

T2000

Jan Herrmann

7. Dezember 2025

Inhaltsverzeichnis

1	Einleitung	4
1.1	Motivation	4
1.2	Problemstellung	4
1.3	Zielsetzung	4
2	Theoretischer Hintergrund	5
2.1	Künstliche Intelligenz	5
2.2	Vektor-Store	7
2.2.1	Abgrenzung Vektor-Store und Vektor-Datenbank	7
2.2.2	Speicherverfahren für Vektordatenbanken	8
2.2.3	Suchverfahren	10
2.3	Neuronale Netze	12
2.4	Large Language Models	12
2.5	Retrieval-Augmented Generation	12
2.6	Schnittstellentechnologie	12
2.7	prompt Engineering	12
3	Anforderungsanalyse	13
3.1	Analyse	13

3.2	Benötigte Daten aus dem PIM System	13
3.3	Vergleich der LLM Modelle	13
4	Konzeption	13
4.1	Architektur	13
4.2	Datenfluss zwischen Vektor Store und Applikation	13
4.3	Schnittstellendesign	13
4.4	Promptdesign	13
5	Implementierung	13
5.1	Überblick über die Systemkomponenten	13
5.2	Umsetzung der Schnittstellen (Vector Store / Applikation / LLM) . .	13
5.3	Datenimport und -export	13
5.4	Integration des LLMs und Promptlogik	13
5.5	Fehlerbehandlung und Parallelität	13
6	Evaluation	13
6.1	Aufbau der Evaluationsbewertung	13
6.2	Bewertungsmetrik/-kriterien	13
6.3	Durchführung	13
6.4	Ergebnis	14
6.5	Diskussion	14

Abkürzungsverzeichnis

KI Künstliche Intelligenz. 4

ML Maschinelles Lernen. 4

Vektor-Store Leichtgewichtiges System zur Speicherung von Embeddings und Ausführung semantischer Ähnlichkeitssuchen. 4, 5

Vektordatenbank Datenbanksystem für Embeddings mit Mechanismen wie Sharding, Partitionierung, Caching und Replikation. 5

1 Einleitung

1.1 Motivation

1.2 Problemstellung

1.3 Zielsetzung

2 Theoretischer Hintergrund

2.1 Künstliche Intelligenz

Bevor eine präzise Definition von künstlicher Intelligenz möglich ist, muss geklärt werden, welches Ziel ein intelligentes System verfolgen soll. Russell und Norvig zeigen, dass gängige KI-Definitionen in der wissenschaftlichen Literatur entlang zweier zentraler Dimensionen variieren (vgl. [1], Kap. 1.1):

- **Mensch vs. Rationalität** Soll ein System wie ein Mensch denken oder handeln, oder soll es unabhängig vom Menschen ideal rational agieren?
- **Denken vs. Handeln** Soll Intelligenz anhand interner Denkprozesse oder anhand des beobachtbaren Verhaltens beurteilt werden?

Aus diesen beiden Dimensionen ergeben sich vier grundlegende Perspektiven auf KI, die unterschiedliche historische Forschungsrichtungen geprägt haben. Eine Übersicht dieser Einordnung zeigt Tabelle 1

Kategorie	Beschreibung
Systeme, die wie Menschen denken	Fokus auf Nachbildung menschlicher Denkprozesse, z.B. durch kognitive Modelle oder psychologische Theorien.
Systeme, die wie Menschen handeln	Intelligenz wird anhand menschlich ähnlichen Verhaltens beurteilt, unabhängig vom zugrunde liegenden Denkprozess.
Systeme, die rational denken	Fokus auf logische Schlussfolgerungen und formale Wissensrepräsentation.
Systeme, die rational handeln	Intelligente Agenten handeln zielgerichtet und optimal in ihrer Umgebung.

Tabelle 1: Eigene Darstellung in Anlehnung an [1], Kap. 1.1

Abgrenzung von KI und ML

Künstliche Intelligenz (KI) umfasst alle Verfahren, die darauf abzielen, intelligentes Verhalten technisch zu realisieren. Dazu gehören sowohl symbolische Ansätze wie Wissensrepräsentation und logisches Schließen als auch datengetriebene Methoden zur Wahrnehmung oder Sprachverarbeitung (vgl. [1], Kap. 1.1).

Maschinelles Lernen (ML) stellt ein klar abgegrenztes Teilgebiet der KI dar. Russell und Norvig beschreiben es als das Teilfeld der künstlichen Intelligenz, das sich mit Programmen befasst, die aus Erfahrung lernen (vgl. [1], Einleitung zu Teil VI).

Während KI somit als Oberbegriff sämtliche Methoden intelligenter Problemlösung einschließt, konzentriert sich ML ausschließlich auf Verfahren, die Wissen nicht explizit vorgegeben bekommen, sondern selbstständig aus Daten oder Erfahrungen

erschließen. ML bildet damit die Grundlage für viele moderne KI-Anwendungen, insbesondere für datengetriebene Systeme wie neuronale Netze oder Large Language Models.

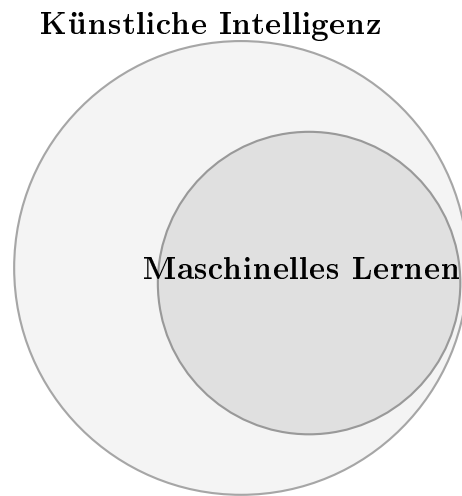


Abbildung 1: Einordnung von Maschinellern Lernen als Teilgebiete der Künstlichen Intelligenz

2.2 Vektor-Store

Vektor-Stores bilden die Grundlage für moderne KI-Anwendungen, die auf semantischer Ähnlichkeitssuche basieren. Sie dienen der Speicherung von Embeddings, also numerischen Repräsentationen von Text- oder Produktdaten, und ermöglichen effiziente Nearest-Neighbor-Abfragen. Damit stellen sie einen zentralen Baustein für Retrieval-Augmented Generation (RAG) und für Systeme dar, die kontextbezogene Informationen an große Sprachmodelle übergeben (vgl. Abschnitt 2.5).

Da der Begriff „Vektorstore“ in der Literatur häufig als Sammelbezeichnung für leichtgewichtige, embeddingbasierte Speichersysteme verwendet wird, ist eine Abgrenzung zu vollwertigen Vektordatenbanken notwendig. Letztere integrieren zusätzliche Mechanismen wie Sharding, interne Partitionierung, Caching oder Replikation und unterscheiden sich damit deutlich in ihrer Komplexität und Skalierbarkeit.

Im Folgenden werden die in [2] vorgestellten Speicher- und Suchmechanismen von Vektordatenbanken zusammengefasst und auf das Anwendungsszenario dieser Arbeit übertragen.

Tabelle 2: Abgrenzung von Vektor-Store und Vektordatenbank

Merkmal	Vektor-Store	Vektordatenbank
Speicherung von Embeddings	✓	✓
k-nächste-Nachbarn-Suche (k-NN)	✓	✓
Sharding über mehrere Knoten	✗	✓
Interne Partitionierung innerhalb eines Knotens	✗	✓
Caching-Mechanismen	✗	✓
Replikation	✗	✓
Approximate Nearest Neighbor (ANN)-Indexe	✗	✓
Metadatenverwaltung und Konsistenzmodelle	✗	✓

2.2.1 Abgrenzung Vektor-Store und Vektor-Datenbank

Der Begriff Vektor-Store wird in der Literatur nicht einheitlich verwendet und dient häufig als Sammelbezeichnung für leichtgewichtige Systeme, die Embeddings speichern und eine grundlegende semantische Ähnlichkeitssuche bereitstellen. Solche Systeme konzentrieren sich in der Regel auf das Einfügen von Embeddings und die Ausführung von k-nächste-Nachbarn-Abfragen, ohne jedoch erweiterte Datenbankfunktionen bereitzustellen.

Unter einer Vektordatenbank werden hingegen vollwertige Datenbanksysteme verstanden, die neben der Speicherung und Suche von Embeddings zusätzliche Verwaltungs- und Infrastrukturmechanismen integrieren. Dazu zählen insbesondere Sharding, interne Partitionierung, Caching, Replikation, Indexstrukturen für Approximate Nea-

rest Neighbor (ANN) sowie Metadatenverwaltung und Konsistenzmodelle. Diese Systeme sind auf hohe Skalierbarkeit, Robustheit und Performanz im produktiven Einsatz ausgelegt.

Für die vorliegende Arbeit wird daher folgende Unterscheidung getroffen:

- **Vektor-Store:** leichtgewichtiges System zur Speicherung von Embeddings und zur Durchführung semantischer Ähnlichkeitssuchen.
- **Vektordatenbank:** vollständiges Datenbankmanagementsystem mit skalierbaren Speicher- und Verwaltungsmechanismen.

Diese definitorische Abgrenzung dient der Klarheit und legt die einheitliche Verwendung der Begriffe im weiteren Verlauf der Arbeit fest.

2.2.2 Speicherverfahren für Vektordatenbanken

Die im Folgenden beschriebenen Mechanismen stellen zentrale Bestandteile moderner, skalierbarer Vektordatenbanksysteme dar. Sie werden im produktiven Umfeld zur Optimierung von Leistung, Verfügbarkeit und Robustheit eingesetzt. Für den im Rahmen dieser Arbeit entwickelten Prototypen spielen diese Konzepte jedoch keine operative Rolle, da ein leichtgewichtiges Single-Node-System ohne verteilte Architektur eingesetzt wird. Die Ausführungen dienen daher der theoretischen Einordnung und sollen den technologischen Kontext verdeutlichen, in dem sich Vektordatenbanken typischerweise bewegen.

horizontale Datenpartitionierung über mehrere Maschinen (auch Sharding genannt) bezeichnet ein Verfahren, bei dem eine Datenbank in mehrere logisch getrennte und auf verschiedene physische Knoten verteilte Teilmengen („Shards“) aufgeteilt wird. Durch diese Aufteilung wird das Gesamtdatenset in kleinere, handhabbarere Einheiten zerlegt, was Skalierbarkeit, Lastverteilung und die Verwaltung großer Datenmengen erleichtert (vgl. [2] S. 3).

horizontale Datenpartitionierung innerhalb einer Maschine die horizontale Datenpartitionierung innerhalb einer Maschine bezeichnet die Aufteilung der in einer einzelnen Datenbankinstanz gespeicherten Vektordaten in mehrere logisch getrennte Teilmengen („Partitionen“). Im Unterschied zum Sharding, das Daten über mehrere physische Knoten verteilt, erfolgt diese Form der Partitionierung ausschließlich innerhalb eines Systems. Ziel ist es, lokale Abfragen effizienter auszuführen, Speicherressourcen besser auszunutzen und parallele Verarbeitung zu ermöglichen. Partitionen können beispielsweise über Wertebereiche (Range-Partitioning), vordefinierte Kategorien (List-Partitioning) oder Hash-Verfahren gebildet werden. Durch diese interne Strukturierung müssen Suchanfragen nur gegen relevante Partitionen ausgeführt werden, was insbesondere bei großen Einbettungsräumen die Latenz der semantischen Ähnlichkeitssuche reduziert (vgl. [2], S. 3f.).

In modernen Vektordatenbanken wird Partitionierung häufig mit Sharding kombiniert (vgl. Abbildung 2).

Während Sharding die Skalierung über mehrere Knoten ermöglicht, optimiert die interne Partitionierung die Datenorganisation innerhalb jedes Shards. Beide Verfahren bilden damit die Grundlage für performante Retrieval-Systeme in Vektor Stores und Vektordatenbanken.

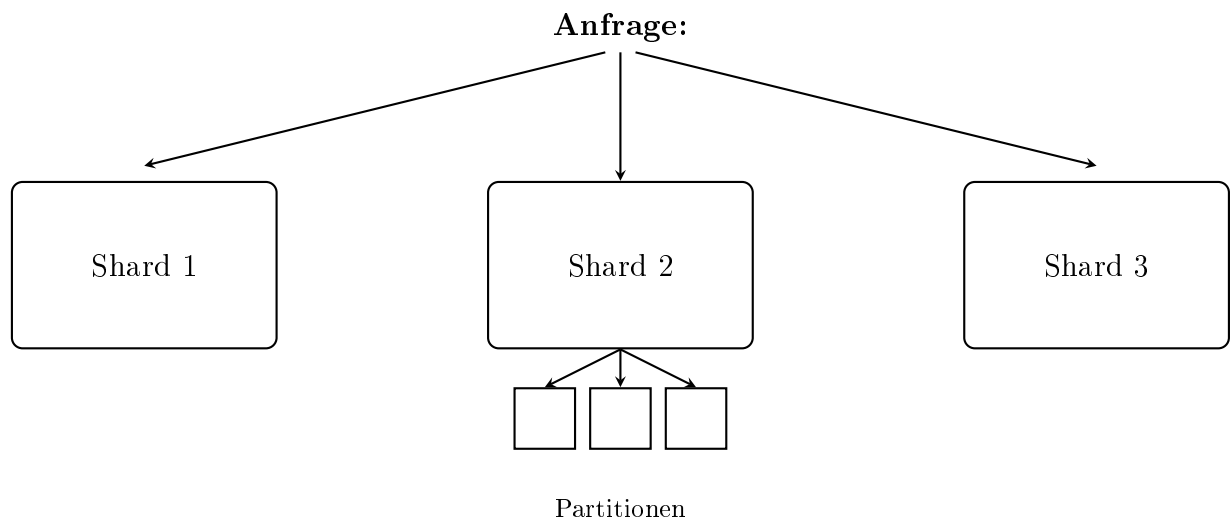


Abbildung 2: Sharding über mehrere Maschinen und interne Partitionierung innerhalb eines Shards

Caching-Mechanismen in Vektordatenbanken beschreiben das Zwischenspeichern häufig oder kürzlich genutzter Daten in besonders schnellen Speichermedien (z.B. RAM), um Zugriffszeiten zu reduzieren und die Last auf der eigentlichen Datenbank zu verringern. Für Vektordatenbanken ist Caching ein zentraler Mechanismus, da semantische Ähnlichkeitssuchen typischerweise rechenintensiv sind und von wiederholten Abfragen profitieren (vgl. [2], S. 4f.).

Im Gegensatz zu klassischen Schlüssel-Wert-Caches (etwa Redis) ist das Caching in VDBs herausfordernder, da Embeddings hochdimensionale Vektoren darstellen und identische Anfragen selten auftreten. Daher kommen allgemeine Cache-Strategien zum Einsatz, die unabhängig vom genauen Vektorinhalt arbeiten, da der Vergleich hochdimensionaler Vektoren zur Bestimmung von Cache-Schlüsseln ungeeignet ist.

Zu den gängigen Verfahren zählen:

- **First-In First-Out (FIFO):** entfernt das älteste Element im Cache; einfach, aber ohne Berücksichtigung der Zugriffshäufigkeit.
- **Least Recently Used (LRU):** löscht den am längsten nicht genutzten Eintrag; gut geeignet für Arbeitslasten mit zeitlicher Lokalität.
- **Most Recently Used (MRU):** entfernt das zuletzt genutzte Element; sinnvoll bei einmaligen Zugriffsmustern.
- **Least Frequently Used (LFU):** bevorzugt das Entfernen seltener genutzter Einträge; vorteilhaft bei stabilen Zugriffshäufigkeiten.

Einige Systeme nutzen zudem **partitioniertes Caching**, bei dem Vektordaten in Gruppen (z.B. nach Kategorien oder Zugriffsmustern) getrennt gecacht werden, um Ressourcen gezielt zu optimieren. Insgesamt trägt Caching wesentlich zur Reduktion der Abfragelatenz und zur Stabilisierung der Systemlast bei, insbesondere bei wiederkehrenden Ähnlichkeitsanfragen.

Replikation bezeichnet das Anlegen und Verteilen mehrerer Kopien von Vektordaten auf unterschiedliche Knoten eines verteilten Systems, um Ausfallsicherheit, Verfügbarkeit und Leselastverteilung zu erhöhen. Während sie für die semantische Ähnlichkeitssuche nicht unmittelbar leistungsbestimmend ist, stellt sie einen zentralen Mechanismus für die betriebliche Robustheit moderner Vektordatenbanken dar (vgl. [2], S. 5f.).

Insgesamt trägt Replikation wesentlich zur Robustheit und Verfügbarkeit verteilter Vektordatenbanksysteme bei, steht jedoch weniger im direkten Zentrum der Ähnlichkeitssuche als vielmehr ihrer infrastrukturellen Betriebsstabilität.

2.2.3 Suchverfahren

Diese Arbeit befasst sich mit zwei Kategorien von Suchverfahren, deren grundlegende Funktionsweise in Abbildung ?? schematisch gegenübergestellt ist. Die Darstellung der approximierenden Suche stellt hierbei ein vereinfachtes, heuristisches Beispiel dar: Es illustriert das Prinzip, dass ANN-Verfahren den Suchraum gezielt einschränken, ohne jeden Punkt im Vektorraum zu prüfen. In der Praxis existieren verschiedene ANN-Algorithmen (z.B. HNSW, IVF, LSH), die dieses Prinzip auf unterschiedliche Weise umsetzen. Die Abbildung dient somit der didaktischen Veranschaulichung des Grundkonzepts.

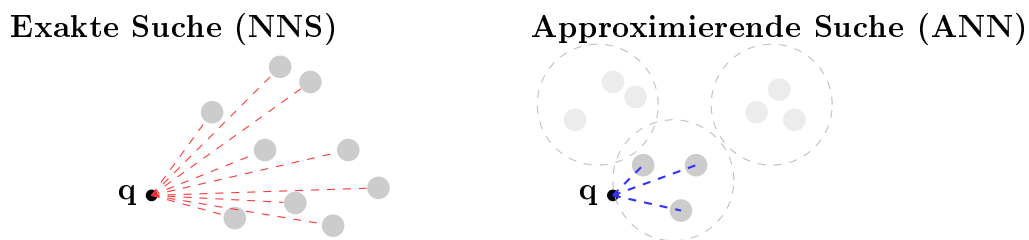


Abbildung 3: Vergleich zwischen exakter Nearest-Neighbor-Suche (NNS) und approximierender Suche (ANN). Bei exakter Suche berechnet das System die Distanz zwischen der Anfrage und *jedem* Punkt im Raum (rote Linien). ANN-Verfahren durchsuchen hingegen nur relevante Regionen des Suchraums (blauer Pfad).

Exakte Nearest-Neighbor-Suche (NNS) bezeichnet das Verfahren, für einen gegebenen Anfragevektor denjenigen Vektor in einer Menge zu bestimmen, der gemäß einem definierten Distanzmaß (z.B. euklidische Distanz) am nächsten liegt. In einfachster Form erfolgt dies durch eine lineare Suche, bei der für jedes gespeicherte Embedding die Distanz zum Anfragevektor berechnet und anschließend der Vektor mit dem geringsten Abstand ausgewählt wird (vgl. [2], S. 6f.).

Da die Laufzeit dieses Ansatzes linear ($O(n)$) mit der Anzahl der gespeicherten Em-

beddings wächst, stoßen Systeme bei großen Datenmengen schnell an ihre Leistungsgrenzen. In solchen Fällen können approximierende Verfahren eingesetzt werden. Die Grundidee besteht darin, eine Lösung zu liefern, die zwar nicht exakt optimal ist, dafür jedoch eine deutlich bessere Laufzeit und geringeren Speicherbedarf ermöglicht. Die leichte Abweichung vom exakten Ergebnis wird somit bewusst zugunsten einer höheren Effizienz in Kauf genommen.

Approximierende Nearest-Neighbor-Suche (ANNS) verzichtet auf eine vollständige Durchmusterung aller Embeddings und nutzt stattdessen probabilistische oder heuristische Verfahren, um den Suchraum gezielt einzugrenzen. Heuristische Verfahren nutzen vereinfachte Entscheidungsregeln, die eine schnelle, aber nicht zwingend optimale Abschätzung erlauben. Ziel ist es, möglichst schnell einen Vektor zu finden, der dem Anfragevektor sehr ähnlich ist, ohne zwingend den exakt nächsten Nachbarn bestimmen zu müssen. Dadurch lassen sich deutliche Laufzeitgewinne erzielen, insbesondere bei großen Datenmengen (vgl. [2], S. 6f.).

Während exakte Verfahren insbesondere in kleinen Vektor-Stores mit geringem Datenvolumen eingesetzt werden können, sind approximierende Verfahren bei skalierbaren Vektordatenbanken mit einer großen Anzahl von Embeddings sinnvoll.

- 2.3 Neuronale Netze**
- 2.4 Large Language Models**
- 2.5 Retrieval-Augmented Generation**
- 2.6 Schnittstellentechnologie**
- 2.7 prompt Engineering**

3 Anforderungsanalyse

3.1 Analyse

3.2 Benötigte Daten aus dem PIM System

3.3 Vergleich der LLM Modelle

4 Konzeption

4.1 Architektur

4.2 Datenfluss zwischen Vektor Store und Applikation

4.3 Schnittstellendesign

4.4 Promptdesign

5 Implementierung

5.1 Überblick über die Systemkomponenten

5.2 Umsetzung der Schnittstellen (Vector Store / Applikation / LLM)

5.3 Datenimport und -export

5.4 Integration des LLMs und Promptlogik

5.5 Fehlerbehandlung und Parallelität

6 Evaluation

6.1 Aufbau der Evaluationsbewertung

6.2 Bewertungsmetrik/ -kriterien

6.3 Durchführung

6.4 Ergebnis

6.5 Diskussion

LLM ohne RAG halluziniert → schlecht für PIM-Daten (Faktentreue).

Prompt mit Rollenbeschreibung („Du bist ein Marketing-Experte...“) besser als generischer Prompt.

Techniktexte brauchen präzisere Struktur -> Template hilft.

Marketingtexte profitieren von höherer Kreativität → Temperatur-Einstellungen diskutieren.

Literatur

- [1] S. J. Russell und P. Norvig, *Artificial intelligence: a modern approach* (Prentice Hall series in artificial intelligence). Upper Saddle River: Prentice Hall, 1995, 932 S., ISBN: 978-0-13-103805-9 978-0-13-360124-4.
- [2] L. Ma u. a., *A Comprehensive Survey on Vector Database: Storage and Retrieval Technique, Challenge*, 16. Juni 2025. DOI: 10.48550/arXiv.2310.11703. arXiv: 2310.11703[cs]. besucht am 3. Dez. 2025. Adresse: <http://arxiv.org/abs/2310.11703>.