



Universität Regensburg
UNIVERSITÄTSBIBLIOTHEK

Repositorien und KI
Dr. Gernot Deinzer

Vorüberlegungen

KI – LARGE LANGUAGE MODELLS (LLMS)

Funktionsweise LLM

Zerlegung des Textes in Bausteine (Tokens)

- Wörter oder Silben

Einbettung

- Zuordnung von Wort zu Vektor
- Neuronalen Netzen (Transformator Modell)

Vorhersage des nächsten Tokens

- Neuronale Netze
- Erkennung von Eingaben
- Berechnung von Wahrscheinlichkeiten
- Voraussetzung Training

Dekodierung

- Auswahl eines Tokens
- Wahrscheinlichkeitsberechnung

Zerlegung eines Textes in Bausteine

Repositorien spielen eine wichtig Rolle in der KI-Landschaft.

Erstellt mit <https://tiktokenizer.vercel.app/>

Funktionsweise LLM

Zerlegung des Textes in Bausteine (Tokens)

- Wörter oder Silben

Einbettung

- Zuordnung von Wort zu Vektor
- Neuronalen Netzen (Transformator Modell)

Vorhersage des nächsten Tokens

- Neuronale Netze
- Erkennung von Eingaben
- Berechnung von Wahrscheinlichkeiten
- Voraussetzung Training

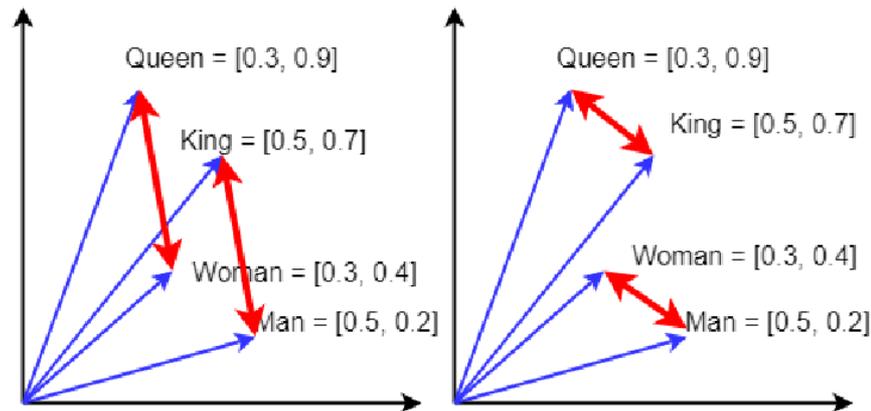
Dekodierung

- Auswahl eines Tokens
- Wahrscheinlichkeitsberechnung

Funktionsweise LLM

Zuordnung Wort – Vektor

Wichtig Beinhaltung semantischer Strukturen



Funktionsweise LLM

Zerlegung des Textes in Bausteine (Tokens)

- Wörter oder Silben

Einbettung

- Zuordnung von Wort zu Vektor
- Neuronalen Netzen (Transformator Modell)

Vorhersage des nächsten Tokens

- Neuronale Netze
- Erkennung von Eingaben
- Berechnung von Wahrscheinlichkeiten
- Voraussetzung Training

Dekodierung

- Auswahl eines Tokens
- Wahrscheinlichkeitsberechnung

Problematik

Halluzination

Realistisch wirkende Information

Abweichung von Input

Fehlende Übereinstimmung (faithfulness)
Mangelnde Richtigkeit (factualness)

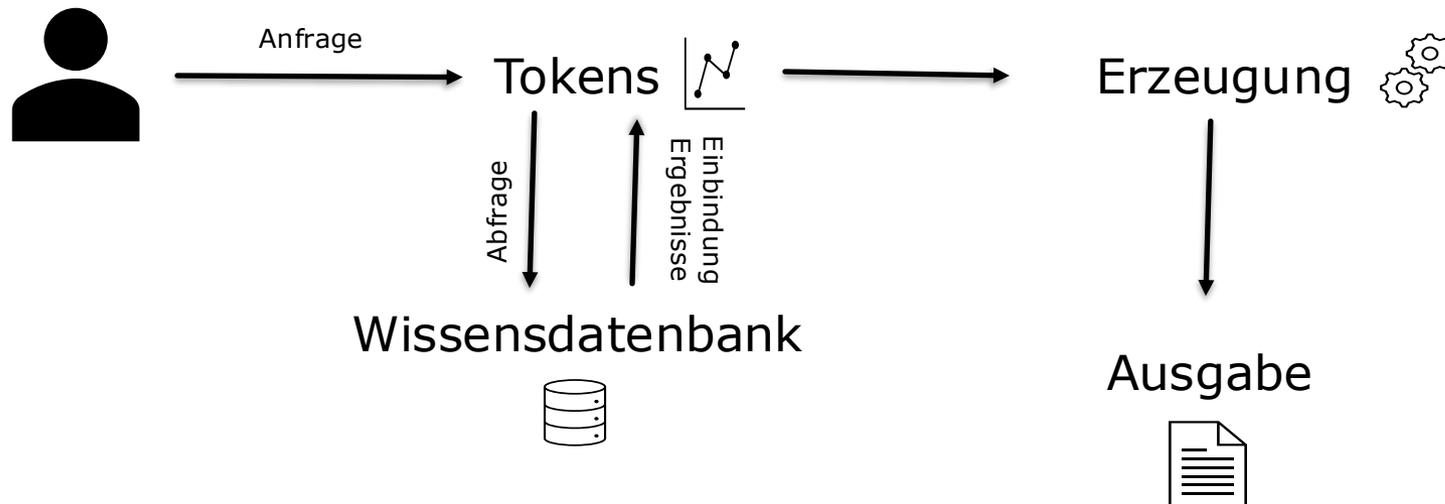
Erkennen:

Unsicherheit messen

Mit Faktendatenbanken vergleichen

Retrieval Augmented Generation (RAG)

Ergänzung LLMs mit einer Wissensdatenbank





Universität Regensburg

Dr. Gernot Deinzer

Leitung Abteilung IT- und Publikationsdienste
Universitätsbibliothek Regensburg

Datenquelle für KI

REPOSITORIEN

Repositorien - LLMs

Repositorien

Qualitätsgesicherte Inhalte im Internet

Dauerhafte Texte

Ergebnisse zu verbessern und richtige Ergebnisse
anzuzeigen

**Ziel: Repositorieninhalte in LLMs zu
berücksichtigen**

Zusätzlich:

Möglichkeit einer Wissensdatenbank

Inhalte besser erzeugen und kontrollieren

Halluzinationen vermindern

Policies

Leitlinien für die Nutzung Inhalte und Metadaten

- **Maschinenlesbar**
- Internationalität

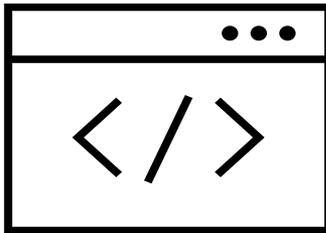
Metadaten frei und nachnutzbar (CC-0)

Barcelona Declaration on Open Research

Information (<https://barcelona-declaration.org/>)

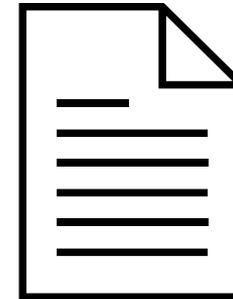
(9 Unterzeichner Deutschland, 1 Bayern)

Präsentation



Landing Page

Alle Information zum
Dokumen (Metadaten)



Dokument (pdf)

Volltext des Beitrags

**Problem: Verlinkung nur von Landing Page zu
Dokument, nicht beidseitig**

Dateiformate- pdf

Standardformat

Gut menschenlesbar

Maschinenlesbarkeit ?

Fehlende explizite Struktur

Komplexe und fehleranfällige Extraktion

Wichtig

Text hinterlegt

Metadaten in Eigenschaften ausgefüllt

Dateiformate - HTML

Natives Format für Internet

Strukturierte und semantische Information

Semantische Tags, Listen, Überschriften

Trennung von Inhalt und Layout

CSS getrennt von html

Direkte Verarbeitung

HTML-Dokument bekannt aus Pretraining

Effizientere Extraktion relevanter Inhalte

z. B. Entfernung Navigation, Werbung

Bessere Ergebnisse in RAG-Systemen

Aufgrund semantischer Tiefe



Von Krauss - Eigenes Werk, CC BY-SA 3.0,
<https://commons.wikimedia.org/w/index.php?curid=32998484>

Dateiformate – JATS

Präzise und standardisierte Struktur:

Hohe Interoperabilität und Wiederverwendbarkeit

Metadaten im Format strukturiert

Maschinenlesbarkeit und Suchfähigkeit

Erfolgreicher Einsatz mit LLMs

Potenzielle Probleme

Komplexität

Tokelimitierungen bei LLMs

Metadaten

Siehe Vortrag

Frech Andreas, KI &
Metadaten - Generierung,
Extraktion und Bereitstellung
(Andreas Frech, UB der LMU)

Wichtig:
Maschinelle Repräsentation

The screenshot shows a web interface for exporting bibliographical data. At the top, there is a dark red button labeled "Export bibliographical data". Below it, a dropdown menu is open, showing a list of export formats. The current selection is "ASCII Citation". The list includes: Data Cite XML, DataCite XML OAI, Dublin Core, Dublin Core (LZV), EndNote, HTML Citation, JSON, JSON LD (highlighted with a red box), METS, MODS, OPENAIRE, OpenAPC, RDF+N-Triples (highlighted with a red box), RDF+N3 (highlighted with a red box), RDF+XML (highlighted with a red box), Refer, and Reference Manager. To the right of the dropdown, there is an "Export" button. The background of the page shows a snippet of a scientific article with text about miRNA and CCR4-NOT dead endonucleases.

Suchmaschinenoptimierung

Verändertes Suchverhalten

Fragen an KI-gestützte Systeme – Antworten ohne Klick (Zero-Click ergebnisse)

Neue Anforderungen an Inhalte

Nicht ausschließlich Keywords
vermehrt semantischer Relevanz und Kontext

-> Hochwertiger, tiefgehender, nutzerzentrierter Content wird wichtiger als reine Keyword-Dichte

LLM-SEO als neue Disziplin

klassische Suchmaschinen als auch für KI-Modelle attraktiv



Universität Regensburg

Dr. Gernot Deinzer

Leitung Abteilung IT- und Publikationsdienste
Universitätsbibliothek Regensburg

Knackpunkte KI

PROBLEME

KI-Bots

Training von großen LLMs

Scannen des Internets

analog klassischen Suchmaschinen

Aber:

häufiger
tiefer (mittels KI-Technologie)

Hoher Zugriff auf Repositorien

Sehr viele LLMs

Auch vermehrt in in Suchanfragen

Probleme

Schwer zu filtern

Ignoranz gegenüber herkömmlicher Tools (robots.txt)



Generated with leonardo.ai

¹ Siehe <https://coar-repositories.org/news-updates/open-repositories-are-being-profoundly-impacted-by-ai-bots-and-other-crawlers-results-of-a-coar-survey/>

Siehe AI bots are destroying Open Access,
<https://go-to-hellman.blogspot.com/2025/03/ai-bots-are-destroying-open-access.html>

KI- Bots

Lösungsmöglichkeiten

✓ I'm not a robot

CAPTCHAS

Sperrung gegenüber Zugriffen spezieller Muster

Problem:

Fehlender Inhalt in LLMs

Bedeutung Repositorien bezüglich LLMs

Open Access

Kein freier Zugriff auf content



Universität Regensburg

Dr. Gernot Deinzer

Leitung Abteilung IT- und Publikationsdienste
Universitätsbibliothek Regensburg

Nur gemeinsam haben wir eine Chance

VERNETZUNG

DINI

Arbeitsgruppe

Elektronisches Publizieren (E-Pub)

<https://dini.de/e-pub>

Ziel: Unterstützung und konzeptionelle Weiterentwicklung des elektronischen Publizierens an wissenschaftlichen Einrichtungen.

Bsp.: DINI-Zertifiakt

-> Vortrag Daniel Beucke

Confederation of Open Access Repositories (COAR)

**Notwendigkeit einer starken Zusammenarbeit,
um KI und Repositorien bestmöglich zu
verknüpfen**

Internationales Netzwerk

<https://coar-repositories.org>



13 Mitglieder aus Deutschland

1 Mitglied aus Bayern

Action Plan for Open Repositories in Europe: IMPACT-REPO



<https://coar-repositories.org/wp-content/uploads/2025/03/IMPACT-REPO-2025.pdf>

Powered by AI-ready infrastructure

Structure metadata and repositories for AI-driven research, discovery, and automation.

- Adopt machine-readable metadata, linked data practices, and text/data mining capabilities to ensure repository resources remain accessible and valuable in an AI-driven research ecosystem.



Danke für Ihre Aufmerksamkeit

Fragen

Kontakt:

Dr. Gernot Deinzer
Open Access Beauftragter
Abteilungsleitung IT- und Publikationsdienste
Geschäftsführung UR Data Hub
Fachreferent Mathematik, Physik
93042 Regensburg
email: gernot.deinzer@ur.de
<http://www.uni-regensburg.de/bibliothek>

