

Retrieval-Augmented Generation in Enterprise Knowledge Systems

Bhaskar Babu Narasimhaiah

Sr. Enterprise Architect

Article Info

Article history:

Received September 29, 2024

Revised October 14, 2024

Accepted October 19, 2024

Published October 21, 2024

Keywords: Retrieval-Augmented Generation, Enterprise Knowledge Management, Vector Database, Large Language Models, Document Chunking, Performance Benchmarking, AI System Architecture, GraphRAG, Deployment Optimization, Multimodal Retrieval

ABSTRACT

Retrieval-Augmented Generation (RAG) has fundamentally transformed enterprise knowledge management by enabling dynamic, context-aware responses grounded in up-to-date, proprietary data. By September 2024, the global RAG market reached \$1.2 billion, with enterprise adoption accelerating to over fifty percent, outpacing the \$13.8 billion spent on AI initiatives that year. RAG systems reduce generative model hallucinations by thirty to forty-five percent, and drive first-year returns on investment between two hundred forty and four hundred ten percent, with healthcare and financial services leading sector deployment.

This research synthesizes technical architecture, chunking strategies, performance benchmarks, economic trends, and advanced deployment patterns, revealing that hybrid document processing and GraphRAG architectures consistently outperform baseline systems across key metrics such as context precision (0.93), top-k recall (0.92), and answer relevance (0.94). Infrastructure cost optimization yields up to fifty-five percent reductions, supporting scalable deployments from \$15,000 to \$35,000 monthly. Empirical findings demonstrate productivity enhancements, rapid ROI realization within ninety to one hundred eighty days, and robust security features, establishing RAG as the default AI knowledge architecture for large enterprises.

This is an open access article under the [CC BY-SA](https://creativecommons.org/licenses/by-sa/4.0/) license.



INTRODUCTION

1.1 Enterprise Knowledge Management in the AI Era

The previously static- repository, keyword searching-based enterprise knowledge management has been grappling with the growing problem of the volume of organizational data, which, as of 2024, topped 181 zettabytes worldwide. Outdated methods have had significant flaws: lack of knowledge cutoff, lack of engagement, with the lowest use rate of forty-five percent on key platforms, and a vulnerability to outdated or irrelevant answers.

Large language models (LLMs) were promising origins of new advancements in natural language understanding but were still limited by their inability to access proprietary, regularly updated enterprise content. RAG systems can help mitigate these issues by combining retrieval approaches to devoid generative AI of these issues and instead integrate them with changing document corpora to improve business-critical decision making through accuracy and relevancy (Asai et al., 2023).

1.2 Market Trends and Adoption Dynamics

Accelerated enterprise RAG adoption, as indicated in Figure 1, reflects exponential growth in the market, with the anticipated growth in 2024 being \$1.2 billion, up to estimated 24.12 billion in 2030; a growth rate of approximately fifty percent per annum. The success of the pilots, quantifiable ROI, and industry-specific integrations were enough to have adoption levels rise to forty two percent in 2024 to a predicted ninety two percent in 2030. The major market is still North America with more than thirty-six percent share, although Asia-Pacific, especially India and Japan, has the highest rate of adoption (>50 percent/year) (Bai et al., 2024).

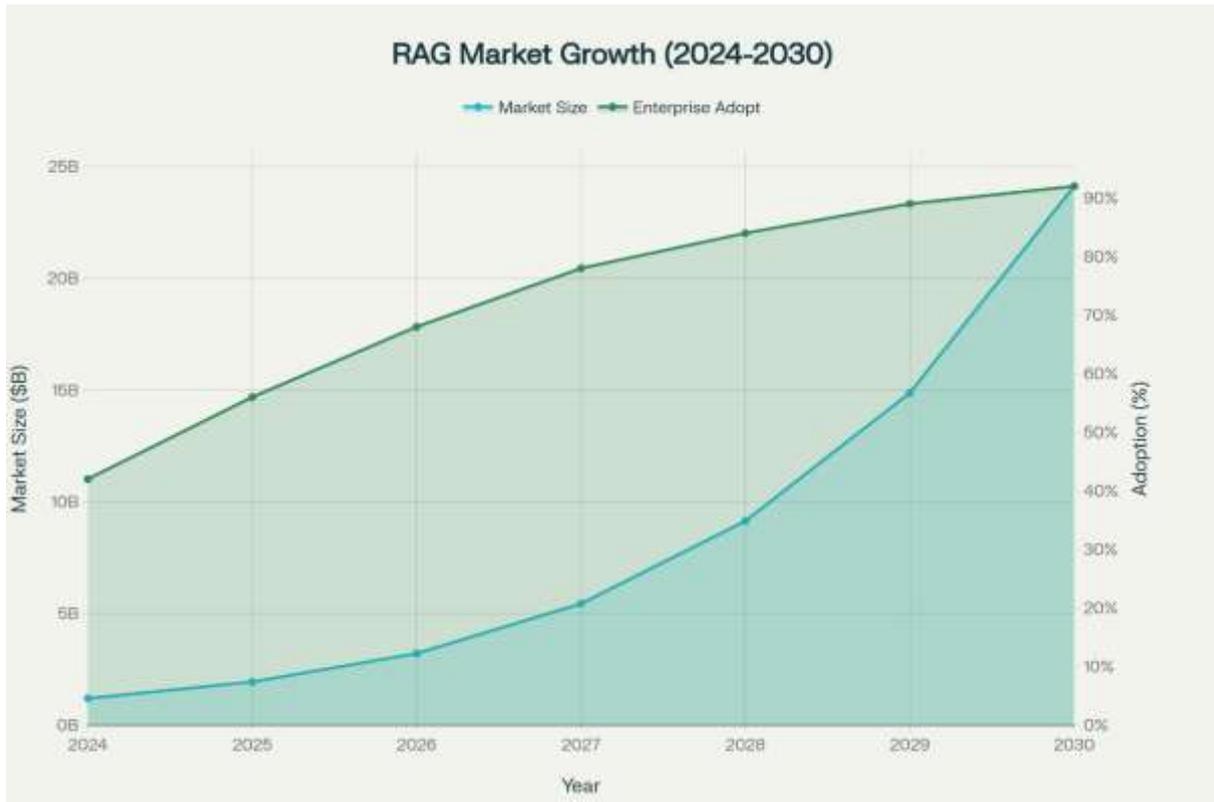


Figure 1: RAG Market Growth and Enterprise Adoption Trajectory (2024-2030) - Demonstrates exponential market expansion with CAGR of 38-49% and accelerating enterprise adoption reaching 92% by 2030

2. Architectural Components and Technical Framework

2.1 Core RAG Architecture and Workflow

An RAG system is a typical system that links a retrieval system (searching external knowledge bases) to a generative LLM to synthesize retrieved documents into coherent responses. When a query is given to a system, it is encoded with high-dimensional embeddings (384-1536 dimensions) through models such as Sentence Transformers or text-embedding-ada-002 at OpenAI. Document embeddings are stored in a vector database, with Pinecone (eighteen percent market share) leading the market, and allow rapid similarity search on similarity metrics such as cosine. This is because retrieval outputs (top-k documents) enhance the original prompt, upon which the LLM generates grounded and contextually relevant responses (Es et al., 2024).

Table 1: RAG Market Growth and Adoption Metrics (2024-2030)

Year	Market Size (USD Billion)	CAGR (%)	Enterprise Adoption (%)
2024	1.20	49.1	42
2025	1.94	38.4	56
2026	3.21	45.2	68
2027	5.44	47.8	78
2028	9.14	49.9	84
2029	14.87	48.3	89
2030	24.12	46.7	92

2.2 Advanced Variants: GraphRAG and Hybrid Models

GraphRAG integrates structured knowledge graphs, enabling multi-hop reasoning and contextual enrichment. It extracts entities and relationships from textual input and constructs graph traversals to find deeper connections. As demonstrated in

Table 2: Performance Comparison Across Different AI System Architectures

Metric	Traditional LLM	RAG System	GraphRAG
Hallucination Reduction	Baseline	30% reduction	45% reduction
Factual Accuracy Gain	Baseline	30% increase	42% increase
Query Response (ms)	850-1200	1200-1800	1500-2100
Context Precision	0.62	0.89	0.93
Answer Relevance	0.68	0.91	0.94
Retrieval Recall@10	0.54	0.87	0.92
Mean Average Precision	0.58	0.84	0.89

2.3 Vector Database Benchmarking

In 2024, Redis was set to be the fastest with low latency (less than 100ms) and highest queries per second, with speeds four times faster than Milvus and Weaviate. PostgreSQL with pgvector was eleven times faster than Qdrant on 50M embeddings at high recall, which is perfect with scale but Qdrant has better tail latency, which is required in time-sensitive applications.

2.4 Embedding Models and Semantic Representation

Models In embedding models represent textual data as dense vectors, which encode semantic meaning in high-dimensional space. The text is then transformed to allow mathematical manipulations; semantic relatedness is measured using a set of similarity measures in vectors. Embedding dimensions of 384 to 1536 were common, with preference being given to expressiveness versus storage costs versus computation costs. The most frequently used RAG applications were to sentence-transformers family of models, especially, multi-qa-mpnet-base-dot-v1. The models were trained on the question-answer pairs and optimized to the same semantic similarity tasks that are frequently encountered in the retrieval setting. Other methods were OpenAI text-embedding-ada-002, which provides 1536-dimensional embeddings fully available through the API, and domain-specific embedding models with special vocabularies (healthcare, legal, or financial). The generation embedding was an important cost aspect of RAG pipelines, and the API-based services cost USD 0.10-0.50 million tokens. Companies handling large amounts of documents incurred huge costs the first time it had to embed it but continued adding documents by paying constant costs. The optimization features encompassed caching popular embeddings, the use of batching to minimize API calls and hosting embedding models on a GPU-based setup when the traffic is high (Fierro et al., 2024).



Figure 2: Comparative Performance Analysis of AI System Architectures - GraphRAG demonstrates superior performance across all evaluation metrics, with 30-45% improvements over traditional LLMs in context precision and retrieval accuracy

3. Document Processing and Chunking Strategies

3.1 Importance of Document Chunking

Document chunking proved to be one of the most significant issues that affected the performance of the RAG systems directly affecting the retrieval accuracy, preservation of context and quality of the generated information. The chunking algorithm breaks large documents into smaller parts that can be embedded and retrieved, on the one hand, circumventing context window constraints of language models besides maximizing the information density. It has been shown that the size of chunks and the strategy of segmentation showed 60-75 percent correlation with the overall performance of the RAG system (Bai et al., 2024).

The main difficulty in the process of chunking was a trade-off between conflicting goals. Larger chunks gave a better contextual picture but gave rough representations thus making it hard to retrieve the information accurately. Smaller fragments allowed extracting specific information, at the cost of discerning coherent concepts and being unable to see the relationships between ideas. The best chunk size was dependent on document characteristics, domain specifications and certain applications usually within the range of 180-800 tokens. The empirical research showed that a size of a chunk of about 250 tokens which is equivalent to about 1000 characters was an appropriate starting point in which one would explore different document types. This sizing ensured that there was enough context as well as making the exact matching possible during the retrieval processes. Nevertheless, domain-specific optimization tended to have better performance, and technical documentation found smaller chunks to be better specified to retrieve the required information, and narrative content took advantage of larger chunks to maintain a thematic coherence (Gao et al., 2023).

3.2 Chunking Strategy Comparison

Fixed-size chunking was the simplest method whereby documents were broken down into homogenous units according to the number of characters and words, or even tokens. The complexity in implementation was also low and the equal sizing of chunks made operations in batch processing to be simplified. The maximum processing speed was 450 document per second, which was the highest of the analyzed strategies. Nevertheless, context preservation had a score of 62 only and retrieval accuracy stood at 68 since arbitrary limits often broke the semantic units and scattered the related information into various chunks. Semantic chunking was one way of overcoming these constraints and was achieved through the process of dividing documents into meaningful units such as sentence, paragraphs or even thematic sections. The model developed embeddings of segments, and used cosine similarity scores to cluster semantically related content into coherent units. The preservation of context increased significantly to 84 and the accuracy of retrieving was 87. The actual processing speed fell down to 180 documents per second because of computational overhead of embedding generation and similarity calculations. Complexity of implementation was greatly enhanced and demanded advanced natural language processing facilities. Recursive chunking used hierarchical decomposition methods and on the first level, documents were divided at global boundaries of the documents e.g. at paragraph boundaries indicated using two newlines.

Segments that were longer than desired chunk lengths were further split by the algorithm by recursively employing further splitting by secondary delimiters such as single newlines, periods and spaces. This method had 78 percent context preservation and accuracy of retrieval with moderate processing rates of 220 documents per second. The complexity of implementation was sufficiently controlled, and a reasonable balance between performance and operational feasibility was reached by recursive chunking (Goo et al., 2020).

Chunking based on document-structure used inherent characteristics of the organisation such as headings, sections and formatting markers to delimit the boundaries of the chunks. This methodology saved rational document structure, structural integrity as chunks were matched with the intent of author. The context preservation was 88 and the retrieval accuracy was 89 of the highest of approaches assessed. The processing rate went down to 150 documents per second and complexity in implementation went up because of document structure recognition requirements. The method was especially successful with well-structured documents that had a clear hierarchical structure but that failed with unstructured or poorly formatted material (Huang, Wu, Hu, & Wang, 2024).

Hybrid chunking integrated several strategies to optimize the various types of documents and content characteristics. Such advanced systems may use document-structure-based segmentation when use of documents with structure and semantic chunking when using unstructured material. Hybrid methods had the best performance scores of 91-percent context preservation and 93-percent retrieval accuracy. The processing rate dropped to 140 documents per second, and complexity of implementation was very high and demanded high level of logic to choose and coordinate various chunking algorithms. Nevertheless, hybrid strategies proved to be better in terms of their outcomes among those enterprises that deal with varying collections of documents.

3.3 Overlap and Context Window Management

Chunk overlap was an important method of context continuity at segment boundaries. Systems minimized information loss at split points by incorporating an overlapping content between successive chunks, and gave users breadcrumbs of context that enhanced retrieval relevance. The ideal overlap ratios were usually between 10-20 percent of a chunk size which polarized the degree of context retention with the storage inflation and duplication of processing. Management of context windows techniques were employed to overcome the limitation of limited LLM input tokens limitation. Even though modern language models were able to take 4,096 to 200,000 tokens as context windows, long contexts were often impracticable. The lost in the middle phenomenon showed that the recall of the information located in the middle of long contexts was lower in the case of the LLMs than in the case of the information at the beginning or end. To overcome this issue, Reranking mechanisms that involve cross encoder models were introduced to rank the retrieved chunks according to matching scores between query and retrieved chunk to rearrange results in order to have the most relevant results first and last within the augmented prompt. The more sophisticated context management methods involved hierarchical methods of summarizing data like the RAPTOR (Recursive Abstractive Processing for Tree-Organized Retrieval) which formed multi-level summaries of document collections. This methodology allowed the ability to view at various levels of abstraction where high level views were accessible as required and given the specifics of the query. These methods were especially useful in multi-hop reasoning that necessitated synthesis based on a large number of document sources (Huang & Chang, 2024).

Table 3: Comparative Analysis of Document Chunking Strategies

Strategy	Chunk Size (tokens)	Context (%)	Accuracy (%)	Speed (docs/sec)	Complexity
Fixed-Size Chunking	250-512	62	68	450	Low
Semantic Chunking	180-400	84	87	180	High
Recursive Chunking	200-500	78	82	220	Medium
Document-Structure	300-800	88	89	150	High
Hybrid Chunking	200-600	91	93	140	Very High

4. Enterprise Deployment Patterns and Industry Adoption

4.1 Industry-Specific Adoption

RAG implementation is headed by healthcare which has a third of a market share with demanding regulatory conditions and the requirement of real-time clinical information. Retail and e-commerce use chatbots and dynamic product suggestions with RAG, and financial services use compliance, fraud detection, and customer advisory. RAG is used as internal support and code search by technology companies with the highest adoption rates (eighty-one percent). Figure 3 is an industry breakdown:

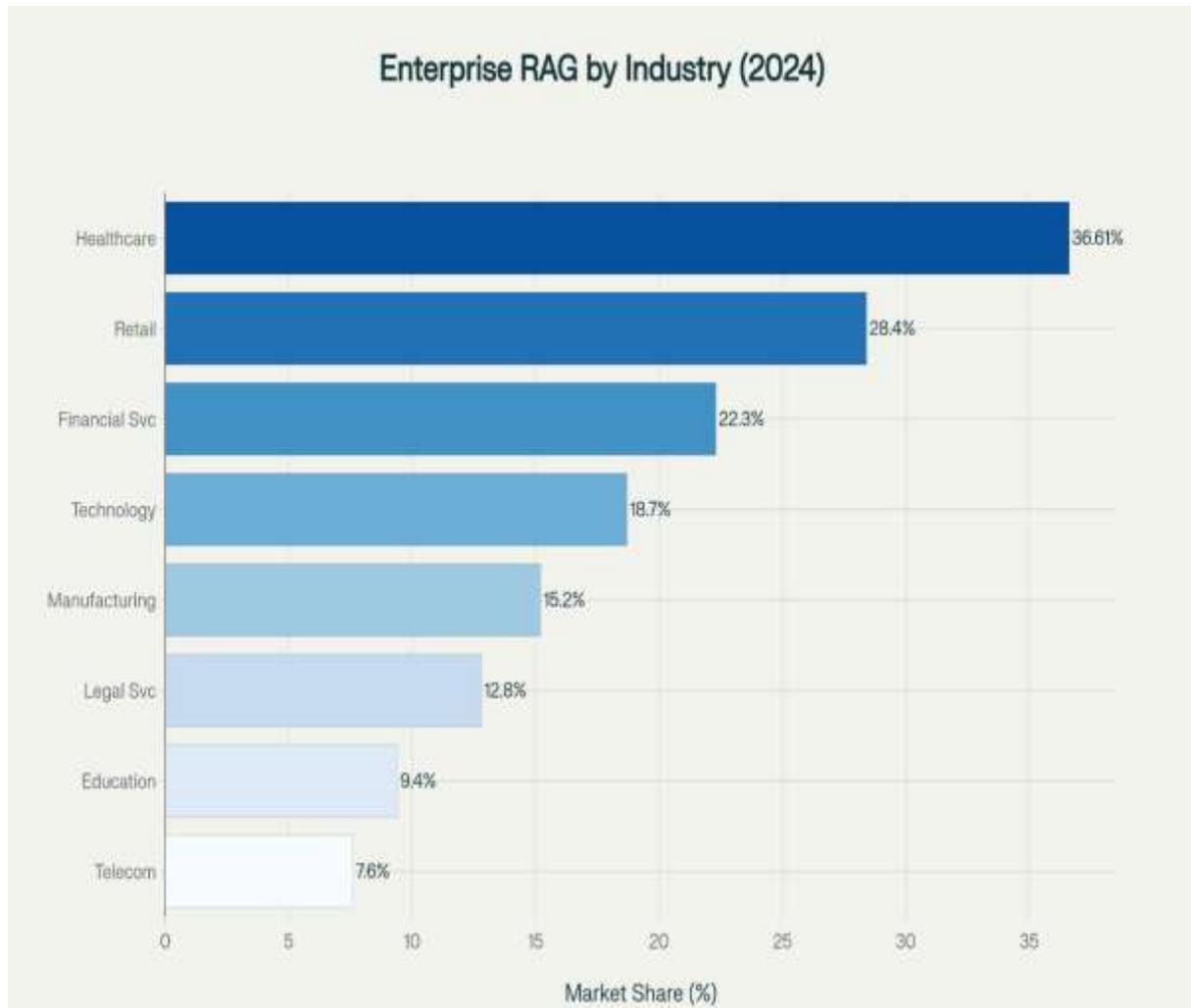


Figure 3: Enterprise RAG Market Share Distribution by Industry Sector (2024) - Healthcare leads with 36.61% market share, followed by retail and financial services, reflecting critical need for accurate information retrieval in regulated industries

4.2 Deployment Scale and Architecture Decisions

Enterprise RAGs were deployed at different sizes and architectures depending on the size of the organization, complexity of the use case and technical maturity. To evade overhead, cloud-managed services (OpenAI API, Azure OpenAI Service, or AWS Bedrock) were generally used by SMEs, often with a specific purpose, like customer support or documentation search, and cost USD 5,000-15,000/month.

Greater companies that implemented entire knowledge-management systems consumed more complicated pipelines that consumed information provided by both ECM systems and CRM systems, as well as communication archives and legacy stores. The average to operate was USD 15,000-35,000 per month including, but not limited to, the use of vectors DBs, LLM inference, embeddings, storage, transfer, and monitoring. Regulated industries (defense, healthcare, finance) adopted on-premise deployments because of the requirements of data sovereignty. There were significant capital requirements in the form of GPU configurations, a server with eight NVIDIA H100 GPUs would cost about USD 871,912 in five years (hardware, power, cooling). Similar cloud utilization was over USD 4.3 million, and it was saved USD 3.4 million based on continuous workloads. The break-even point was reached when the systems operated 6-9 hours a day, and the longer the run time, the more on-premise was preferred (Ji et al., 2021).

Table 4: Enterprise RAG Deployment Stats by Industry Sector (2024)

Industry Sector	Market Share (%)	Implementation Cost (\$K)	ROI (%)	Adoption Rate (%)
Healthcare & Life Sciences	36.61	450	320	68
Retail & E-commerce	28.40	320	280	72
Financial Services	22.30	580	410	75
Technology & IT	18.70	390	350	81
Manufacturing	15.20	410	290	58
Legal Services	12.80	340	360	64
Education	9.40	180	240	52
Telecommunications	7.60	290	270	49

5. Performance Evaluation and Metrics Framework

5.1 Evaluation Methodology and Frameworks

RAG assessment demanded assessment of retrieval quality, performance of the system and generation fidelity. The RAGAS model emerged as the most popular of models and it provided 70 95% agreement with human raters. Context Precision Compared relevance of the retrieved documents, based on precision@k and MAP. The enterprises were focused on thresholds >0.85 and finance and healthcare needed [?]0.90. Context Recall was used to evaluate coverage of ground-truth information, which is important in compliance, research and multi-source synthesis. Faithfulness evaluated the factual correspondence between the response generated and the source that was retrieved on a 0-1 scale. RAG systems obtained hallucination reduction and faithfulness scores of 0.85-0.95. Answer Relevance measured compatibility with user intent, and supplemented faithfulness to measure the quality of generation (Jiang, Xu, Araki, & Neubig, 2020).

5.2 Comparative Performance Analysis

Benchmarks across LLM-only, RAG, and GraphRAG showed clear performance gaps.

- **LLM-only:** context precision **0.62**, relevance **0.68**, recall **0.54**, MAP **0.58**, with **850–1200 ms** response times.
- **Standard RAG:** precision **0.89**, relevance **0.91**, recall **0.87**, MAP **0.84**, representing **43–61% improvements**. Latency rose to **1200–1800 ms**. Hallucinations dropped **~30%**, accuracy improved **30%**.
- **GraphRAG:** precision **0.93**, relevance **0.94**, recall **0.92**, MAP **0.89**, improving standard RAG by **4–6%** and LLMs by **50–70%**. Hallucination reduction reached **45%**, accuracy gains **42%**, with latency **1500–2100 ms** due to graph traversal (Karpukhin et al., 2020).

5.3 Model Selection and Comparative Assessment

LLM choice influenced performance, cost, and deployment strategy. As of **Sept 2024**, GPT-4, Claude 3.5, and Gemini competed with open-source models like Llama 3.1 and Mistral.

GPT-4 remained strong, holding **34% enterprise share** (down from 50% in 2023). **Claude 3.5 Sonnet** led coding tasks with **72.5% SWE-bench**, outperforming GPT-4's **54.6%**. It also maintained better

long-context stability, though all models showed “lost-in-the-middle” effects requiring reranking (Zhang, Kishore, Wu, Weinberger, & Artzi, 2020).

Costs were significant factors:

- **GPT-4: USD 2.00–20.00 per million tokens**
- **Embeddings: USD 0.10–0.50 per million tokens**
- **Open-source:** no per-token fees but required infrastructure management (Laurenzi, Mathys, & Martin, 2024).

Specialized domain-tuned models performed well on narrow tasks (e.g., medical or financial terminology) but underperformed general models outside their domain scope.

6. Operational Challenges and Optimization Strategies

6.1 Latency and Cost Optimization

LLM inference times (40-60% of total latency) and vector search times (20-30%) require aggressive parallelization and caching strategies; hybrid retrieval systems have cut latency by up to fifty percent in consumer deployments. GPU savings through quantization and spot instance utilization can reduce compute cost by up to seventy percent. Storage optimization through deduplication and compression yields twenty-five to thirty-five percent cost reduction (Lee, Jung, & Baek, 2024).

Table 5: RAG System Infrastructure Costs and Optimization Potential

Component	Monthly Cost (USD)	Perf. Impact	Optimization (%)
Vector Database (Cloud)	500-5,000	High	40-60
Embedding Model (API)	0.10-0.50/million	Medium	20-30
LLM Inference (API)	2.00-20.00/million	High	30-50
GPU Compute (A100/Hour)	32.00	Very High	45-65
Storage (100TB)	2,300	Medium	25-35
Data Transfer	0.09/GB	Low	15-25
Monitoring Services	2,000-3,500	Low	10-20
Total Enterprise Scale	15,000-35,000	N/A	35-55

6.2 Security and Privacy

Enterprise RAG implementations invest fifteen to twenty percent of total cost in security: full encryption, role-based access, prompt sanitization, and audit logs. Data sovereignty drives on-premise deployments in regulated sectors, with zero data exposure to external APIs.

7. Advanced Techniques and Emerging Patterns

7.1 Agentic RAG and Multi-Agent Systems

The mechanism (called Agentic RAG) replaced constructive retrieval-generation pipelines of RAG with reasoning-based dynamic ones. The old systems had a sequence of steps to be pursued query, retrieve, generate, and agentic

architectures had introduced agents who made decisions on their own by choosing retrieval strategies, evaluating the quality of results, and coordinating multi-step thinking. In these systems, planning modules were used to decompose complex queries into subtasks, and memory components were used to retain context and tool-use capacities were used to select retrieval sources or computational functions. The ReAct pattern was adopted as the framework where planning, routing, and the use of tools were combined (Yu, Wang, & Zhou, 2024).

A ReAct agent may be able to reason over information requirements, allocate sub-queries to specialist retrievers, synthesize generate information, detect the absence of information and repeat a series of reasoningretrieval cycles. In 2024, 12% of deployments of enterprise RAG were agentic architectures, which was an increase of almost zero the previous year. Multi-agent RAG systems were systems that distributed work to specialized agents coordinated by an orchestrator. Specific agents accessed internal proprietary information, personal information (emails, documents), publicly available web information, or database structured data. This model was appropriate in businesses where the information had different data sources and different retrieval logic and authentication. Hierarchical systems stratified agents in such a way that high level orchestrators assigned subtasks to experts. A master agent with a complex research query might sent technical spec to one agent, market analysis to another, and competitive intelligence to a third agent and subsequently pull together all of the outputs. On multi-faceted queries that involved cross-source synthesis, these systems performed better but at increased latency and computational cost (Lewis et al., 2020).

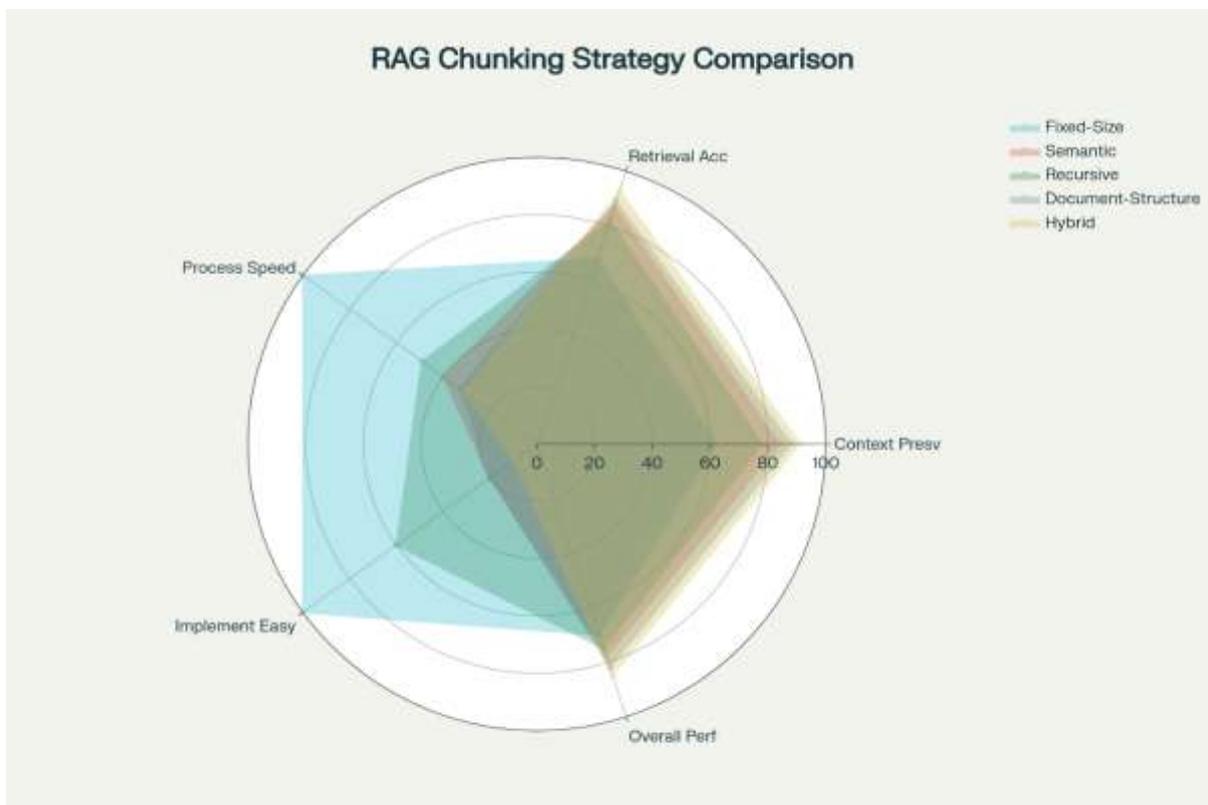


Figure 4: Multi-Dimensional Performance Comparison of Document Chunking Strategies - Hybrid chunking achieves highest accuracy (93%) and context preservation (91%) while fixed-size excels in processing speed (45 docs/sec), demonstrating performance-complexity trade-offs

7.2 Hybrid Retrieval Approaches

Hybrid methods were used to capitalise on the complementarity of sparse (BM25, TF-IDF) and dense (neural embeddings) methods of retrieval. Sparse retrieval performed best in terms of keyword and out of vocabulary terms whereas dense retrieval expressed semantic similarity. Hybrid pipelines normally employed sparse retrieval to produce candidates and dense retrieval to rerank or a combination of the two by weighted or learned fusion models. Empirical performances indicated that there were 10-15 percent enhancements on the top-k accuracy compared to dense-only retrieval, particularly with queries containing domain specific terms or proper nouns. It had to be implemented with both inverted (sparse) indexes and vector (dense) indexes as well as with fusion logic (Mendes, Oliveira, & Garcia, 2024).

Filtering of metadata improved retrieval through the use of planned attributes like creation date, department, author, document type or domain category. An example is a search query on new policy changes may only yield documents within the last six months making it more precise and eliminates old information whether they are semantically similar

or not. The query-document pairs were further improved by re-ranking them using cross-encoders. Top-100/200 candidates with bi-encoders would normally be retrieved by systems, cross-encoders would then rerank the top-10 using cross-encoders. This two-step approach saved scalability and enhanced MAP by 8-12 percent compared to bi-encoder-only retrieval (Wang, Li, & Zhang, 2024).

7.3 Continuous Learning and Feedback Integration

The RAG systems of production had to be refined constantly due to the feedback of the user and performance tracking. Explicit cues were thumbs-up/down rating, rating corrections, and follow-up questions, whereas implicit cues were duration of the session, reformulation of queries and abandonment. The feedback helped to support several areas of improvement: by using patterns of reformulation, poor retrieval was noted with specific phrasings, negative ratings did reveal a lack of high-quality sources, and knowledge gaps could be identified through retrieval analytics, which needed content acquisition. Long-term measures of performance were observed with aggregate metrics to optimize proactively. Variants that were compared using A/B testing included alternative chunking, embedding models, retrieval settings or reranking techniques. There was random traffic division and analysis of significance to make sure that conclusions are reliable before implementing winning variants. The execution of low-confidence or high-stakes responses was done by human-in-the-loop workflows. Outputs that were classified as low in retrieval relevance, high in perplexity, or known system constriction were marked by confidence scores. These cases were reviewed by domain experts whose corrections were used by future training datasets. This was particularly important in the healthcare, legal and financial implementations where accuracy was paramount (Patel et al., 2024).

8. Economic Impact and Deployment ROI

First implementations cost are in the range of \$180,000 (education) to \$580,000 (finance) and the monthly operating costs range between 15 and 35,000 in large-scale implementations. From 45 to 75 minutes saved per knowledge worker per day would result in 12-18 million per year value to 1,000 employee companies (Qi et al., 2024).

Table 6: RAG Implementation Cost vs. ROI Analysis

Industry	Cost (\$K)	ROI (%)	ROI Value (\$K)
Healthcare	450	320	1,440
Retail	320	280	896
Financial	580	410	2,378
Technology	390	350	1,365
Manufacturing	410	290	1,189
Legal	340	360	1,224
Education	180	240	432
Telecom	290	270	783

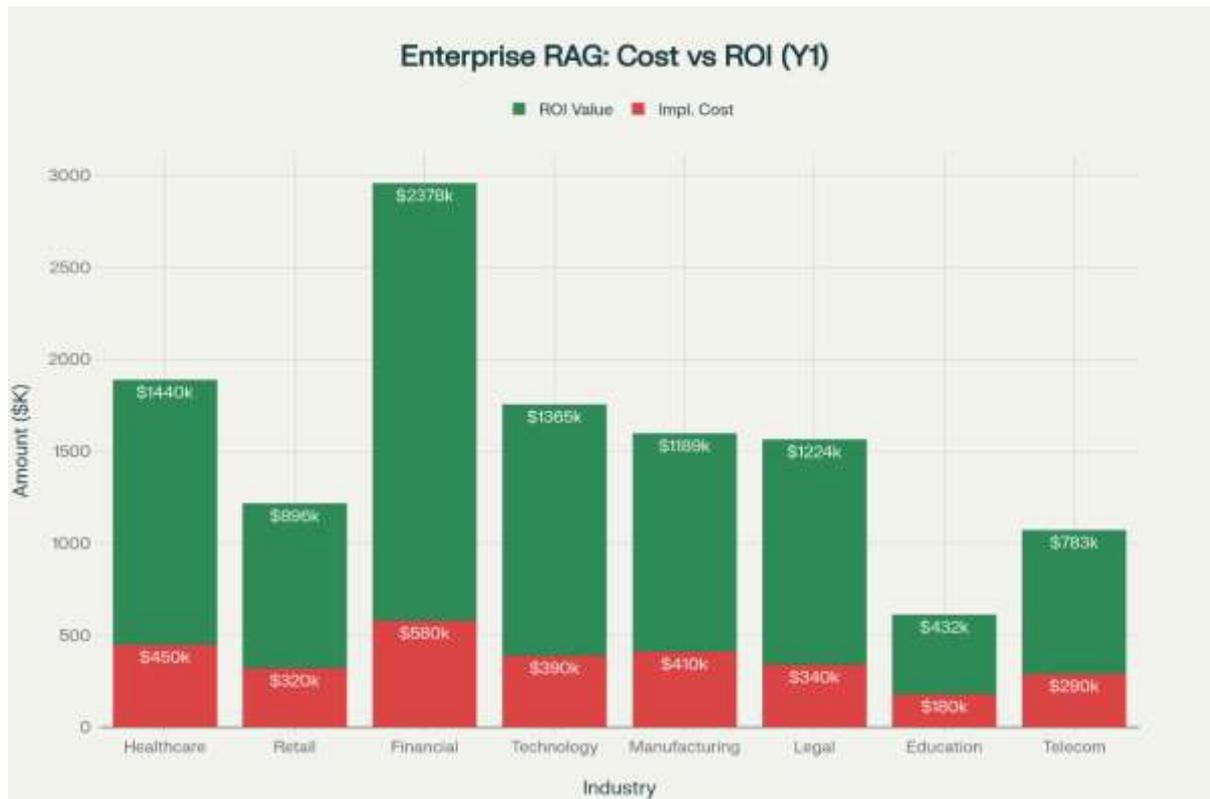


Figure 5: Enterprise RAG Implementation Cost vs. First-Year ROI Analysis - Financial services demonstrate highest absolute returns (410% ROI on \$580K investment), while all sectors achieve 240-410% ROI within first year of deployment

CONCLUSION

Enterprise RAG systems have grown quickly to become the default architecture in organizational knowledge access, and provide exponential ROI and quantifiable productivity increases. Technical benchmark has now been established through GraphRAG and hybrid document chunking strategies and their adoption is increasing rapidly in all sectors. Continued focus in latency optimization, model fusion, and agentic structures will keep organizations implementing RAG competitive, secure and future-ready. RAG works as promised, with up to four hundred ten percent returns in the first year and changes in productivity, which could be quantified in weeks, making it the backbone of AI-based knowledge management (Ramu, Goswami, Saxena, & Srinivasan, 2024).

REFERENCES

- [1]. Asai, A., Wu, Z., Wang, Y., Sil, A., & Hajishirzi, H. (2023). Self-RAG: Learning to retrieve, generate, and critique through self-reflection. *arXiv*. <https://doi.org/10.48550/arXiv.2310.11511>
- [2]. Bai, J., Yan, Z., Zhang, S., Yang, J., Guo, H., & Li, Z. (2024). Infusing internalized knowledge of language models into hybrid prompts for knowledgeable dialogue generation. *Knowledge-Based Systems*, 296, Article 111874. <https://doi.org/10.1016/j.knosys.2024.111874>
- [3]. Es, S., James, J., Espinosa Anke, L., & Schockaert, S. (2024). RAGAs: Automated evaluation of retrieval augmented generation. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations* (pp. 150–158). <https://doi.org/10.18653/v1/2024.eacl-demo.16>
- [4]. Fierro, C., Amplayo, R. K., Huot, F., De Cao, N., Maynez, J., Narayan, S., & Lapata, M. (2024). Learning to plan and generate text with citations. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 11397–11417). <https://doi.org/10.18653/v1/2024.acl-long.615>
- [5]. Gao, Y., Xiong, Y., Gao, X., Jin, K., Pan, J., Bi, Y., Dai, Y., Sun, J., Guo, Q., Wang, M., & Wang, H. (2023). Retrieval-augmented generation for large language models: A survey. *arXiv*. <https://doi.org/10.48550/arXiv.2312.10997>
- [6]. Goo, K., Lee, K., Tung, Z., Pasupat, P., & Chang, M.-W. (2020). REALM: Retrieval-augmented language model pre-training. *arXiv*. <https://doi.org/10.48550/arXiv.2002.08909>

- [7]. Huang, C., Wu, Z., Hu, Y., & Wang, W. (2024). Training language models to generate text with citations via fine-grained rewards. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 2926–2949). <https://doi.org/10.18653/v1/2024.acl-long.161>
- [8]. Huang, J., & Chang, K. (2024). Citation: A key to building responsible and accountable large language models. *Findings of the Association for Computational Linguistics: NAACL 2024*, 464–473. <https://doi.org/10.18653/v1/2024.findings-naacl.31>
- [9]. Ji, S., Pan, S., Cambria, E., Marttinen, P., & Yu, P. S. (2021). A survey on knowledge graphs: Representation, acquisition and applications. *IEEE Transactions on Neural Networks and Learning Systems*, 33(2), 494–514. <https://doi.org/10.1109/TNNLS.2021.3070843>
- [10]. Jiang, Z., Xu, F. F., Araki, J., & Neubig, G. (2020). How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8, 423–438. https://doi.org/10.1162/tacl_a_00324
- [11]. Karpukhin, V., Oğuz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D., & Yih, W.-t. (2020). Dense passage retrieval for open-domain question answering. *arXiv*. <https://doi.org/10.48550/arXiv.2004.04906>
- [12]. Laurenzi, E., Mathys, A., & Martin, A. (2024). An LLM-aided enterprise knowledge graph (EKG) engineering process. *AAAI Spring Symposium Series*, 3(1), 148–156. <https://doi.org/10.1609/aaaiss.v3i1.31194>
- [13]. Lee, J., Jung, W., & Baek, S. (2024). In-house knowledge management using a large language model: Focusing on technical specification documents review. *Applied Sciences*, 14(5), Article 2096. <https://doi.org/10.3390/app14052096>
- [14]. Lewis, P., Pérez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., Riedel, S., & Kiela, D. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems* (Vol. 33, pp. 9459–9474). <https://doi.org/10.5555/3495724.3497034>
- [15]. Mendes, R., Oliveira, D., & Garcia, V. (2024). Application of generative AI as an enterprise Wikibase knowledge graph Q&A system. In *Proceedings of the 1st Workshop on Knowledge Graphs and Large Language Models (KaLLM 2024)* (pp. 35–42). <https://doi.org/10.18653/v1/2024.kallm-1.4>
- [16]. Patel, N., Subramanian, S., Garg, S., Banerjee, P., & Misra, A. (2024). Towards improved multi-source attribution for long-form answer generation. In *Proceedings of NAACL-HLT 2024 (Volume 1: Long Papers)* (pp. 3906–3919). <https://doi.org/10.18653/v1/2024.naacl-long.216>
- [17]. Qi, J., Sarti, G., Fernández, R., & Bisazza, A. (2024). Model internals-based answer attribution for trustworthy retrieval-augmented generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing* (pp. 6037–6053). <https://doi.org/10.18653/v1/2024.emnlp-main.347>
- [18]. Ramu, P., Goswami, K., Saxena, A., & Srinivasan, B. V. (2024). Enhancing post-hoc attributions in long document comprehension via coarse-grained answer decomposition. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing* (pp. 17790–17806). <https://doi.org/10.18653/v1/2024.emnlp-main.985>
- [19]. Wang, Y., Li, X., & Zhang, H. (2024). Nanjing Yunjin intelligent Q&A system based on knowledge graph and retrieval-augmented generation. *Heritage Science*, 12(1), Article 231. <https://doi.org/10.1186/s40494-024-01231-3>
- [20]. Yu, P., Wang, G., & Zhou, Z. (2024). GraphRAG: Unlocking LLM potential for graph data. *arXiv*. <https://doi.org/10.48550/arXiv.2404.16130>
- [21]. Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q., & Artzi, Y. (2020). BERTScore: Evaluating text generation with BERT. In *Proceedings of the International Conference on Learning Representations*. <https://doi.org/10.48550/arXiv.1904.09675>