



Retrieval-Augmented Generation (RAG):

Wissen gezielt nutzen,
Antworten präzise generieren

Unternehmenseigenes Wissen mit der Sprachkompetenz von Large Language Models verbinden.

Künstliche Intelligenz hat mit Large Language Models beeindruckende Sprachfähigkeiten erreicht. Doch für den Unternehmenseinsatz fehlt oft das spezifische Domänenwissen. Retrieval-Augmented Generation (RAG) schließt diese Lücke elegant: Es verbindet die Stärken generischer Sprachmodelle mit unternehmenseigenem Wissen – ohne aufwendiges Training und mit voller Kontrolle über die eigenen Daten.

Begriffserklärung und Abgrenzung

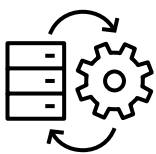
Retrieval-Augmented Generation (RAG) ist ein technologischer Ansatz, der die Anfragen an ein beliebiges Large Language Model (LLM) um eigene Daten und Informationen anreichert, um relevantere und präzisere Antworten zu erzeugen. Ein RAG-System verbindet dabei zwei Welten: die Informationsabfrage (Retrieval) und die Textgenerierung (Generation).

Im Kern ermöglicht RAG einem LLM, auf Informationen zuzugreifen, die nicht Teil der ursprünglichen Trainingsdaten waren. Anstatt sich ausschließlich auf das Wissen aus dem Training zu verlassen, kann ein RAG-System während der Anfrage dynamisch zusätzliche Informationen aus externen Quellen abrufen und in seine Antworten integrieren. Diese Quellen umfassen Dokumente, Datenbanken (durch einen semantischen Layer), das Web oder APIs zu Tools und Systemen.

RAG bildet damit die Grundlage, um mit eigenen Unternehmensdaten innerhalb von KI-Prozessen zu arbeiten – wie ein beleseener Assistent, der sich vor jeder Antwort noch einmal gezielt in der Unternehmensbibliothek informiert.

Wie funktioniert Retrieval-Augmented Generation?

Die Funktionsweise von RAG in fünf zentralen Schritten



Datenquellen erschließen

1.



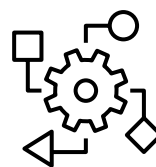
Vektorisierung und Indexierung

2.



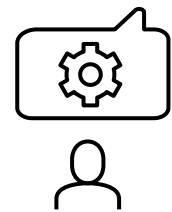
Intelligente Suche (Retrieval)

3.



Kontext-Anreicherung (Augmentation)

4.



Präzise Antwortgenerierung (Generation)

5.

1. Datenquellen erschließen

Unternehmenseigene Dokumente, Datenbanken, Wikis und andere Wissensquellen werden als Grundlage verfügbar gemacht. Diese bleiben dabei unter vollständiger Kontrolle des Unternehmens.

2. Vektorisierung und Indexierung

Die Inhalte werden in semantische Vektoren umgewandelt und in einer Vektordatenbank gespeichert. Dieser Prozess ermöglicht es, inhaltlich ähnliche Informationen schnell zu finden.

3. Intelligente Suche (Retrieval)

Bei einer Nutzeranfrage durchsucht das System die Vektordatenbank nach den relevantesten Informationen. Dabei werden nicht nur Schlüsselwörter, sondern semantische Ähnlichkeiten berücksichtigt.

4. Kontext-Anreicherung (Augmentation)

Die gefundenen relevanten Textabschnitte werden als zusätzlicher Kontext an das LLM übergeben. Das Modell erhält damit spezifisches Wissen für die konkrete Anfrage.

5. Präzise Antwortgenerierung (Generation)

Das LLM generiert auf Basis des angereicherten Kontexts eine Antwort, die sowohl die allgemeine Sprachkompetenz als auch das spezifische Domänenwissen vereint.

Warum sollten Unternehmen RAG-Systeme einsetzen?

Unternehmen wählen RAG für 30–60% ihrer KI-Anwendungsfälle (Vectara, 2024)¹, insbesondere wenn hohe Genauigkeit, Transparenz und verlässliche Ergebnisse mit eigenen Daten gefordert sind. Die Adoption von RAG-Architekturen für Enterprise-Anwendungen liegt bei 51% (Menlo Ventures, 2024)². Standard-LLMs sind oft zu generisch und kennen die Business-Domäne nicht, während selbst trainierte Modelle hohe Wartungsaufwände bedeuten.

RAG vereint das Beste aus beiden Welten: Es nutzt die Sprachkompetenz von Standard-LLMs und reichert sie mit eigenem Domänenwissen an – ohne aufwendiges Training. In der Praxis berichten Unternehmen von erheblichen Effizienzgewinnen: Juristische Recherchen werden deutlich beschleunigt, medizinische Diagnoseunterstützung wird präziser, und Wissensarbeiter sparen signifikant Zeit. Die Datenhoheit bleibt gewahrt, sensible Informationen werden nicht ins Modell trainiert, und neue Datenquellen lassen sich flexibel hinzufügen.

Das Ergebnis: Präzise, aktuelle Antworten mit Quellenangaben, die Halluzinationen reduzieren und echten Mehrwert schaffen. Laut Deloitte berichten bereits 42% der Organisationen von signifikanten Produktivitäts-, Effizienz- und Kostenvorteilen durch GenAI (Deloitte, 2024)³.

Wo sind die konkreten Einsatzgebiete?



Wissensmanagement

Investigativer Zugang zu Unternehmenswissen für alle Hierarchieebenen



KI-Agentensysteme

Spezifisches Wissen für autonome Agenten im Unternehmenskontext



Datenanalyse

Vergleichen, Klassifizieren und Visualisieren per natürlicher Sprache



Kundensupport

Präzise Antworten aus aktuellen Produktdokumentationen und FAQs



Compliance

Zugriff auf Regularien und Verträge mit quellenbasierten Antworten



F&E

Peer-reviewte Studien zeigen deutliche Effizienzgewinne bei systematischen Literaturreviews durch den Einsatz von Generative AI – insbesondere im Screening und in der Datenauswertung werden Arbeitsaufwände um bis zu 50 % reduziert (peer-reviewed Studien, 2024–2025)⁴.



Braucht es ein RAG oder reicht ein LLM? Die richtige Wahl

Ein **klassisches LLM** mit großem Kontextfenster (bis 128k Tokens) eignet sich für strukturierten Text in festem Ablauf – Skripte, Reports, Gesprächsverläufe. RAGs dagegen brillieren bei großen, heterogenen Wissensbasen – Knowledge Bases, Archive, Rechtsdokumente.

Im Vergleich zum **Fine-Tuning**, bei dem LLMs mit Unternehmens-Know-How nachtrainiert werden und das Modell dauerhaft verändert wird (was hohe Wartungsaufwände nach sich ziehen kann), trennen RAG-Systeme Modell und Wissen. Patrick Lewis, der den Begriff RAG prägte, betont, dass die Methode mit nur fünf Zeilen Code implementierbar ist – schneller und kostengünstiger als Neutrainung (NVIDIA Interview, Januar 2025)⁵.

Gilt es also, spezielles Wissen zu integrieren, das sich womöglich häufig ändert, ist meist ein RAG-System die bessere Wahl. Gibt es dieses spezielle Wissen nicht, oder ist es nur wenig Veränderung unterworfen, kann ein klassisches LLM, eventuell mit Fine-Tuning, die richtige Wahl sein.

Und wo sind (aktuell noch) Grenzen?

RAG ist kein Allheilmittel. Die Antwortqualität hängt von der Retrieval-Güte ab, zusätzliche Suchvorgänge erhöhen die Latenz, und optimales Chunking erfordert Expertise. Datenschutz und Sicherheit bleiben kritische Herausforderungen – Studien zeigen, dass die meisten Unternehmen bereits GenAI-bezogene Sicherheitsvorfälle erlebt haben. Auch mit RAG bleiben Kontextfenster-Limits bestehen, und bei widersprüchlichen Quellen muss das System Prioritäten setzen. Deshalb gilt: Sorgfältige Planung, Data Governance und kontinuierliche Optimierung sind essentiell.

Praktische Implementierung in Unternehmenssysteme

Ein neues RAG-System entsteht selten auf der grünen Wiese, sondern muss sich in eine bestehende IT-Landschaft integrieren. Dazu empfiehlt sich unbedingt ein planvolles Vorgehen, um Insellösungen zu vermeiden:

Technische Architektur

Enterprise-RAG-Plattformen bieten produktionsreife Infrastruktur mit SOC 2-, GDPR- und ISO 27001-Compliance (Firecrawl, 2025)⁶. Eine bestehende gute Datenqualität vorausgesetzt, können die meisten Unternehmen grundlegende RAG-Workflows innerhalb weniger Wochen implementieren.

Sicherheit und Datenschutz

73% der Organisationen nennen Sicherheitsbedenken als Haupthürde für KI-Implementierung (Ragwalla, 2025)⁷. Erfolgreiche RAG-Systeme setzen daher auf:

Ende-zu-Ende-Verschlüsselung
für Datenübertragung und -speicherung

Föderierte Architekturen,
bei denen Rohdaten lokal bleiben

GDPR-konforme Verarbeitung
mit Pseudonymisierung und Recht auf Löschung

Attribute-basierte Zugriffskontrolle (ABAC) für granulare Berechtigungen

Integration bestehender Systeme

Moderne RAG-Implementierungen integrieren direkt mit CRM-, ERP- und anderen Unternehmenssystemen für Echtzeiteinblicke. APIs ermöglichen nahtlose Anbindung an bestehende Workflows, während Middleware-Lösungen Legacy-Systeme einbinden.

Data Governance etablieren

Klare Richtlinien für Datenqualität, Zugriffsrechte und Aktualisierungsprozesse definieren. Eine wiederholbare Daten-Vorbereitungspipeline ist essentiell: Filtern, Segmentieren, Standardisieren, Metadaten extrahieren.

Vektorisierung verstehen

Die richtige Wahl von Embedding-Modellen (z.B. Sentence-BERT) und Vektordatenbanken (Pinecone, Weaviate, Chroma) ist entscheidend für die Performance.

Chunking-Strategien entwickeln

Die optimale Aufteilung von Dokumenten beeinflusst maßgeblich die Antwortqualität. Methoden wie Overlap, Sliding Window oder Semantic Chunking müssen evaluiert werden.

Pilot-Projekte starten

Mit fokussierten Use Cases beginnen, wo hochwertige, strukturierte Daten verfügbar sind: Kundensupport, internes Wissensmanagement, Compliance-Dokumentation.

Metriken definieren

Erfolg messbar machen durch Relevanz-Scores, Antwortzeiten und Nutzerzufriedenheit. A/B-Tests zwischen RAG und reinem LLM etablieren.

Risiken und Mitigation

Datenschutz-Herausforderungen

RAG-Systeme müssen GDPR, HIPAA und andere Regularien einhalten, was explizites Consent-Management und Datenminimierung erfordert. 97% der Unternehmen hatten mindestens einen GenAI-bezogenen Sicherheitsvorfall (Capgemini, 2024)¹¹. Lösungsansätze umfassen:

- **Differential Privacy:** Hinzufügen von kontrolliertem Rauschen zum Schutz individueller Datenpunkte
- **Homomorphe Verschlüsselung:** Berechnungen auf verschlüsselten Daten ohne Entschlüsselung
- **Föderiertes Lernen:** Modelle, die Daten lokal verarbeiten und nur Modell-Updates teilen

Prompt Injection und Sicherheit

Adversariale Inhalte in Nutzeranfragen oder abrufbaren Dokumenten können Modellantworten manipulieren. Gegenmaßnahmen:

- **Input-Validierung** und Sanitization aller Nutzeranfragen
- **Sandboxing** der Retrieval-Pipeline
- **Kontinuierliche Sicherheitstests** in CI/CD-Pipelines

Evolution und Meilensteine



Ausblick

Die Entwicklung von RAG-Systemen steht erst am Anfang einer vielversprechenden Evolution. Fünf zentrale Trends zeichnen sich ab:

1. Multimodale und Echtzeit-RAG

RAG entwickelt sich über textbasiertes Retrieval hinaus zu Bildern, Videos und Audio, mit dynamischen Echtzeit-Feeds und Hybrid-Search-Techniken.

2. On-Device und Edge-RAG

Lokale Verarbeitung für besseren Datenschutz und reduzierte Latenz, besonders relevant für sensible Bereiche wie Gesundheit und Finanzen.

3. Automatisierte Optimierung

Selbstlernende Systeme optimieren Chunking-Strategien, Retrieval-Parameter und Kontext-Zusammenstellung automatisch.

4. Förderierte und Privacy-Preserving RAG

Neue Architekturen, bei denen Rohdaten niemals die lokale Umgebung verlassen und nur verschlüsselte Modell-Updates geteilt werden.

5. Standardisierung und Demokratisierung

Einheitliche Frameworks und Best Practices vereinfachen die Implementierung. Open-Source-Modelle und Cloud-Lösungen senken Einstiegshürden für KMUs.

Erfolgsmessung

Wie jedes implementierte System, sollte auch ein RAG mit einem kritischen Blick kontinuierlich beobachtet werden. Nur so können Fehler und Probleme rechtzeitig erkannt, oder Verbesserungsmöglichkeit sinnvoll geplant werden.

Dabei sollten sowohl die Retrieval-Qualität, als auch die Generierungsqualität mit geeigneten Methoden gemessen werden. Dazu haben sich erste Standards etabliert: NDCG Score, Recall-Rate, Answer Relevance Score, Faithfulness Score, Hallucination Rate und andere.

Da ein RAG-System niemals Selbstzweck ist, sollten auch Business-Metriken (bringt das RAG-System wirklich die Vorteile, für die es eingerichtet wurde?) eingeführt werden.

Fazit – RAG als Brücke zwischen Sprach-KI und Unternehmenswissen

Retrieval-Augmented Generation markiert einen Paradigmenwechsel im Unternehmenseinsatz von KI. Mit über 1.200 wissenschaftlichen Publikationen allein 2024 (ArXiv, 2025)⁹ ist RAG keine experimentelle Technologie mehr, sondern produktionsreife Realität.

Für Unternehmen bedeutet das: Der Zeitpunkt, sich mit RAG auseinanderzusetzen, ist jetzt. Wer die Technologie versteht, gezielt einsetzt und kontinuierlich optimiert, kann einen echten Wettbewerbsvorteil aufbauen. Die Datenhoheit bleibt gewahrt, die Kosten überschaubar, die Ergebnisse einzigartig.

RAG ist mehr als eine technische Lösung – es ist der Schlüssel, um die Kluft zwischen allgemeiner KI und spezifischem Unternehmenswissen zu überbrücken. In einer Welt, in der Daten zum wichtigsten Asset werden, macht RAG dieses Asset intelligent nutzbar.

Quellenverzeichnis

- 1 **Vectara** (2024): Enterprise RAG Predictions for 2025. <https://www.vectara.com/blog/top-enterprise-rag-predictions>
 - 2 **Menlo Ventures** (2024): The State of Generative AI in the Enterprise. <https://menlovc.com/2024-the-state-of-generative-ai-in-the-enterprise/>
 - 3 **Deloitte** (2024): GenAI Survey, zitiert in Vectara Enterprise RAG Predictions.
 - 4 **Peer-reviewte Studien zu systematischen Reviews** (2024–2025): Enhancing Systematic Literature Reviews with Generative Artificial Intelligence (JAMIA, 2024). <https://academic.oup.com/jamia/article/32/4/616/8045049>
 - 5 **NVIDIA Blog** (Januar 2025): Interview mit Patrick Lewis. <https://blogs.nvidia.com/blog/what-is-retrieval-augmented-generation/>
 - 6 **Firecrawl** (2025): Best Pre-Built Enterprise RAG Platforms in 2025. <https://www.firecrawl.dev/blog/best-enterprise-rag-platforms-2025>
 - 7 **Ragwalla** (2025): The Complete Guide to Enterprise AI Security. <https://ragwalla.com/docs/guides/the-complete-guide-to-enterprise-ai-security-rag-agents-compliance-in-2025>
 - 8 **Lewis, P.** et al. (Mai 2020): Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. ArXiv:2005.11401. <https://arxiv.org/abs/2005.11401>
 - 9 **Systematic Review of Key RAG Systems** (Juli 2025). ArXiv:2507.18910. <https://arxiv.org/html/2507.18910v1>
 - 10 **RAGFlow** (Juli 2025): RAG at the Crossroads – Mid-2025 Reflections. <https://ragflow.io/blog/rag-at-the-crossroads-mid-2025-reflections-on-ai-evolution>
 - 11 **Capgemini** (2024): Studie zu GenAI-Sicherheitsvorfällen, zitiert in verschiedenen Branchenanalysen. <https://www.capgemini.com/de/news/pressemitteilung/cybersicherheit-von-unternehmen-neue-bedingungen-durch-ki-generative-ai/>
- Implementation and Evaluation of an Additional GPT-4-based Reviewer** in PRISMA-based Medical Systematic Literature Reviews (Int. J. Med. Informatics, 2024). <https://publica-rest.fraunhofer.de/server/api/core/bitstreams/60db6241-Oc51-4a79-b3c2-06815dbc0c63/content>
- Do It Faster with PICOS:** Generative AI-Assisted Systematic Review Screening (J. Biomed. Informatics, 2025). <https://www.sciencedirect.com/science/article/pii/S1532046425000899>

Projektverantwortliche

Stephan Clasen, Technical Director, denkwerk

Katharina Jäger, Head of Innovation & Technology, BVDW

Bundesverband Digitale Wirtschaft (BVDW) e.V.

Der Bundesverband Digitale Wirtschaft (BVDW) e. V. ist die Interessenvertretung für Unternehmen, die digitale Geschäftsmodelle betreiben oder deren Wertschöpfung auf dem Einsatz digitaler Technologien beruht. Mit seinen Mitgliedern aus der gesamten Digitalen Wirtschaft gestaltet der BVDW bereits heute die Zukunft – durch kreative Lösungen und modernste Technologien. Als Impulsgeber, Wegweiser und Beschleuniger digitaler Geschäftsmodelle setzt der Verband auf faire und klare Regeln und tritt für innovationsfreundliche Rahmenbedingungen ein. Dabei hat der BVDW immer Wirtschaft, Gesellschaft und Umwelt im Blick. Neben der DMEXCO, der führenden Fachmesse für Digitales Marketing und Technologien, und dem Deutschen Digital Award richtet der BVDW auch den CDR-Award, die erste Preisverleihung im DACH-Raum für Digitale Nachhaltigkeit und Verantwortung sowie eine Vielzahl von Fachveranstaltungen aus. Mehr Informationen finden Sie unter

Working Group Künstliche Intelligenz

Die gewinnbringende und verantwortungsvolle Nutzung von künstlicher Intelligenz (KI) in der deutschen digitalen Wirtschaft steht im Fokus der Arbeit der Working Group. Ziel ist es, Fragen rund um die Veränderungen der Wertschöpfungskette der digitalen Wirtschaft zu beantworten und Lösungsansätze für die ethischen, sozialen und rechtlichen Herausforderungen durch KI zu bieten.

Kontakt

Katharina Jäger, Head of Innovation & Technology, jaeger@bvdw.org

Bundesverband Digitale Wirtschaft (BVDW) e.V.

Obentrautstraße 55, 10963 Berlin

www.bvdw.org

