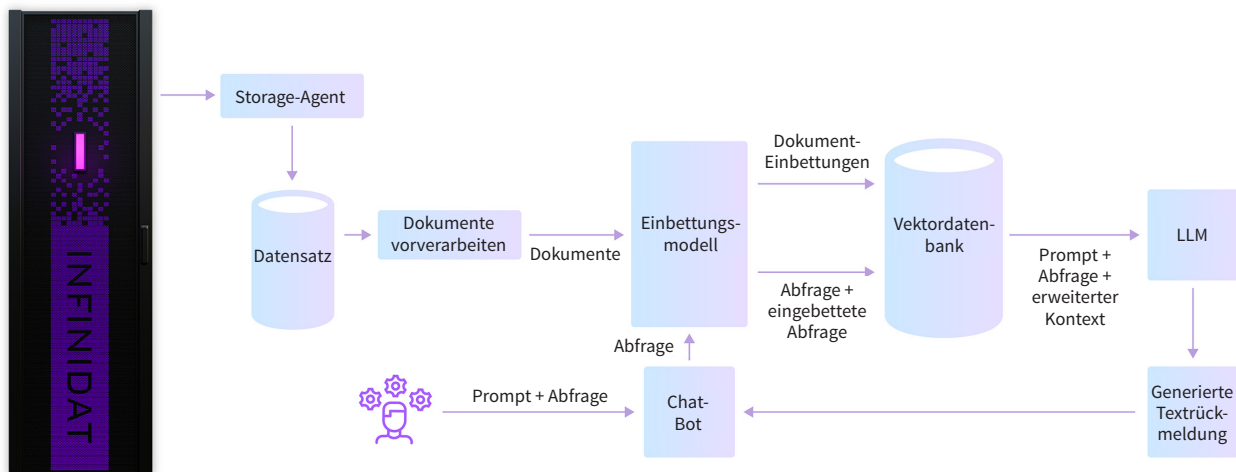




Ein RAG-Workflow kann leicht aus vorhandenen Open-Source-Produkten und Daten erstellt werden, die sich bereits im Rechenzentrum eines Unternehmens befinden. Die Infinidat-Entwickler erstellten eine RAG-Workflow-Architektur, die den Prozess skizziert.

### Infinidat RAG Workflow Architektur<sup>1</sup>



Die RAG-Workflow-Architektur von Infinidat läuft auf einem Kubernetes-Cluster. Anwender, die RAG mit Daten vor Ort betreiben wollen, aber keine GPU-Ressourcen zur Verfügung haben, können eine schnelle und bequeme Lösung über die Cloud nutzen. Unser Ansatz verwendet einen Kubernetes-Cluster als Grundlage für die Ausführung der RAG-Pipeline, was hohe Verfügbarkeit, Skalierbarkeit und Ressourceneffizienz ermöglicht. Mit AWS Terraform vereinfachen wir die Einrichtung eines RAG-Systems erheblich, indem wir die gesamte Automatisierung mit nur einem Befehl ausführen. Der gleiche Kerncode, der zwischen der InfiniBox vor Ort und der InfuzeOS™ Cloud Edition läuft, macht die Replikation zu einem Kinderspiel. Innerhalb von zehn Minuten ist ein voll funktionsfähiges RAG-System auf der InfuzeOS Cloud Edition für die Arbeit mit Ihren Daten bereit.

Die Entwicklung einer RAG-Pipeline ist von Natur aus ein iterativer Prozess, der aktualisiert und gepflegt werden muss, um die Genauigkeit aufrechtzuerhalten. Indem sie mit den neuesten Fortschritten Schritt halten und ihre RAG-Pipeline kontinuierlich verfeinern, können Unternehmen die Genauigkeit und Praxistauglichkeit ihrer KI-gestützten Erkenntnisse erheblich verbessern und die Wettbewerbsvorteile dieser neuen Technologie maximieren.

#### Der InfiniBox-Vorteil für RAG

Bestehende Infinidat-Kunden können davon ausgehen, dass die InfiniBox in ihrem Rechenzentrum bereits über das „Big Data“-Repository an Datensätzen verfügt, um die notwendigen Dokumente für die Einbettung in ein Sprachmodell zu erstellen. Die Auswahl einer Datenbank zum Speichern von Einbettungen (Vektoren) ist für Ihr RAG-System von grundlegender Bedeutung. Glücklicherweise unterstützen die aktuellen Versionen gängiger Datenbank-Engines wie Oracle, Postgres, MongoDB und DataStax Enterprise das Speichern und Abrufen von Vektordaten und sind damit RAG-fähig. Viele unserer aktuellen Kunden setzen ihre InfiniBox in geschäftskritischen Anwendungen ein, die Oracle verwenden.

Unsere solide Datenserviceschicht, die von InfuzeOS bereitgestellt wird, ermöglicht die Steuerung aller Infrastrukturkomponenten durch Software - eine sehr leistungsfähige Fähigkeit. InfuzeOS nutzt unseren innovativen Neural Cache zum dynamischen, vorausschauenden Optimieren der Datenanordnung. Die aktivsten Daten – Hot Data – werden im DRAM-Cache gespeichert, damit der Datenabruf von SSD- oder HDD-Laufwerken auf ein Minimum reduziert wird. InfuzeOS sorgt für eine kontinuierlich optimierte Leistung bei der Skalierung, eine erstklassige niedrige Latenz und eine kompetente Betriebsunterstützung für unterschiedliche moderne Workloads in On-Premises- und Cloud-Umgebungen.

Der Standard für Enterprise Storage

**Fortune Business Insights geht davon aus, dass der Markt für Big-Data-Technologien von 349,40 Milliarden Dollar im Jahr 2023 auf 1.194,35 Milliarden Dollar im Jahr 2032 wachsen wird.**

INFINIDAT

Für einen RAG-Workload bedeutet dies, dass InfiniBox während der Einbettungsphase, in der hochdimensionale Vektoren in einer Datenbank gespeichert werden, eine blitzschnelle Leistung mit sehr geringer Latenz erreicht, die den Prozess nicht behindert.

Wenn ein Nutzer eine Frage stellt (z. B. ChatGPT), wird Ihre Anfrage in eine Einbettung umgewandelt, die sich im gleichen Raum wie die bereits vorhandenen Einbettungen in der Vektordatenbank befindet. Bei der Ähnlichkeitssuche ermitteln Vektordatenbanken schnell die Vektoren, die der Anfrage am nächsten liegen und diese beantworten. Auch hier ermöglicht die extrem niedrige Latenz der InfiniBox schnelle Reaktionen für GenAI-Workloads.

Darüber hinaus bietet die Speicherung der Vektordatenbank auf der InfiniBox einige Vorteile. Das Einlesen der Daten kann ressourcen- und zeitintensiv sein, wobei die Ergebnisse in der Vektordatenbank gespeichert werden. Diese Datenbank kann als Snapshot gespeichert und/oder repliziert werden, sodass eine Vektordatenbank zur Verfügung steht, die in andere RAG-Workflow-Implementierungen integriert werden kann und den Ressourcenbedarf in anderen Implementierungen reduziert.

Sobald die Daten vorverarbeitet und in einem effizienten, schnell abrufbaren Datenspeichersystem (InfiniBox) gespeichert sind, kommt das Large Language Model (LLM) ins Spiel. LLMs geben kohärente und kontextbezogene Antworten. Nach dem Abruf der relevanten Daten aus der Datenbank kombiniert das LLM das abgerufene Wissen mit seinem Verständnis, um eine umfassende Antwort zu geben.

Schließlich bietet die InfuzeOS Cloud Edition für AWS und Azure Cloud-basierte Speichervorgänge, die denen vor Ort ähneln, sodass Speicheradministratoren die gleichen leistungsstarken InfuzeOS-Funktionen des Rechenzentrums in der Cloud nutzen können, um Datensätze schnell auf KI-Ressourcen zu replizieren. Diese Speicherstandardisierung reduziert die Kosten und die Komplexität der Speicherverwaltung in allen Umgebungen.

Der Vorteil von InfiniBox gegenüber RAG liegt in erster Linie darin, dass wir mit Innovationen wie Neural Cache die leistungsstärkste Speicherplattform für diesen Workload anbieten können. Die Lösung von Infinidat kann eine beliebige Anzahl von InfiniBox-Plattformen umfassen und ermöglicht die Erweiterbarkeit auf Speicherlösungen von Drittanbietern über dateibasierte Protokolle wie NFS. Aber vergessen wir nicht, dass wir auch 100-prozentige Verfügbarkeit und Cyber-Recovery garantieren, alles wesentliche Bestandteile einer effizienten KI-Infrastruktur.

Der Einsatz von GenAI und Retrieval-Augmented Generation (RAG) wird wesentlich dazu beitragen, dass KI-Modelle genauer und relevanter werden. Unternehmen können ihre bereits getätigten Investitionen in die InfiniBox-Architektur zusammen mit dem oben beschriebenen RAG-Workflow nutzen, um schnelle und reaktionsfähige KI-Modelle für ihre bestehenden privaten Datensätze zu erstellen, neue Erkenntnisse zu gewinnen und Arbeitsabläufe zu optimieren.



<sup>1</sup> Entwurf und Entwicklung des RAG-Arbeitsablaufs, der Andrew Wang zuzuschreiben ist