

# Embeddings und RAG: Revolutionäre Ansätze für intelligente Informationssuche

Embeddings und Retrieval-Augmented Generation (RAG) setzen neue Maßstäbe in der Informationssuche und -verarbeitung. Durch die direkte Abrufung relevanter Informationen aus den Embeddings entfällt das zeit- und ressourcenintensive Neutrainieren großer Sprachmodelle (LLMs). Diese effiziente Methode ermöglicht schnelle, präzise und kontextrelevante Antworten, maßgeschneidert für spezifische Anwendungsfälle.

# Was sind Embeddings?

Embeddings sind Vektoren, die Wörter oder Konzepte in einem mehrdimensionalen Raum repräsentieren. Diese Vektoren erlauben es uns, Beziehungen und Ähnlichkeiten zwischen Wörtern zu erkennen.

## Beispiel

Stell dir vor, wir haben die Wörter "König", "Mann" und "Königin". In einem embedding System sieht das ungefähr so aus:

- König (King) hat eine bestimmte Position im Raum.
- Mann (Man) hat auch eine Position.
- Königin (Queen) hat ebenfalls eine Position.

Die Beziehung zwischen diesen Wörtern kann als Vektoren dargestellt werden. Wenn wir von "König" den Vektor von "Mann" abziehen und den Vektor von "Frau" hinzufügen, landen wir bei "Königin". Mathematisch sieht das so aus:

$\text{König} - \text{Mann} + \text{Frau} = \text{Königin}$



## Grafische Darstellung

Stellen wir uns diese Vektoren grafisch vor:

1. "König" wird durch einen Punkt in einem Raum repräsentiert.
2. "Mann" wird ebenfalls durch einen Punkt repräsentiert.
3. "Frau" wird durch einen anderen Punkt repräsentiert.
4. "Königin" wird durch einen Punkt repräsentiert, der sich ergibt, wenn wir vom "König" zum "Mann" und dann von dort zur "Frau" gehen.

Dieses Prinzip ermöglicht es, komplexe Beziehungen und Analogien zwischen Wörtern zu erfassen und in einem mehrdimensionalen Raum darzustellen.

## Darstellung von Wortbeziehungen durch Vektoroperationen



*König - Mann + Frau = Königin*



# Die Darstellung von Embeddings

Der blaue Vektor repräsentiert "König", der rote Vektor repräsentiert "Mann", der grüne Vektor repräsentiert "Frau" und der violette Vektor repräsentiert "Königin". Diese Visualisierung zeigt, wie die Beziehungen zwischen den Wörtern dargestellt werden können.

## **Erklärung:**

- Der Vektor von "König" zu "Mann" repräsentiert das Attribut "männlich".
- Der Vektor von "Königin" zu "Frau" repräsentiert das Attribut "weiblich".
- Durch Subtraktion des "Mann"-Vektors von "König" und Hinzufügen des "Frau"-Vektors erhalten wir die Position von "Königin".

Dieses Prinzip ermöglicht es, komplexe Beziehungen und Analogien zwischen Wörtern darzustellen.

# Embeddings in Aktion

1

## Daten sammeln

Wir sammeln alle relevanten Daten in Form von Texten, Textabschnitten, Textelementen.

2

## Embeddings erstellen

Diese Daten werden in Embeddings umgewandelt, also in eine Form, die das System leicht durchsuchen kann.

3

## Beispiel Suche

Wenn du eine Frage stellst zu einem Text stellst, durchsucht das System die Embeddings und findet die passenden Informationen in deinen Daten.

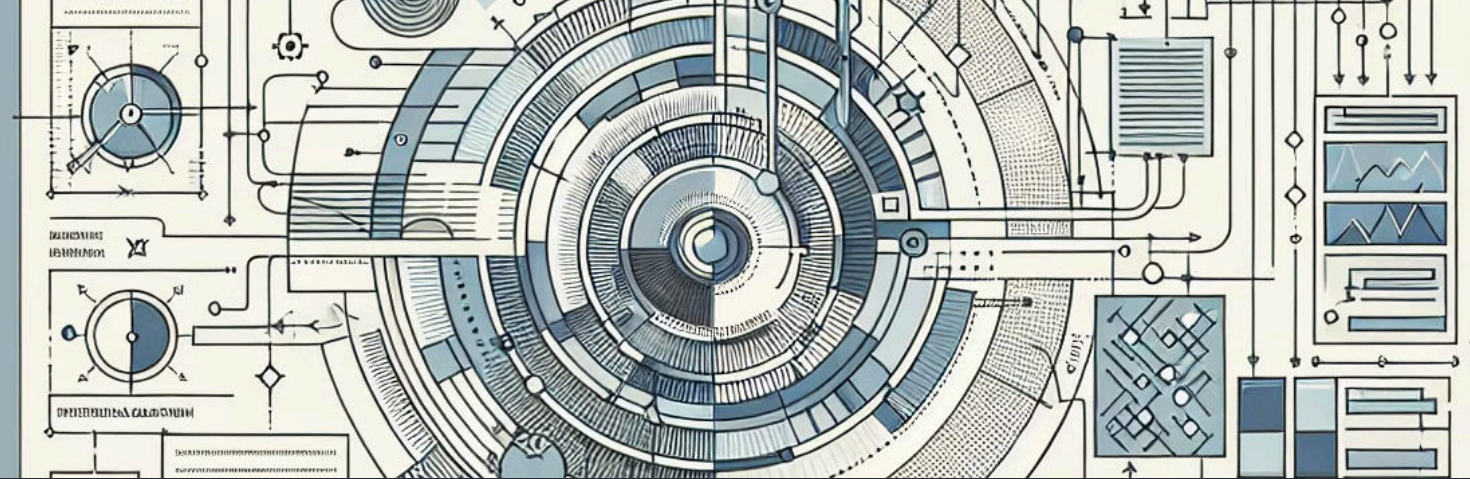
4

## LLM-Integration

Die Suchergebnisse aus den Embeddings werden in den Prompt an das LLM eingebaut. Dieses generiert dann eine präzise und gut aufbereitete Antwort.

# Vorteile von Embeddings

- Embeddings ermöglichen eine effiziente Darstellung und Verarbeitung von komplexen Beziehungen zwischen Wörtern und Konzepten.
- Sie können verwendet werden, um **Ähnlichkeiten** zwischen Objekten zu erkennen und **Analogien** zu ziehen.
- Embeddings können in **Suchfunktionen** eingesetzt werden, um relevante Inhalte basierend auf Bedeutung und Kontext zu finden.
- Sie ermöglichen eine **intuitive Visualisierung** von Wort- und Konzeptbeziehungen in einem mehrdimensionalen Raum.
- Embeddings können in **maschinellen Lernmodellen** verwendet werden, um die Leistung bei Aufgaben wie Textklassifikation oder Übersetzung zu verbessern.



# Embeddings und Retrieval-Augmented Generation (RAG)

Verbindung von Embeddings und fortschrittlichen Sprachmodellen wie RAG zur Erstellung leistungsstarker Suchanwendungen.

# Leistungsstarke Suchanwendungen

RAG ist eine Technologie, bei der Embeddings mit Großen Sprachmodellen (LLMs) kombiniert werden, um präzisere und kontextrelevantere Antworten zu generieren.

Das System durchsucht zunächst die Embeddings, um relevante Informationen zu finden. Diese werden dann in den Prompt für das LLM eingebaut, um eine maßgeschneiderte Antwort zu erstellen.



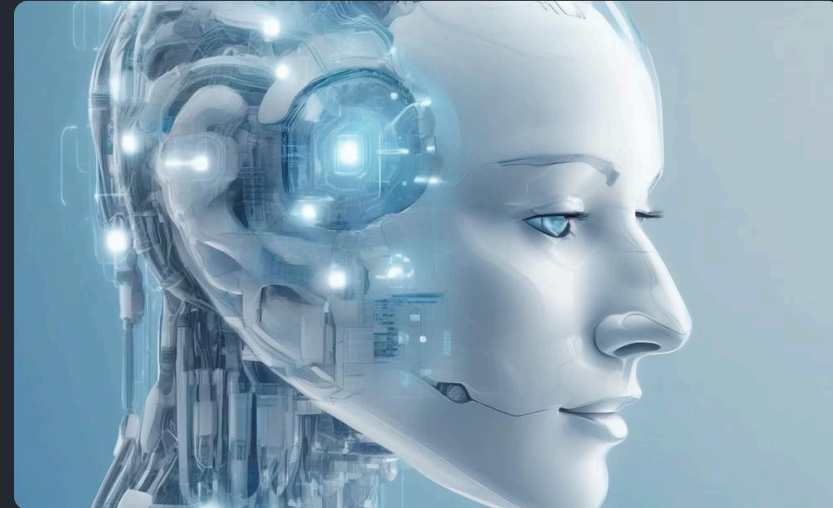


# Vorteile und Flexibilität von RAG: Effizienz und Anpassungsfähigkeit



## Vorteile von RAG

RAG erspart das komplette Neutrainieren eines LLMs, da die nötigen Informationen direkt aus den Embeddings abgerufen werden. So können Antworten schnell, präzise und kontextrelevant generiert werden.



## Flexibilität von RAG

Die Einbindung von Embeddings ermöglicht es, RAG-Systeme an spezifische Anwendungsfälle anzupassen und ständig zu verbessern, ohne das LLM komplett neu trainieren zu müssen.

# Zusammenfassung

Dieses Paper untersucht die Kombination von Embeddings und Retrieval-Augmented Generation (RAG) zur Entwicklung effizienter Informationssuchsysteme. Embeddings erfassen semantische Beziehungen zwischen Wörtern, während RAG diese nutzt, um präzise und kontextrelevante Antworten zu generieren. Dadurch entfällt das aufwändige Neutrainieren großer Sprachmodelle (LLMs). RAG ermöglicht flexible und anpassbare Suchanwendungen, die schneller und effizienter arbeiten. Erfahren Sie die Grundlagen, Vorteile und praktischen Anwendungen dieser Technologien in der modernen Informationsverarbeitung.

