

Retrieval-Augmented Generation: Bridging Parametric Knowledge and External Evidence in Large Language Models

A. Naziya Inamdar, B. Janhavi Gurav, C. Hitu Bharal, D. Janhvi Indapurkar

¹Naziya Inamdar, AISSMS Institute of Technology, Pune

²Janhavi Gurav, AISSMS Institute of Technology, Pune

³Hitu Bharal, AISSMS Institute of Technology, Pune

⁴Janhvi Indapurkar, AISSMS Institute of Technology, Pune

Abstract—Large language models (LLMs) encode substantial knowledge within their parameters during pre-training; however, this knowledge is inherently static, bounded by a fixed training cutoff, and prone to factual inconsistencies commonly termed hallucinations. Retrieval-Augmented Generation (RAG) has emerged as a principled framework to address these limitations by dynamically conditioning a generative model on externally retrieved documents at inference time. This paper provides a structured examination of the RAG paradigm, covering its theoretical motivation, architectural variants, indexing and retrieval strategies, generation-side integration mechanisms, evaluation protocols, and open research challenges. We further compare RAG against fine-tuning and in-context learning, and outline directions for future investigation. A discussion of real-world deployments in enterprise, medical, legal, and scientific domains is also presented. Our analysis establishes that RAG occupies a distinct and complementary role in the LLM ecosystem, offering a practical balance between factual grounding, updatability, and computational efficiency.

Index Terms—dense retrieval, hallucination mitigation, knowledge grounding, large language models, neural information retrieval, question answering, retrieval-augmented generation.

I. INTRODUCTION

The rapid progression of transformer-based language models has produced systems capable of performing complex reasoning, translation, summarisation, and open-ended dialogue with remarkable fluency. Despite these achievements, a persistent limitation remains: LLMs are trained once on a large but finite corpus, after which their internal knowledge is frozen. Any fact or development arising after the training cutoff is inaccessible to the model unless externally supplied. Moreover, even within the training distribution, models frequently generate confident yet incorrect statements—a phenomenon widely studied as hallucination [1].

Retrieval-Augmented Generation (RAG) offers a compelling response to both problems. Rather than relying exclusively on compressed parametric representations, a RAG system augments the generative process with dynamically retrieved passages drawn from a live or frequently updated corpus. The intuition is straightforward: before generating a response, the model first retrieves relevant reference material—much as a professional consults documentation before rendering a judgment. This hybrid approach preserves the fluency and reasoning capacity of the underlying LLM while grounding its outputs in verifiable, up-to-date evidence [2].

The significance of RAG extends beyond academic interest. Enterprise deployments

increasingly demand systems that can reason over proprietary document collections—legal corpora, medical records, technical manuals, and internal knowledge bases—without the prohibitive cost of continuously retraining foundation models. RAG provides a natural architectural fit for such use cases. The framework is particularly attractive in regulated industries where auditability and citation of sources are compliance requirements.

From a system design perspective, RAG decouples the knowledge store from the model weights, enabling modular updates and multi-source integration. This separation of concerns aligns well with software engineering principles: components can be developed, tested, and optimised independently. As a result, RAG has rapidly transitioned from a research prototype to a production-grade pattern deployed across industry.

This paper offers a structured treatment of the RAG design space, spanning retrieval strategy, fusion mechanism, and evaluation methodology. We cover key architectural choices, discuss open challenges, and situate RAG relative to complementary paradigms such as fine-tuning and tool-augmented generation. The remainder of this paper is organised as follows: Section II provides background on parametric knowledge, transformer models, and neural IR; Section III formally defines the RAG framework; Sections IV and V detail retrieval and generation strategies; Section VI surveys architectural variants; Section VII addresses evaluation; Section VIII surveys real-world applications; Section IX compares RAG with alternative paradigms; and Section X discusses open research challenges before the conclusion in Section XI.

II. BACKGROUND

A. Parametric vs. Non-Parametric Knowledge

A central distinction in knowledge-intensive NLP is between parametric knowledge—facts encoded

in model weights during training—and non-parametric knowledge—information stored explicitly in an external, retrievable memory. Pure neural language models are entirely parametric: the only information they can access is what was implicitly captured during gradient-based optimisation. This confers impressive generalisation but limits factual precision and makes knowledge inflexible to updates.

Non-parametric systems store information in retrievable form and can be updated without retraining. Traditional information retrieval (IR) systems exemplify this approach: a search index can be modified by adding or removing documents without altering the retrieval function. RAG combines both paradigms—a parametric model provides linguistic fluency and reasoning, while a non-parametric store provides factual grounding. This hybrid design philosophy has antecedents in memory-augmented neural networks and knowledge-base-integrated dialogue systems, but RAG uniquely scales this idea to modern, billion-parameter LLMs.

B. Transformer-Based Language Models

Modern LLMs are built on the transformer architecture, which employs self-attention to model long-range dependencies in text. Pre-training on large web-scale corpora yields models that exhibit strong zero-shot and few-shot performance across diverse tasks. However, the quadratic complexity of self-attention with respect to sequence length imposes practical constraints on how much contextual information can be processed in a single forward pass. This constraint becomes especially relevant when injecting retrieved documents into the prompt.

Instruction-tuned variants of these models—such as GPT-4, Claude, Gemini, and LLaMA-based models—have further improved adherence to user-specified tasks and formats, making them more tractable as generators in RAG pipelines. The rapid cadence of model releases has also benefited RAG architectures: as base model

quality improves, the downstream quality of RAG-generated responses often improves in tandem without changes to the retrieval infrastructure.

C. Neural Information Retrieval

Contemporary IR has shifted from purely lexical methods—such as BM25, which scores documents by normalised term frequency and inverse document frequency—toward dense retrieval, wherein queries and documents are encoded as vectors in a shared embedding space. Dense retrieval captures semantic similarity beyond exact token overlap, enabling retrieval of documents on cardiac arrest even when the query contains the phrase “heart attack”. Approximate nearest-neighbour (ANN) algorithms allow this search to scale to millions of documents with sub-linear query latency [3].

The evolution of embedding models—from word2vec and GloVe to contextual sentence encoders such as Sentence-BERT, E5, and GTE—has progressively improved the quality of dense representations. Modern embedding models are fine-tuned on large-scale passage-ranking datasets using contrastive objectives, yielding representations that generalize well across domains. Vector database systems such as FAISS, Qdrant, Weaviate, and Chroma provide efficient ANN search infrastructure that integrates seamlessly into RAG pipelines.

D. Hallucination in LLMs

Hallucination refers to the generation of plausible-sounding but factually incorrect or unsupported text. Studies have categorised hallucination into intrinsic hallucination (contradicting source documents) and extrinsic hallucination (generating unverifiable claims). Hallucination arises partly from the probabilistic nature of language model decoding and partly from memorisation of spurious correlations during training. RAG directly addresses extrinsic hallucination by anchoring responses to retrieved

documents, though it does not fully eliminate intrinsic hallucination, as the generator may still misinterpret retrieved content.

III. THE RAG FRAMEWORK

A. Formal Definition

Let $D = \{d_1, d_2, \dots, d_n\}$ denote a corpus of N documents, and let q denote an input query. A RAG system operates in two stages. In the retrieval stage, a function R selects the k most relevant documents $D_k \subset D$. In the generation stage, a generative model G conditions on both q and D_k to produce an output y , maximising the conditional probability $P(y | q, D_k; \theta)$ over model parameters θ . The key design choices concern the architecture of R , how D_k is presented to G , and whether R and G are trained jointly [1].

Formally, the retriever R can be characterised as a scoring function $s(q, d)$ that assigns a relevance score to each document d given query q , and returns the top- k documents by score. In sparse retrieval, $s(q, d)$ is the BM25 score; in dense retrieval, $s(q, d) = \text{enc}^T(q) \cdot \text{enc}^D(d)$ is the dot product between query and document embeddings. The generator G is typically a sequence-to-sequence or decoder-only LLM that produces y autoregressively given the concatenation of q and retrieved passages.

B. Pipeline Overview

The canonical RAG pipeline proceeds as follows. At index time, source documents are chunked into passages, encoded into dense vectors, and stored in a vector database alongside their metadata. At query time, a user query is encoded by the retrieval model and used to perform ANN search over the pre-built vector index. The top- k retrieved passages are concatenated with the query and provided to the LLM as context. The LLM then generates a response that is grounded in the retrieved evidence and may cite specific passages.

This two-phase structure—offline indexing and online retrieval-generation—provides a clear separation between knowledge curation and inference. The index can be updated incrementally as new documents arrive, without interrupting the online serving infrastructure. Most production RAG systems also incorporate metadata filters that allow retrieval to be scoped to specific document categories, time ranges, or access control groups, enabling fine-grained information governance.

C. Latency and Throughput Considerations

A common concern with RAG systems is inference latency: retrieving from a vector index and loading retrieved context adds overhead compared to a simple LLM query. In practice, modern ANN indexes (e.g., HNSW in FAISS) achieve sub-10 ms retrieval latency for million-scale corpora on commodity hardware. The dominant latency bottleneck in RAG systems is LLM generation time, which scales with the length of the context window. Careful passage selection and compression can mitigate this. Asynchronous architectures decouple retrieval and generation, enabling pipelined execution that further reduces end-to-end latency.

IV. RETRIEVAL AND INDEXING STRATEGIES

A. Sparse Retrieval

Sparse retrieval methods, rooted in the bag-of-words paradigm, assign non-zero weights only to terms present in the document. BM25 remains a highly competitive baseline. Its term-weighting scheme combines normalised term frequency with inverse document frequency, governed by free parameters k_1 and b . Sparse methods are fast, interpretable, and require no GPU infrastructure, making them practical for resource-constrained deployments. Learned sparse retrieval methods such as SPLADE extend BM25 by using neural models to expand the sparse representation with semantically related terms.

B. Dense Retrieval

Dense retrieval encodes queries and documents into continuous vector representations using neural encoders. Models such as DPR (Dense Passage Retrieval) [3], Contriever, and E5 are trained with contrastive objectives that pull relevant query-document pairs closer together and push irrelevant pairs apart in the shared embedding space. Dense retrieval handles vocabulary mismatch gracefully and can be fine-tuned on domain-specific data. Its main drawback is the upfront cost of encoding the entire corpus before deployment. Incremental encoding strategies and model distillation can reduce this cost in practice.

C. Hybrid Retrieval

Hybrid approaches combine sparse and dense signals to obtain complementary strengths. A common strategy uses Reciprocal Rank Fusion (RRF), which aggregates rank positions from multiple retrievers using a smoothed reciprocal formula: $RRF(d) = \sum_r 1 / (k + \text{rank}_r(d))$, where k is a smoothing constant typically set to 60. Hybrid retrieval consistently outperforms either method in isolation on standard open-domain benchmarks [4], balancing the lexical precision of BM25 with the semantic coverage of dense encoders. Learned fusion methods train a linear or neural combiner on labelled data to optimally weight each signal.

D. Chunking and Indexing

Documents must be split into passages before indexing. Fixed-size chunking divides text into segments of a fixed token count with optional overlap. Sentence or paragraph chunking respects natural linguistic boundaries. Semantic chunking groups sentences by topical coherence using embedding similarity. Hierarchical indexing maintains both fine-grained chunks and coarser document summaries, enabling multi-granularity retrieval. Chunk size and overlap involve a fundamental trade-off: smaller chunks improve

retrieval precision but may lack sufficient context for the generator; larger chunks provide richer context but reduce retrieval focus. Empirical studies suggest chunks of 128–512 tokens with 10–20% overlap achieve a good balance for most RAG applications.

E. Re-ranking

A two-stage pipeline first retrieves a large candidate set using a fast bi-encoder, then re-ranks it with a computationally heavier cross-encoder that attends jointly to the query and each candidate document. Re-ranking substantially improves precision at the top of the ranked list, at the cost of additional inference compute. The cross-encoder score $s(q, d) = f_{\theta}([q; d])$ is computed by concatenating the query and document and passing them through the encoder [4]. Models such as MS-MARCO-trained cross-encoders and ColBERT provide strong re-ranking performance. LLM-based re-ranking—where the generator itself scores candidate passages—is an emerging alternative that leverages the generator’s richer world model.

F. Query Expansion and Transformation

Query expansion techniques improve retrieval by enriching or transforming the user’s query before retrieval. Pseudo-relevance feedback appends high-frequency terms from an initial retrieval round to the query. HyDE (Hypothetical Document Embeddings) [6] prompts the LLM to generate a hypothetical answer, which is then used as the retrieval query; since hypothetical answers often share vocabulary with relevant documents, this can substantially improve dense retrieval recall. Step-Back Prompting reformulates the query at a higher level of abstraction, enabling retrieval of background knowledge that supports answering the original question.

V. GENERATION-SIDE INTEGRATION

A. Prompt-Based Fusion

The simplest integration strategy concatenates retrieved passages with the query in the LLM context window. This approach requires no modification to the underlying model and benefits directly from the LLM’s in-context learning ability. The primary limitation is context window length: as the number of retrieved passages or their individual lengths increase, useful information may be truncated and the model’s attention may diffuse across irrelevant content. System prompts can instruct the model to preferentially use retrieved context over parametric knowledge, improving faithfulness.

B. Lost in the Middle

Empirical studies have documented a systematic effect in which LLMs assign disproportionate attention to content at the beginning and end of long contexts, substantially under-attending to information placed in the middle [5]. This has direct practical implications for passage ordering within the prompt—placing the most relevant documents at salient positions rather than in the middle can mitigate this effect and improve answer quality on downstream tasks. Recent positional encoding advances (e.g., RoPE scaling, ALiBi) partially address this issue by improving length generalisation.

C. Attention-Based Fusion

More sophisticated integration approaches fuse retrieved representations at the attention level, allowing the model to attend directly over encoded document vectors rather than prepended text tokens. This reduces the effective input length and allows the generator to selectively query the memory store across multiple layers. However, such methods require architectural modifications and re-training, limiting their applicability to open-weight models where internal activations are accessible.

D. Passage Compression and Summarisation

To mitigate context length limitations, passage compression techniques distil retrieved content into more compact representations before providing them to the generator. Extractive compression selects the most relevant sentences from each passage. Abstractive compression uses a dedicated model to summarise passages, potentially combining information from multiple sources. The RECOMP framework trains a compressor jointly with the generator to produce evidence summaries that maximise downstream answer quality while minimising token count. Compression introduces a risk of information loss, and the compressor must be carefully calibrated to preserve answer-critical details.

E. Citation Generation

A distinctive advantage of RAG over purely parametric generation is the ability to produce answers with inline citations linking to retrieved source passages. Citation generation can be approached as a post-hoc attribution task—identifying which retrieved passages support each claim in the generated output—or as an integrated generation task where the model is trained to emit citation markers alongside answer text. Faithful citation generation remains an active research problem; models may produce plausible-sounding citations that do not accurately support the claimed content.

VI. ARCHITECTURAL VARIANTS

A. Naive RAG

The simplest instantiation—Naive RAG—involves a fixed retriever and a generator that are not jointly trained. Retrieved chunks are inserted verbatim into the prompt and the model generates a response conditioned on them. This approach is straightforward to implement and deploy but is susceptible to retrieval noise: irrelevant or partially relevant passages in the context can degrade generation quality and introduce factual errors. Despite its simplicity, Naive RAG establishes a strong baseline that many advanced

methods struggle to substantially surpass on standard benchmarks.

B. Advanced RAG

Advanced RAG introduces pre-retrieval and post-retrieval enhancements to improve both retrieval quality and generation faithfulness. Pre-retrieval techniques include query rewriting, which reformulates the user's query to better match retrieval vocabulary; HyDE [6], where the LLM generates a hypothetical answer used as the retrieval query; and query decomposition for multi-hop questions. Post-retrieval techniques include re-ranking, passage compression, and selective passage inclusion based on relevance scores. Together, these enhancements substantially improve both retrieval precision and generation accuracy.

C. Modular RAG

Modular RAG decomposes the system into independently swappable components: a query planner, one or more retrievers, a fusion module, and a generator. This modularity enables complex reasoning patterns such as iterative retrieval, where the model retrieves documents, generates an intermediate answer, and retrieves again based on the updated information need. Modular designs also facilitate domain-specific customisation by replacing individual components without retraining the full pipeline [7]. IRCOT (Interleaving Retrieval with Chain-of-Thought) is a prominent example, interleaving retrieval steps with chain-of-thought reasoning to solve multi-step problems.

D. Self-RAG

Self-RAG [8] introduces a mechanism for the model to decide whether to retrieve at all, and to critique retrieved documents for both relevance and factual support. Special reflection tokens are generated inline with the output to signal retrieval decisions and quality assessments. Candidate outputs are scored and filtered according to these

self-assessments. This approach improves both factual accuracy and output fluency by avoiding unnecessary retrieval on straightforward queries where parametric knowledge suffices, and by rejecting low-quality retrieved passages before they influence generation.

E. Graph RAG

Graph RAG [9] extends the document corpus with an explicit knowledge graph, enabling the system to retrieve structured relational information alongside unstructured text. Entity and relation extraction pipelines populate the graph at index time, and graph traversal is combined with dense retrieval to answer multi-hop factual queries more reliably. Community detection algorithms additionally allow the method to generate high-level summaries of large document collections for query-focused synthesis tasks. Graph RAG is particularly effective for domains with rich relational structure, such as biomedical research and legal case law.

TABLE I
Comparison of RAG Architectural Variants

Variant	Joint Train	Iterative Ret.	Comple xity
Naive RAG	No	No	Low
Advanced RAG	No	No	Medium
Modular RAG	Optional	Yes	High
Self-RAG	Yes	Yes	High
Graph RAG	Optional	Yes	High

VII. EVALUATION

A. Retrieval Metrics

Retrieval performance is assessed with three principal metrics. Recall@K measures the fraction of relevant documents appearing in the top-K retrieved set—capturing whether the

correct evidence exists anywhere in the retrieved context. Mean Reciprocal Rank (MRR) computes the average reciprocal of the rank of the first relevant document, rewarding systems that place relevant documents at the top of the list. Normalised Discounted Cumulative Gain (nDCG) is a graded relevance metric that rewards placing highly relevant documents at the top of the ranked list, discounting contributions by log-rank. MAP (Mean Average Precision) averages precision at each rank position where a relevant document is retrieved, providing a comprehensive single-number summary of ranking quality.

B. End-to-End Metrics

End-to-end performance varies by downstream task. Exact Match (EM) and token-level F1 are standard for extractive question answering benchmarks such as Natural Questions and TriviaQA. ROUGE-L and BLEU measure n-gram overlap for summarisation and translation tasks respectively. BERTScore provides semantic similarity via contextual embeddings, capturing paraphrastic adequacy better than surface-level n-gram metrics. FactScore decomposes generated text into atomic claims and verifies each claim against a reference corpus, providing a fine-grained faithfulness measure suitable for knowledge-intensive tasks.

C. RAGAS Framework

The RAGAS framework [10] proposes four metrics specifically designed for RAG evaluation: (i) Faithfulness—whether the generated answer is entailed by the retrieved context; (ii) Answer Relevance—whether the answer addresses the user’s query; (iii) Context Precision—whether retrieved passages contain the necessary answer; and (iv) Context Recall—whether all information required to answer is present in the retrieved set. These metrics can be computed automatically using an LLM-as-judge approach, enabling scalable evaluation without

human annotation. RAGAS's LLM-based scoring is broadly correlated with human judgments, though it inherits the biases of the underlying judge model.

D. Benchmarks

Common open-domain benchmarks include Natural Questions (NQ), TriviaQA, and WebQuestions for factoid QA. Multi-hop benchmarks such as HotpotQA and MuSiQue are especially important for assessing RAG variants with iterative retrieval capabilities. Domain-specific evaluations are conducted on MedQA for medicine, LegalBench for law, and FinanceBench for finance. The KILT benchmark provides a unified evaluation platform for knowledge-intensive tasks with a shared Wikipedia snapshot, enabling controlled comparison across retrieval paradigms.

E. Human Evaluation

Automated metrics, while convenient, do not fully capture user-perceived quality. Human evaluation protocols for RAG systems typically assess fluency (grammatical correctness and readability), factual accuracy (consistency with ground-truth references), attribution quality (whether citations accurately support generated claims), and completeness (whether the answer addresses all aspects of the query). Crowd-sourced evaluation via platforms such as Mechanical Turk or expert annotation is costly but remains the gold standard for assessing RAG output quality in high-stakes domains.

VIII. REAL-WORLD APPLICATIONS

A. Enterprise Knowledge Management

RAG enables employees to query internal document repositories using natural language, grounding answers in current company policies, technical documentation, and meeting records. Frequent document updates are accommodated

by refreshing the vector index rather than retraining the model, making RAG a cost-effective solution for dynamic enterprise knowledge bases. Leading enterprise deployments include customer support automation, internal IT helpdesks, and onboarding assistants that synthesise company-specific knowledge on demand.

B. Medical and Clinical Decision Support

In healthcare settings, RAG systems retrieve evidence from clinical guidelines, drug databases, and research literature to assist clinicians with differential diagnosis and treatment planning. The ability to cite specific source passages preserves interpretability and supports clinical accountability requirements. Systems such as Med-PaLM 2 and clinical RAG pipelines have demonstrated performance approaching specialist-level on medical licensing examinations when augmented with curated medical corpora. Regulatory compliance requires careful curation of the document store and rigorous evaluation of citation accuracy.

C. Legal Research

Legal professionals use RAG to query large collections of case law and statutory text. The ability to retrieve and cite specific passages makes RAG especially well-suited to legal contexts where provenance is essential for professional accountability and compliance. Commercial legal AI systems have adopted RAG architectures to provide real-time access to jurisdiction-specific case law, legislative records, and regulatory guidance, substantially reducing the time required for legal research tasks.

D. Scientific Literature Synthesis

RAG assists researchers by retrieving and summarising papers relevant to a query, enabling rapid literature review across corpora too large for manual inspection. Integration with structured databases of abstracts and citation graphs further

enhances coverage and precision. Systems built on Semantic Scholar's Open Research Corpus and PubMed demonstrate the feasibility of RAG for biomedical literature synthesis, helping researchers identify relevant prior work, extract key findings, and identify open research gaps.

E. Customer-Facing Conversational AI

Many commercial conversational AI products are RAG-powered, retrieving product documentation, FAQ entries, and policy documents to ground their responses. RAG dramatically reduces hallucination rates in customer-facing settings where accuracy is critical to brand reputation. Hybrid architectures that combine RAG with dialogue state tracking allow these systems to maintain conversational context across multi-turn interactions while continuously grounding each response in retrieved evidence.

IX. RAG VERSUS ALTERNATIVE PARADIGMS

A. RAG vs. Fine-Tuning

Fine-tuning encodes domain knowledge directly into model weights and can produce strong domain-appropriate responses without retrieval overhead at inference time. However, fine-tuning is computationally expensive, requires labelled data, and does not easily support updates—incorporating new information requires additional training cycles. Parameter-efficient fine-tuning methods (LoRA, QLoRA) reduce the compute burden but do not eliminate it. RAG, by contrast, is updated by modifying the document corpus without touching model weights, making it far more practical for rapidly changing knowledge domains [1]. In practice, RAG and fine-tuning are often combined: a model is fine-tuned for domain-appropriate style and instruction-following, while factual grounding is delegated to the retrieval system.

B. RAG vs. Long-Context LLMs

Long-context LLMs can, in principle, process an entire document collection within a single forward pass, obviating the need for explicit retrieval. In practice, this approach is computationally prohibitive for large corpora, and attention quality degrades at very long contexts due to positional aliasing and the lost-in-the-middle effect [5]. RAG provides a principled mechanism for selecting the most relevant information subset, maintaining both computational efficiency and retrieval precision. As context windows grow—current frontier models support up to 1 million tokens—the boundary between long-context LLMs and RAG becomes blurred; hybrid approaches that apply RAG for corpus-scale retrieval and long-context LLMs for passage synthesis are increasingly common.

C. RAG vs. Tool-Augmented LLMs

Tool-augmented LLMs can invoke APIs, execute code, and perform web searches at runtime. RAG can be viewed as a specialised form of tool use where the tool is a retrieval engine over a document collection. Both paradigms share the goal of extending LLMs with external knowledge; the distinction lies primarily in the structured, indexed nature of RAG's knowledge store versus the more general and dynamic nature of tool APIs. In practice, RAG and tool use are often combined within agentic architectures, where an orchestrating agent decides whether to retrieve from an internal corpus, call an external API, or use both.

D. RAG vs. Knowledge Graphs

Structured knowledge graphs store facts as typed triples (subject, relation, object) and support formal reasoning via graph traversal and SPARQL queries. Graph-based systems offer high precision on factoid queries expressible in the schema but struggle with natural language questions and domain gaps. RAG complements knowledge graphs by providing unstructured

evidence that covers facts not explicitly encoded in the graph schema. Hybrid systems that query both unstructured vector stores and structured knowledge graphs represent a promising direction for comprehensive factual QA.

X. OPEN RESEARCH CHALLENGES

A. Multi-Hop Reasoning

Multi-hop reasoning is challenging for standard RAG, as many real-world queries require synthesising information across multiple documents. Current RAG pipelines struggle with this unless using iterative or graph-based retrieval. Reliably decomposing complex queries into retrievable sub-questions remains an unsolved issue. Joint optimisation of decomposition, retrieval, and generation for multi-hop queries is an active area of research, with promising early results from methods such as IRCOT and DecompRAG.

B. Retrieval-Generation Alignment

There is no guarantee that the documents ranked highest by the retriever are the most useful for the generator. Joint training of retriever and generator is promising but requires differentiating through a discrete retrieval step, introducing optimisation challenges such as variance in gradient estimates. EMDR² and REALM address this via expectation-maximisation, while RALM uses in-context learning-based alignment. Closing the gap between retrieval relevance and generation utility remains a key open problem.

C. Adversarial Robustness

The retrieval stage introduces a new attack surface—poisoning the document store with adversarially crafted passages can induce targeted model failures or bias responses. Prompt injection via retrieved content is a particular concern in open-web retrieval settings. Defending against corpus poisoning and adversarial retrieval manipulation requires robust

filtering at both the document ingestion and passage selection stages. Certified robustness guarantees for RAG systems remain largely unexplored.

D. Multilingual RAG

Most RAG research has focused on English-language corpora and queries. Cross-lingual retrieval—where the query and document languages differ—and the handling of code-switched or low-resource language documents represent significant unsolved challenges that limit global deployment of RAG systems. Multilingual embedding models such as mE5 and LaBSE provide a foundation, but performance degrades substantially for low-resource languages. Developing robust multilingual RAG evaluation benchmarks is a prerequisite for measuring progress in this area.

E. Evaluation Standardisation

The lack of agreed-upon evaluation protocols makes it difficult to compare results across papers in the literature. Different works use different corpora, chunking strategies, retriever configurations, and generator models, making it nearly impossible to attribute performance differences to specific design choices. Standardised benchmarks, reproducible baselines, and open-source evaluation toolkits—such as BEIR for retrieval and RAGAS for end-to-end RAG—would substantially accelerate community-wide progress on RAG research.

F. Privacy and Data Governance

RAG systems that index sensitive documents introduce privacy risks: the retrieval system may surface confidential information to unauthorised users, or retrieved content may be memorised and reproduced by the generator. Access control at the document and passage level, differential privacy mechanisms for embedding models, and secure multi-party computation for federated RAG are nascent research directions with

significant practical importance in regulated industries.

XI. CONCLUSION

Retrieval-Augmented Generation represents a mature and practically valuable paradigm for augmenting language models with external knowledge. By separating knowledge storage in an updatable retrieval index from language understanding in a pre-trained model, RAG achieves a favourable balance among factual accuracy, updatability, interpretability, and deployment cost. The architectural design space is rich and actively evolving, with advances in dense retrieval, hybrid fusion, iterative reasoning, and self-reflective generation progressively expanding the capability frontier.

Significant challenges remain, particularly around multi-hop reasoning, retrieval-generation alignment, and robust evaluation. As LLMs themselves grow more capable, the relationship between retrieval and generation will continue to evolve—but the fundamental insight that external grounding improves factual reliability shows no sign of diminishing importance. RAG is expected to remain a central component of production NLP systems for the foreseeable future.

X. REFERENCES

- [1] Y. Gao et al., “Retrieval-Augmented Generation for Large Language Models: A Survey,” arXiv preprint arXiv:2312.10997, 2023.
- [2] P. Lewis et al., “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” in Advances in Neural Information Processing Systems (NeurIPS), vol. 33, pp. 9459–9474, 2020.
- [3] V. Karpukhin et al., “Dense Passage Retrieval for Open-Domain Question Answering,” in Proc. EMNLP, pp. 6769–6781, 2020.
- [4] G. Izacard and E. Grave, “Leveraging Passage Retrieval with Generative Models

for Open Domain Question Answering,” in Proc. EACL, pp. 874–880, 2021.

[5] N. F. Liu et al., “Lost in the Middle: How Language Models Use Long Contexts,” Transactions of the Association for Computational Linguistics, vol. 12, pp. 157–173, 2024.

[6] L. Gao et al., “Precise Zero-Shot Dense Retrieval without Relevance Labels,” in Proc. ACL, pp. 1762–1777, 2023.

[7] H. Trivedi et al., “Interleaving Retrieval with Chain-of-Thought Reasoning for Knowledge-Intensive Multi-Step Questions,” in Proc. ACL, pp. 10014–10037, 2023.

[8] A. Asai et al., “Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection,” in Proc. ICLR, 2024.

[9] D. Edge et al., “From Local to Global: A Graph RAG Approach to Query-Focused Summarization,” arXiv preprint arXiv:2404.16130, 2024.

[10] S. Es et al., “RAGAS: Automated Evaluation of Retrieval Augmented Generation,” in Proc. EACL, 2024.