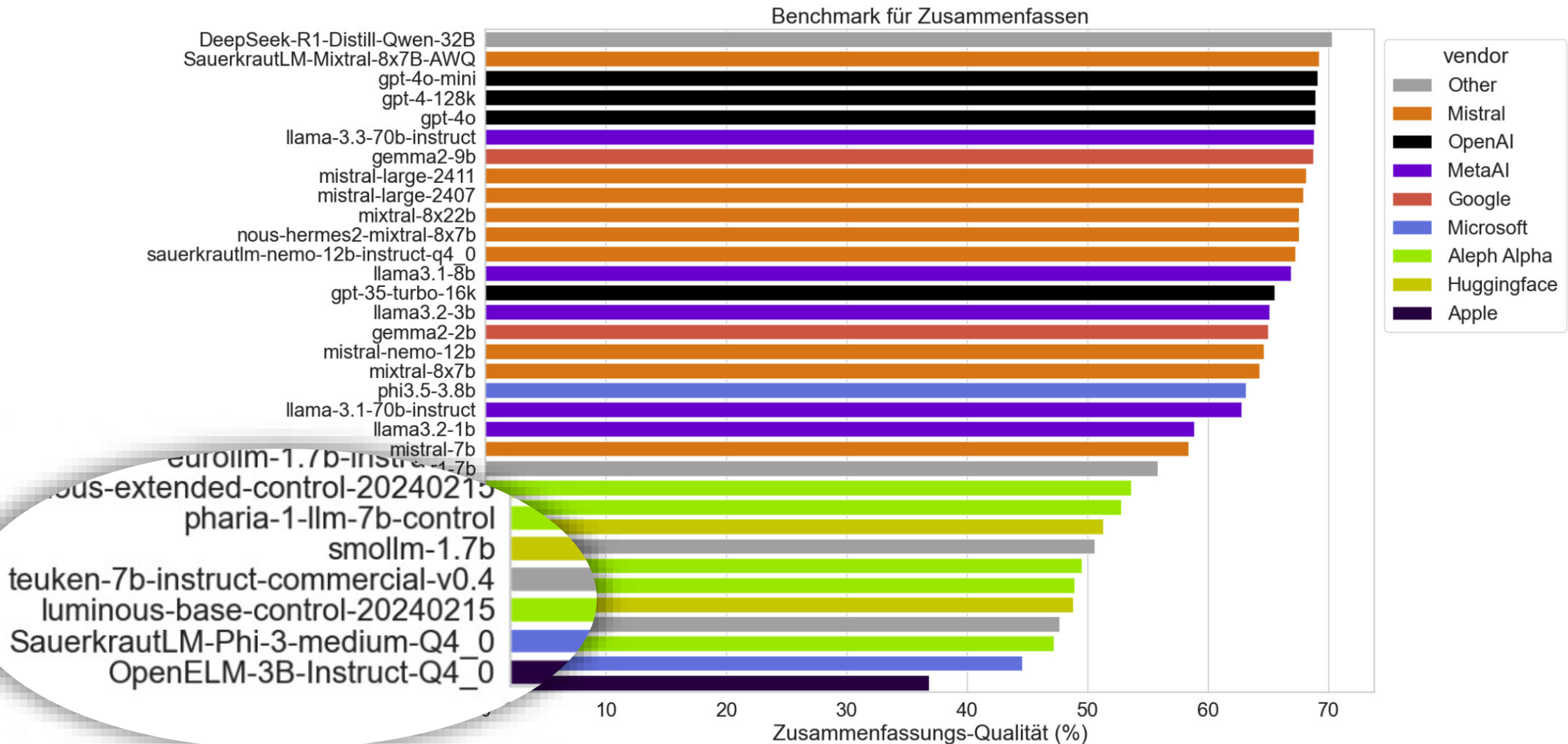
Ergebnisse: Zusammenfassen



Ergebnisse: Zusammenfassen

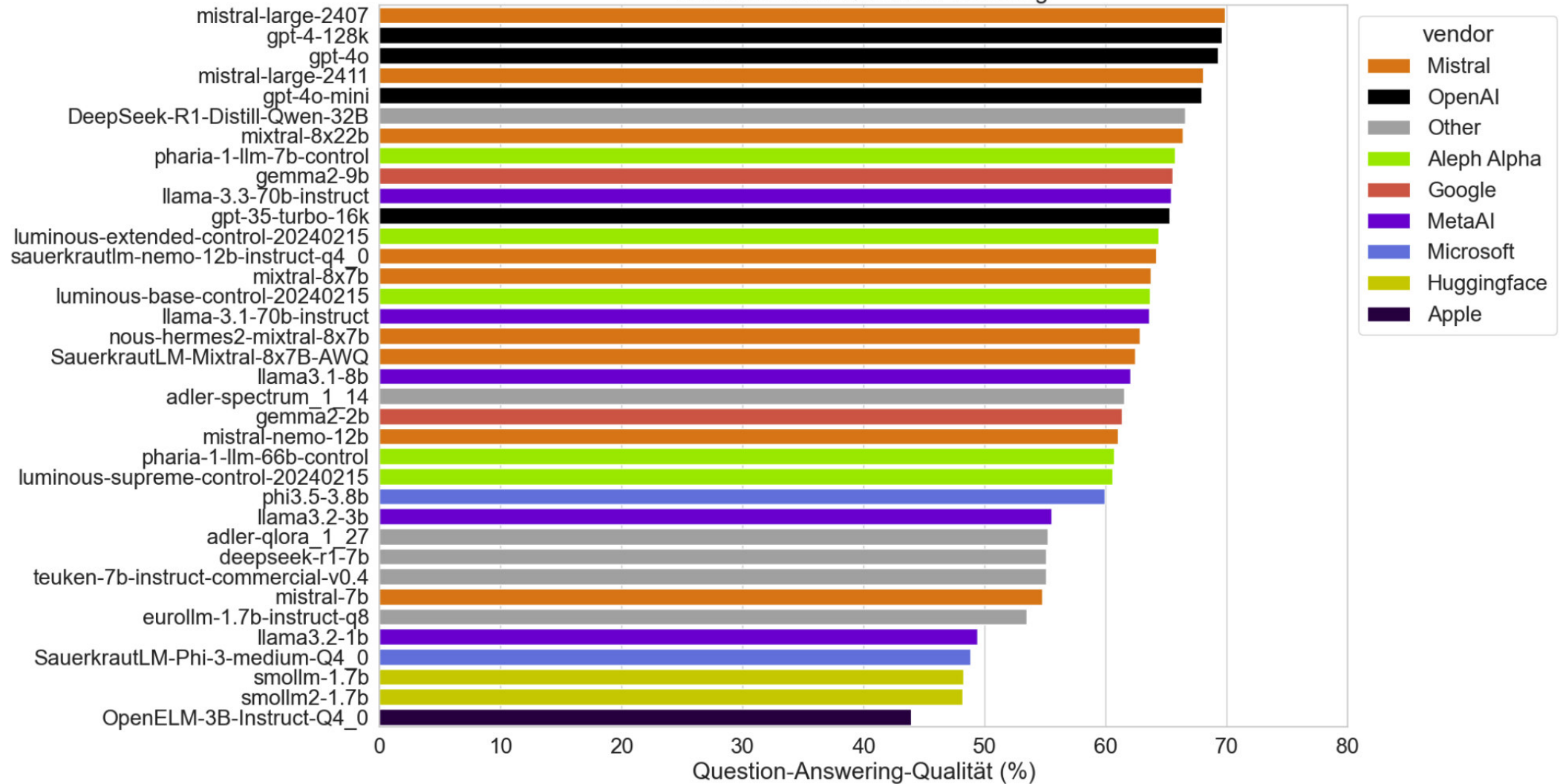


Ergebnisse: Zusammenfassen



Ergebnisse: Question Answering

Benchmark für Question Answering



Governance-Kriterien

Soziale Fairness

- Biases (z.B. Misogynie)
- Daten
 - En: Fairness (DecodingTrust), Microaggressions (Social Bias Inference Corpus)
 - De: German NLG Bias

Politischer Bias

- Extrahieren von politischen Ansichten aus den Modellen
- Sehr relevant u.a. wegen der politischen Entwicklungen in den USA
- Daten
 - GermanPartiesQA, Political Compass Test

Andere Governance Kriterien

- Potentielle Kooperationen (TÜV AI.Lab, TU/DFKI)

Performance-Kriterien

Keyword extraction

- Händisch erstellter Goldstandard-Datensatz verfügbar
- Große Nachfrage bei Ministerien

Übersetzen

- Finden von öffentlich verfügbaren, ministeriumsnahen Datensätzen
- Fokus: Deutsch <-> Englisch
- Austausch mit dem BAMF zu weiteren Sprachen

Reasoning Modelle

- Evaluierung von Reasoning Modellen
 - OpenAI (o1, o1-mini, o3, o3-mini)
 - Deepseek (R1 & distills)
 - Alibaba Qwen (QwQ)

Mehr Interesse an KI?

KI-Report



DER KI-REPORT DES KI-KOMPETENZCENTERS

Hallo in der schönen neuen KI-Welt

Serien, Mobilität, Texte: Künstliche Intelligenz verändert viele Lebensbereiche. Wir wollen die Potenziale von KI nutzen, um auch der Öffentlichkeit Verwaltung das Leben zu erleichtern. Immer im Fokus: die Nutzenden. Dafür schauen wir ganz genau hin – gemeinsam mit den Mitarbeitenden der ÖV: Für welche Anwendungsfälle eignet sich KI? An welcher Stelle hilft sie, Routineaufgaben effizienter zu bewältigen und Mitarbeitende zu entlasten? Daraus entwickeln wir KI-Prototypen, erproben und evaluieren sie in der Praxis.



MAXIM SCHNJAKIN
Projektleiter des KIKC

Halluzinationen – was sind das?

KI-Halluzinationen sind Antworten von Sprachmodellen, die faktisch nicht korrekt sind. Sie "halluzinieren" sozusagen. Das geschieht deshalb, weil das Modell darauf trainiert wird möglichst "passende" Texte zu erstellen – z.B. Texte die von Menschen präferiert werden – aber nicht faktisch korrekte Texte.



WIE KOMMT ES ZU HALLUZINATIONEN?

Beim "Fine-Tuning" (dem Training von Verhalten des Modells) muss das Verhalten eintrainiert werden: nur faktische Informationen bzw. in einer Wissensquelle enthaltenen Informationen sollen ausgegeben werden. Die Qualität und Menge der Daten für dieses Training ist entscheidend für die Anfälligkeit zu halluzinieren.

Bedarfsanalyse

Anwendungsfälle identifizieren



Lernen

Aufbauseminar KI-Kompetenz

3 Tage

Praxisnahes Seminar für
Bundesbehörden zur
Stärkung von
KI-Kompetenz





Danke

**für Ihre
Aufmerksamkeit**