# RAG Interview Questions & Answers

100+ questions with detailed answers

Kalyan KS

### Q1. Explain the requirement of RAG when LLMs are already powerful.

LLMs are powerful, as they are trained on large volumes of data using sophisticated techniques. However, LLMs because of knowledge cutoff (static knowledge), struggle to answer queries related to the latest events or the data not present in their training corpus.

RAG addresses this challenge by retrieving relevant context from external knowledge sources, which allows LLMs to provide accurate responses. This is why RAG is essential for LLM-based applications that need to be accurate. Otherwise, LLMs alone might provide you answers that are incomplete or outdated.

### Q2. Is RAG still relevant in the era of long context LLMs?

RAG is still important even with long context LLMs. This is because long-context LLMs without RAG have three big problems: "lost in the middle,", high API costs, and increased latency.

Long-context LLMs often struggle to find the most relevant information in large contexts, which hurts the quality of generated responses. Furthermore, processing lengthy sequences in each API call results in high latency and high API costs.

RAG addresses these issues by providing the most relevant information from external knowledge sources. So, you still need RAG to get accurate and cost-efficient responses, even with long context LLMs.

### Q3. What are the fundamental challenges of RAG systems?

RAG is powerful, but it has to deal with the following challenges:

Scalability: Searching and retrieving from large, dynamic knowledge sources quickly and efficiently requires a lot of computing power and well-optimized indexing, which can be expensive or take a long time.

Latency - The two-step process (retrieval then generation) can cause delays, making it less suitable for real-time applications without careful optimization.

Hallucination Risk - Even with retrieval, the model might generate plausible but unsupported details if the retrieved data is ambiguous or insufficient.

Bias and Noise - Retrieved content might carry biases, errors, or irrelevant noise from the web or other sources, which can propagate into the output.

## Q4. What are effective strategies to reduce latency in RAG systems?

Caching, embedding quantization, selective query rewriting, and selective re-ranking are some of the ways to reduce RAG latency. Caching stores retrieved results or generated responses to avoid redundant computation. Embedding quantization to lower bit precision reduces memory and computational load, speeding up retrieval.

Selective query rewriting enhances recall and relevance by refining queries prior to retrieval, primarily utilized for complex or ambiguous queries. Selective re-ranking is only used for complicated queries, which cuts down on unnecessary computation for simpler ones.

## Q5. Explain R, A, and G in RAG.

RAG stands for Retrieval-Augmented Generation. The "R" or Retrieval, refers to the process of searching and fetching the most relevant information from external knowledge sources for the given user query.

The "A" or Augmented, involves including the retrieved relevant context in the LLM prompt having the user query and instructions so that the LLM can generate a response based on the provided context.

Finally, the "G" or Generation is the phase during which the generator LLM processes the prompt having instructions, a query, and context to generate a response that is coherent, accurate, and contextually relevant.

## Q6. How does RAG help reduce hallucinations in LLM generated responses?

Without RAG, LLM answers user questions based on what it learned from the training corpus, which may not be up-to-date or complete. This could lead to hallucinated responses, which are answers that sound right but are wrong.

Retrieval-Augmented Generation (RAG) helps cut down on hallucinations in LLM-generated responses by adding an external retrieval system that pulls relevant, factual information from trusted, up-to-date external knowledge sources.

By combining retrieval with generation, RAG ensures that answers are more accurate, contextually relevant, and less prone to fabrications or false information, significantly enhancing the reliability of the output.

### Q7. Why is re-ranking important in the RAG pipeline after initial document retrieval?

The top K chunks fetched by the RAG retriever may have irrelevant chunks ahead of relevant ones. Passing these results directly to the LLM hurts the quality of the answers because LLMs mostly look at the top-ranked chunks that are given as context.

Re-ranking uses cross-encoder models to deeply measure the semantic relevance of query-chunk pairs and then brings relevant chunks ahead of irrelevant chunks. This reduces the noise and helps the generator LLM to generate more accurate and coherent answers.

### Q8. What is the purpose of character overlap during chunking in a RAG pipeline?

In a RAG pipeline, chunk overlap during chunking ensures contextual continuity and prevents loss of information at the boundaries of chunks. This improves the retrieval accuracy and maintains coherence in the text fed to the LLM.

Typically, an overlap of about 10-20% of the chunk size is used to strike a balance between preserving context and computational efficiency in RAG applications.

### Q9. What role does cosine similarity play in relevant chunk retrieval within a RAG pipeline?

Cosine similarity measures how similar the query embedding is to the embeddings of chunks in the vector database. It finds the cosine of the angle between two vectors and provides a score that shows how closely related the query is to each chunk. Higher scores mean that the chunk is more relevant.

This enables the RAG system to retrieve the most relevant chunks for the query, which is then used by the generator LLM to generate accurate answers.

### Q10. Can you give examples of real-world applications where RAG systems have demonstrated value?

AI search engines are a great example of how RAG systems have changed the way people find information online. AI search engines give you accurate, relevant answers by combining information retrieval with generative AI.

For instance, RAG-based AI search platforms like Perplexity AI improve the user experience by fetching the most recent and relevant information from large knowledge bases and then giving it back in the format that the user wants.

### Q11. Explain the steps in the indexing process in a RAG pipeline.

There are four steps in the indexing process of a RAG pipeline: parsing, chunking, encoding, and storing. The parsing step deals with extracting the document content. Then, the chunking step splits the extracted content into smaller pieces called chunks.

The encoding step uses an embedding model to convert chunks into dense numerical vectors called embeddings. Finally, these embeddings are saved in a vector database for efficient search and retrieval.

All these steps in the indexing process are performed offline.

### Q12. Explain the importance of chunking in RAG.

Chunking in Retrieval-Augmented Generation (RAG) is crucial because it breaks down large texts into smaller and semantically coherent segments called chunks. Proper chunking helps to find relevant information efficiently by creating focused chunks that maintain context and avoid irrelevant noise.

Choosing the right chunk size balances detail and context, optimizing both retrieval accuracy and computational efficiency. Ineffective chunking can lead to poor retrieval results and incoherent responses, which makes it a foundational step for successful RAG performance in real-world applications.

### Q13. How do you choose the chunk size for a RAG system?

Choosing the chunk size for a RAG system involves balancing granularity, context completeness, and computational efficiency. Smaller chunks (e.g., 100-200 tokens) allow precise retrieval but may lack sufficient context. Larger chunks (e.g., 500-1000 tokens) provide more context at the cost of increased computational load and potential noise.

The optimal size depends on the use case, document structure, embedding model, and the generator (LLM) model. For example, smaller chunks are suitable for fact-based queries, and more complex queries benefit from larger ones.

### Q14. What are the potential consequences of having chunks that are too large versus chunks that are too small?

Large chunks often mix different topics into one chunk and reduce the chunk's relevance. This can lead to coarse vector representations and less accurate retrieval. Large chunks can also add noise and confuse the model with irrelevant information that isn't important, resulting in a less accurate answer.

Small chunk sizes in RAG systems can lead to fragmented context. This fragmentation often leads to poor retrieval quality because information that is semantically related may be split up into chunks that are not retrieved together. Furthermore, smaller chunks mean that there are more chunks in the vector database, which increases storage costs and slows down the similarity search.

### Q15. Explain the retrieval process step-by-step in a RAG pipeline.

The retrieval process in RAG systems starts by encoding the user query, i.e., converting it into a dense vector representation using an embedding model. This vector representation is then used to search the vector database, which has the embeddings of chunks. Based on the similarity scores, the vector database system returns the most relevant document chunks.

### Q16. What are the key considerations when choosing an LLM for a RAG system?

The key considerations when choosing an LLM for a RAG system are reading-comprehension ability, context window size, and inference speed. Reading-comprehension ability reflects how effectively the model processes the retrieved context to generate accurate responses.

Context window size is crucial, as longer context models enable RAG systems to effectively include more relevant chunks. However, this must be balanced against cost and latency requirements.

Additionally, inference speed, infrastructure compatibility, and licensing terms also play a key role in deployment decisions for real-world RAG solutions.

## Q17. How is the prompt provided to the LLM in a RAG system different from a standard, non-RAG prompt?

The prompt provided to the LLM without a RAG setup includes only the user query and the optional instructions. Here, the LLM generates the response based on its knowledge gained during training.

The prompt provided to the LLM with the RAG setup includes the user query, instructions, and relevant context. Here, the LLM generates the response as per the instructions solely based on the provided relevant context.

## Q18. What are the key hyperparameters in a RAG pipeline?

Chunk size, chunk overlap, embedding dimensionality, retrieval top-k, and retrieval threshold are some of the most important hyperparameters for retrieval in RAG. Temperature and max output length are two important hyperparameters for RAG generation.

The chunk size determines how much text is put into a segment before embedding, influencing the context granularity retrieved. Chunk overlap repeats a set of tokens at chunk boundaries, helping preserve important context across segments. Embedding dimensionality is the vector size used to represent text, which affects retrieval precision and database efficiency.

Retrieval top-k sets the number of most similar chunks returned, directly impacting recall and context diversity in the response. The retrieval threshold is a similarity cutoff that filters retrieved results, ensuring only relevant chunks are selected.

Temperature controls the randomness of generated text, balancing creativity and determinism in model outputs. Max output length limits the number of tokens generated, managing the verbosity and computational cost of responses.

## Q19. What are the popular frameworks to implement a RAG system? Justify your choice of framework.

LangChain, LlamaIndex, and Haystack are the most popular frameworks for RAG implementation. LangChain is great for custom pipelines, and LlamaIndex is great for

efficient document indexing and retrieval. The Haystack framework provides excellent modularity for building RAG systems.

I would recommend LangChain because of its comprehensive ecosystem, extensive documentation, active community support, and flexibility in handling various data sources and LLM integrations.

## Q20. Explain the influence of LLM context window size on RAG hyperparameters.

The size of the LLM context window has a big impact on RAG hyperparameters, like chunk size and the number of chunks that are retrieved. Larger context windows let you feed more retrieved chunks to the LLM, which increases the chance of including more relevant information. This could make the quality of the generated answers better.

But after a certain point, performance gains start to go down because of problems like "lost in the middle" and higher latency.

## Q21. How do you choose values for various LLM inference hyperparameters in a RAG system?

Temperature controls randomness—lower values give more focused and deterministic responses suitable for technical or precise tasks, while higher values make output more creative and diverse.

The max tokens limit the length of the output, making sure that the answers are short or long enough depending on the use case, with a trade-off between completeness and latency. Optimal settings depend on the specific application context and are found through iterative experimentation.

## Q22. Compare reasoning vs. non-reasoning LLMs for RAG systems.

Reasoning LLMs such as GPT-4o1 and DeepSeek R1 are better generators in RAG systems because they have advanced "test-time compute" and chain-of-thought features. These unique abilities allow them to analyze the retrieved information more effectively and do multi-step reasoning to come up with better answers.

But non-reasoning LLMs are still cheaper and faster, which makes them a good choice for many applications. In the end, the choice between reasoning and non-reasoning models depends on the query complexity.

## Q23. What happens with a weak generator LLM in a RAG system?

A weak generator LLM may find it difficult to understand the retrieved context, which could lead to answers that are incomplete or hallucinated. This makes the whole RAG system less useful because the final answers lack coherence and factual correctness, even though the context that was retrieved is good.

So, in a RAG setup, a strong generator LLM is necessary to convert retrieved knowledge into reliable, contextually relevant outputs.

---

## ☕ Support the Author

I hope you found this "RAG Interview Questions and Answers" book highly useful.

I've made this book freely available to help the AI and NLP community grow and to support learners like you. If you found it helpful and would like to show your appreciation, **you can buy me a coffee** to keep me motivated in creating more free resources like this.

👉 **Buy Me a Coffee**

Your small gesture goes a long way in supporting my work—thank you for being part of this journey! 🙏

— Kalyan KS

---

### Q24. How do you handle ambiguous or vague user queries in RAG systems?

Issues with ambiguous or vague user queries in RAG systems include retrieval of irrelevant information, incomplete answers, and increased risk of hallucination due to the lack of specificity.

The most common strategy to handle ambiguous or vague user queries is query rewriting. Query rewriting transforms unclear queries into precise and focused queries, thereby enhancing retrieval quality and leading to more accurate, grounded responses.

### Q25. What are the different query transformation techniques that enhance user queries in RAG?

Different query transformation techniques in RAG include query rewriting, query expansion, query decomposition, and HyDE to enhance retrieval relevance and context precision.

Query Rewriting: Rewrites the initial user query to make it more specific and detailed, boosting retrieval accuracy.

Query Expansion using Step-Back Prompting: Generates a broader, generalized version of the query.

Query Decomposition: Divides complex queries into simpler sub-queries to ensure comprehensive coverage and more precise retrieval for each component question.

HyDE (Hypothetical Document Embedding): Synthesizes a hypothetical answer to the query and uses it as a retrieval query to get more relevant document chunks.

### Q26. What are the pros and cons of query transformation techniques?

Query transformation techniques in RAG systems offer significant advantages, such as improved retrieval accuracy leading to more relevant and contextually accurate responses.

However, their downsides include increased computational cost, added latency, and potential noise from overexpansion. Over expansion risks retrieving noisy or off-topic documents, while complex methods like query decomposition require careful handling to ensure subqueries align with the original intent.

Some strategies may also require substantial prompt engineering and continuous optimization to match diverse query scenarios. Balancing effectiveness and efficiency is critical to avoid diminishing returns.

## Q27. Explain how the HyDE query transformation technique works.

The HyDE (Hypothetical Document Embedding) technique improves RAG retrieval by transforming the user query into a hypothetical answer before embedding it. Rather than directly searching with the query embedding, the HyDE technique utilizes a large language model (LLM) to create a brief, plausible document that could potentially answer the query.

This synthetic document is then encoded into an embedding and used for retrieval, leading to better semantic alignment with actual document chunks in the database. As a result, HyDE enhances retrieval quality, especially for vague or underspecified queries.

## Q28. Explain how the HyPE technique works in RAG.

The HyPE (Hypothetical Prompt Embedding) technique improves retrieval accuracy by addressing the semantic mismatch between user queries and document chunks.

Unlike HyDE, which generates hypothetical answer documents at query time, HyPE precomputes hypothetical questions for each document chunk during the indexing phase. These questions are designed to capture the key concepts in the chunk, transforming retrieval into a "question-to-question" matching process, which reduces latency and improves retrieval.

## Q29. Compare HyPE and HyDE techniques in RAG.

HyDE (Hypothetical Document Embedding) and HyPE (Hypothetical Prompt Embedding) enhance RAG by addressing the semantic gap between user queries and document chunks, but they differ in approach.

Timing: HyPE generates hypothetical questions during indexing, while HyDE generates hypothetical answer documents at query time.

Efficiency: HyPE reduces runtime latency by avoiding LLM calls during retrieval, unlike HyDE, which requires an LLM call per query.

Focus: HyPE focuses on question-question matching, while HyDE focuses on answer-answer matching.

While HyDE is flexible for diverse queries, HyPE's pre-indexed approach is more efficient for real-time applications.

## Q30. To minimize RAG system latency, which pre-retrieval enhancement technique will you choose?

To minimize RAG system latency, I would choose the HyPE (Hypothetical Prompt Embedding) technique. Unlike query transformation techniques such as query rewriting, query expansion, query decomposition, or HyDE, which require LLM calls at query time and increase latency, HyPE precomputes hypothetical questions for document chunks during the indexing phase.

This question-to-question matching approach reduces runtime latency by avoiding real-time LLM calls, making it more efficient for real-time applications while maintaining high retrieval accuracy. By shifting the computational effort to indexing, HyPE ensures faster and more precise document retrieval.

## Q31. What are the different chunk enhancement techniques in RAG?

The different chunk enhancement techniques in RAG are HyPE, Contextual Chunk Header, and Document Augmentation.

HyPE (Hypothetical Prompt Embedding) precomputes hypothetical questions for each document chunk at indexing time, enabling retrieval by question-to-question matching, which improves semantic alignment and retrieval accuracy without adding query-time latency.

Contextual Chunk Header adds relevant contextual information such as document titles or section headings to each chunk before embedding, helping retrieval models understand and rank chunks better when chunk text alone is ambiguous.

Document Augmentation enhances chunks by including additional metadata and enhances retrieval quality.

## Q32. What are the pros and cons of chunk enhancement techniques in RAG?

Chunk enhancement techniques in RAG, such as HyPE, Contextual Chunk Header, and Document Augmentation, improve retrieval accuracy by enhancing semantic alignment, preserving context, and bridging query-document chunk gaps, leading to better generation performance.

HyPE boosts relevance through precomputed question embeddings without query-time latency, Contextual Chunk Header clarifies ambiguous chunks with document or section titles, and Document Augmentation enriches chunks with additional metadata.

However, these methods increase indexing complexity and storage requirements, potentially raising computational costs. Balancing enhanced retrieval quality with resource demands is a key consideration.

## Q33. Explain how the contextual chunk header technique enhances RAG retrieval.

The Contextual Chunk Header technique in RAG enhances retrieval by adding document titles, section headings, or summaries to each chunk before embedding, providing critical context that clarifies ambiguous or isolated chunk content. This additional information helps the retrieval model better understand the chunk's relevance to a query, improving semantic alignment and ranking accuracy.

## Q34. What are some common chunking methods used in RAG?

Common chunking methods used in RAG are fixed-size chunking, recursive chunking, semantic chunking, and agentic chunking.

Fixed-size chunking divides text into uniform segments based on a predefined token or character length, often incorporating overlap to maintain context.

Recursive chunking iteratively splits text using natural separators like paragraphs or sentences to preserve logical boundaries. Semantic chunking groups text based on semantic similarity using embeddings, creating coherent, meaning-based chunks.

Agentic chunking leverages AI agents to dynamically segment text into task-oriented, semantically coherent chunks, often with metadata to enhance retrieval relevance.

## Q35. What are the criteria to choose a specific chunking method in RAG?

The criteria for choosing a specific chunking method in RAG include the nature and structure of the source documents, capabilities of the embedding model, and the specific task or application needs.

For structured or well-formatted data, semantic or agentic chunking ensures logical boundaries and context preservation. The chunk size must balance between being large enough to capture meaningful context and small enough to fit within model constraints for efficient processing.

Task specificity matters since complex tasks may require semantic or agentic chunking for better context and relevance, while simpler cases can use fixed-size chunking.

Ultimately, the chunking method should balance retrieval relevance, context completeness, and computational efficiency.

## Q36. Explain the pros and cons of semantic chunking.

Semantic chunking groups text based on meaning, creating coherent chunks that enhance retrieval relevance and context preservation.

Pros: It aligns chunks with natural topic shifts, improving the quality of retrieved content for complex queries, and reduces information loss across boundaries.

Cons: It is computationally intensive, requiring embedding models. Additionally, it may struggle with highly complex or poorly structured documents where semantic boundaries are unclear.

## Q37. How does the chunking strategy differ when dealing with structured documents (like PDFs with tables and figures) versus plain text documents?

Chunking strategies for structured documents like PDFs with tables and figures differ significantly from plain text chunking due to the need to preserve complex layouts and relationships. For structured documents, the chunking strategy must respect document elements such as tables, figures, headers, and pages to maintain context and semantic meaning.

Agentic and recursive chunking are more suitable for structured documents due to their flexibility in respecting structure and context. Fixed-size and semantic chunking are often better suited for plain text documents where semantic coherence and simplicity are prioritized.

---

### LLM Engineer Toolkit

🤖 This repository contains a curated list of 120+ LLM libraries category wise.



This repository is highly useful for AI/ML Engineers.

---

**Q38. What are the possible reasons for the poor performance of a RAG retriever?**

The possible reasons for the poor performance of a RAG retriever are an outdated or incomplete knowledge base, a weak retrieval model, low-quality embeddings, and lack of domain-specific fine-tuning.

An outdated or incomplete knowledge base prevents the retriever from accessing recent or relevant information, limiting answer accuracy. A weak retrieval model, such as using TF-IDF or BM25 instead of dense vector models, leads to less effective retrieval of relevant context.

Low-quality embeddings reduce the semantic understanding between queries and document chunks, causing mismatches. Lack of domain-specific fine-tuning results in retrieval errors because the embedding model doesn't fully capture the nuances or terminology of the target domain.

**Q39. What happens with a weak retriever in Retrieval-Augmented Generation (RAG) systems?**

A weak retriever in RAG systems leads to the retrieval of irrelevant or noisy document chunks. This can significantly degrade the quality of generated answers, as the RAG generator relies heavily on the retrieved context. The presence of irrelevant or noisy document chunks in the context because of poor retrieval causes the generator model to produce answers that are inaccurate or hallucinated while still appearing fluent.

Therefore, strong retrievers are necessary to provide the most relevant context and ensure factual and relevant outputs in RAG systems.

**Q40. What are the common retrieval approaches used in RAG systems?**

Common retrieval approaches in RAG systems include dense retrieval, sparse retrieval, and hybrid retrieval. Dense retrieval uses embeddings to capture semantic similarity, enabling effective query matching to relevant document chunks. Sparse retrieval relies on traditional methods like TF-IDF or BM25, focusing on keyword-based matching for efficiency.

Hybrid retrieval combines dense and sparse methods to balance semantic understanding and computational speed. These approaches ensure relevant context is retrieved for generating accurate responses in RAG systems.

**Q41. What are some common challenges in RAG retrieval?**

Common challenges in RAG retrieval include ineffective query understanding, scalability issues, context fragmentation, and handling multimodal data. Ineffective query understanding leads to misinterpreting user intent, resulting in irrelevant retrieved document chunks.

Scalability issues arise when large-scale data retrieval slows performance or overwhelms the system. Context fragmentation happens when retrieved chunks lack sufficient context, lowering response quality. Handling multimodal data is challenging due to complexities in integrating text, images, or other formats effectively.

**Q42. What are the key metrics for evaluating retrieval quality in RAG?**

Key metrics for evaluating retrieval quality in RAG systems are precision, recall, Mean Reciprocal Rank (MRR), and Normalized Discounted Cumulative Gain (NDCG). Precision measures the proportion of retrieved document chunks that are relevant, while recall assesses the proportion of relevant document chunks retrieved from the total available.

MRR evaluates the ranking quality by considering the position of the first relevant document chunks, and NDCG accounts for the relevance and ranking of retrieved documents. These metrics collectively ensure the retriever effectively identifies and ranks relevant information.

**Q43. What are embeddings, and how are they utilized in RAG retrieval?**

Embeddings are numerical vector representations of text that capture the semantic meaning and relationships of the data in a high-dimensional space. In Retrieval-Augmented Generation (RAG), embeddings are used to convert both the user query and document chunks into vectors, enabling semantic search by comparing these vectors for similarity.

This process allows RAG systems to retrieve the most relevant and contextually appropriate document chunks from a knowledge base, which are then used as context to generate accurate and grounded responses. Thus, embeddings form the backbone of RAG retrieval by enabling efficient, meaning-driven retrieval beyond simple keyword matching.

## Q44. What are the key considerations when choosing an embedding model for a RAG system?

When choosing an embedding model for a RAG system, key considerations include

(i) the model's domain relevance to ensure it accurately captures domain-specific semantics,
(ii) embedding dimensionality, which balances retrieval precision against computational and storage costs, and
(ii) embedding model performance on the specific dataset to ensure good retrieval quality. This is necessary, as the real-world data often differ from academic datasets.
(iv) Additionally, factors such as embedding model size, API availability, latency, cost implications, and licensing should be considered to align with infrastructure constraints and use case requirements.

Choosing the right embedding model directly impacts the effectiveness and scalability of the RAG system.

## Q45. What is a VectorDB, and how is it utilized in RAG retrieval?

A vector database, or VectorDB for short, is a specialized database designed to store and retrieve high-dimensional vector embeddings. In RAG retrieval, VectorDB is utilized to efficiently perform semantic searches by matching the vector representation of a user query with the closest vectors stored in the database, thereby retrieving the most contextually relevant document chunks.

VectorDBs enable scalable and fast similarity search, which is crucial for the RAG systems.

## Q46. Explain the role of ANN (Approximate Nearest Neighbor) search algorithms in RAG retrieval.

Approximate Nearest Neighbor (ANN) search algorithms play a crucial role in RAG retrieval by enabling fast and scalable search of relevant document chunks within large vector databases. Approximate Nearest Neighbor (ANN) algorithms enable fast search in RAG retrieval by quickly narrowing down the search space to a small subset of candidate vectors instead of scanning all vectors.

This reduces the number of comparisons needed, significantly speeding up retrieval while maintaining good enough accuracy for relevant document chunks matching. This balance of speed and precision is crucial for real-time and large-scale RAG systems.

## Q47. Explain the step-by-step working of ANN algorithms for fast search in RAG retrieval.

ANN algorithms for fast search in RAG retrieval involve four steps namely - Encoding, Indexing, Navigating, Retrieving.

(i) Encoding: Convert document chunks and queries into vector representations.
(ii) Indexing: Organize these vectors into a specialized data structure (like graphs or hash tables) for quick lookup.
(iii) Navigating: Efficiently explore the index to find vectors close to the query without checking all data points.
(iv) Retrieving: Return the closest approximate neighbors that provide relevant information for RAG retrieval.

This approach balances search speed and accuracy, enabling fast retrieval in large-scale RAG systems.

## Q48. What are the typical distance metrics used for similarity search in vector databases, and why are they chosen?

Typical distance metrics used in vector databases for similarity search are Euclidean distance, cosine similarity, and dot product similarity. Euclidean distance measures the straight-line distance between vectors, making it intuitive for geometric closeness. Cosine similarity evaluates the angle between vectors, focusing on their direction (meaning) rather than magnitude.

Dot product similarity considers both magnitude and direction. These metrics are selected based on the data type and the underlying embedding model to ensure effective and accurate retrieval.

## Q49. Explain why cosine similarity is preferred over other distance metrics in RAG retrieval.

Cosine similarity is preferred in RAG retrieval because it measures the angle between vectors, focusing on their direction (meaning) rather than magnitude. This makes it

effective for textual data where the meaning lies more in the direction of the embedding than its length.

Unlike Euclidean distance or dot product, cosine similarity is invariant to vector length, providing stable and interpretable similarity scores. This helps RAG systems retrieve relevant document chunks even when text lengths vary, improving accuracy and consistency in semantic search.

### Q50. Compare keyword-based retrieval and semantic retrieval in RAG systems.

Keyword-based retrieval in RAG systems relies on the exact or partial matching of keywords in a query to fetch relevant document chunks. This offers high precision for queries with specific terms but may miss semantically related information. In contrast, semantic retrieval uses embeddings to understand the meaning behind the query and retrieves conceptually relevant content even when keywords differ.

Combining both methods can balance precision and semantic understanding for effective retrieval in RAG systems.

### Q51. How does hybrid search work in the context of RAG retrieval?

Hybrid search in RAG systems combines keyword-based retrieval and semantic vector search to leverage the strengths of both methods. It uses a weighted formula to balance keyword relevance and semantic similarity scores. This allows precise matching on exact terms while also capturing conceptually related content.

### Q52. When do you opt for hybrid search instead of semantic search?

Hybrid search is preferred over pure semantic search when there is a need to balance exact keyword matches with semantic understanding, especially in scenarios where users require both precision and contextual relevance.

It is ideal for domains where queries may include specific terms, codes, or entities that must be matched exactly, while also benefiting from capturing synonyms or related concepts.

### Q53. How do you balance relevance and diversity when retrieving document chunks for RAG?

The retrieval step in RAG relies on cosine similarity to identify top-k relevant document chunks. However, one downside of this approach is that it can return highly similar document chunks, leading to redundancy.

Balancing relevance and diversity is crucial in RAG retrieval to include contextually important yet diverse document chunks, preventing redundancy and capturing a broader range of perspectives. This balance helps when dealing with complex questions, as different viewpoints or unique insights can improve the answer's quality while still being accurate.

Techniques like Maximal Marginal Relevance (MMR) help to select document chunks that are both highly relevant to the query and diverse from each other, reducing redundancy.

## Q54.  How do sparse embeddings differ from dense embeddings in terms of keyword matching and retrieval interpretability?

Sparse embeddings provide interpretability and excel at exact keyword matching. These embeddings represent text as high-dimensional vectors with many zeros. In this, each dimension corresponds to a specific term or feature, making retrieval results more understandable.

In contrast, dense embeddings are low-dimensional, continuous vectors with mostly non-zero values learned from neural networks, capturing semantic relationships and context beyond exact matches. This makes dense embeddings less interpretable but more effective for retrieving semantically related content where keywords do not exactly overlap.

Thus, sparse embeddings are favored for precise keyword-based retrieval and interpretability, while dense embeddings support richer, context-aware retrieval. Hybrid approaches leverage the strengths of both sparse and dense embeddings to enhance retrieval performance.

## Q55. How can fine-tuning embedding models improve the retriever's performance in RAG?

General embedding models in RAG systems are trained on broad and diverse datasets that capture wide-ranging language patterns. However, they often lack depth in vocabulary and context specific to domains.

Fine-tuning embedding models aligns the embedding space more closely with domain-specific language and context. This allows the embedding model to better represent domain-specific terminology and jargon, which results in more precise and relevant retrieval.

## Q56. Design a retrieval strategy for a RAG system that needs to handle both structured data (knowledge graphs) and unstructured data (text documents) simultaneously.

A retrieval strategy for a RAG system handling both structured (knowledge graphs) and unstructured data (text documents) involves a hybrid approach combining vector-based semantic search with graph-based retrieval techniques.

The system first indexes unstructured text document chunks using vector embeddings for semantic similarity search, while structured data from knowledge graphs is queried using graph traversal methods that leverage explicit entity relationships and schema metadata.

The results from both sources are then fused to ensure factual precision from structured data and contextual richness from unstructured text. This combined approach enhances completeness and reduces hallucinations in generated responses.

## Q57. Discuss the strategies to scale embeddings in RAG retrieval.

To scale embeddings in RAG retrieval, strategies like Matryoshka Representation Learning (MRL) and quantization are highly effective.

MRL enables flexible embeddings by training a single model to produce nested representations of varying sizes, allowing truncation to smaller dimensions (e.g., 64 or 128) with minimal performance loss, achieving up to 14x size reduction and significant retrieval speed-ups.

Quantization reduces memory usage by compressing embeddings into lower-bit formats like float8 or int8. Combining MRL with quantization can yield up to 8x compression, optimizing storage and retrieval efficiency while maintaining high accuracy for large-scale RAG systems.

**Q58. What advantages does quantization offer over dimensionality reduction for scaling embeddings?**

Quantization offers several advantages over dimensionality reduction for scaling embeddings in RAG retrieval. It compresses embeddings by reducing the precision of numerical values (e.g., from float32 to int8 or float8), achieving up to 4x storage reduction with minimal performance loss.

Unlike dimensionality reduction, which may discard important features and degrade accuracy, quantization preserves the full dimensionality of embeddings, maintaining richer semantic information. Additionally, quantization accelerates computation on hardware optimized for lower-precision formats, improving retrieval speed.

This makes it particularly effective for large-scale RAG systems where storage and latency are critical, while dimensionality reduction risks compromising retrieval quality.

**Q59. Explain the pros and cons of quantized embeddings in RAG retrieval.**

Quantized embeddings in RAG systems offer significant benefits such as drastically reduced memory requirements and much faster retrieval speeds. This makes RAG retrieval more efficient and scalable when dealing with large knowledge bases.

However, the trade-off is a slight drop in retrieval accuracy or relevance. Additionally, quantization effectiveness can vary depending on the embedding model.

Overall, quantized embeddings enable cost-effective, high-speed retrieval but require managing a controlled trade-off between resource savings and accuracy.

**Q60. Compare scalar and binary quantization for embeddings in RAG retrieval.**

Scalar quantization in RAG retrieval compresses embeddings by reducing the bit precision (commonly to int8), offering a moderate 4x reduction in memory usage while maintaining a good balance between retrieval accuracy and speed.

Binary quantization, on the other hand, converts embeddings to 1-bit vectors, achieving up to 32x compression and significantly faster retrieval but at the cost of greater accuracy loss.

Overall, scalar quantization suits use cases prioritizing accuracy with some compression, while binary quantization excels in large-scale, speed-critical scenarios where maximal memory efficiency outweighs some loss of precision.

---

## 🚀 **AIxFunda Newsletter (free)**

Join 🚀 AIxFunda free newsletter to get the latest updates and interesting tutorials related to Generative AI, LLMs, Agents and RAG.

✨ Weekly GenAI updates.

📄 Weekly LLM, Agents and RAG paper updates.

📝 1 fresh blog post on an interesting topic every week.



---

## Q61. How does re-ranking differ from the initial retrieval process in RAG?

The initial retrieval process typically uses a bi-encoder that encodes queries and documents independently and then fetches a broad set of candidates quickly.

The re-ranking process reorders the retrieval results by taking the query and each retrieved document chunk as a single combined input, scoring their relevance through deep interaction. This improves the final ranking quality at the cost of higher computational overhead

This two-stage approach balances efficiency and accuracy by separating fast, broad retrieval from slower, more exact reranking.

## Q62. Explain the pros and cons of using re-rankers in RAG.

Re-rankers reorder search results by taking the query and each retrieved document chunk as a single combined input, scoring their relevance through deep interaction within one model pass. This helps to prioritize the most relevant information in the limited context windows in LLMs.

However, re-rankers introduce increased latency and higher computational costs since they perform deep, query-chunk interaction at query time, making them less suitable for real-time or high-traffic applications.

The trade-off between enhanced precision and increased costs makes re-rankers ideal for specialized use cases but less suitable for applications prioritizing speed and cost-efficiency.

## Q63. What are the different types of re-ranker models that can be used in RAG?

The different types of re-ranker models used in Retrieval-Augmented Generation (RAG) are

Cross-Encoder Rerankers: These models jointly encode the query and document chunk pair to produce a highly accurate relevance score, offering nuanced understanding of relationships but with medium computational cost.

Multi-Vector or Late Interaction Models: Such as ColBERT, they encode queries and document chunks separately but perform fine-grained interaction later, balancing efficiency and performance with lower cost.

Large Language Model (LLM) Rerankers: Utilize powerful LLMs to reason about query-document chunk relevance, achieving great accuracy but incurring high computational overhead.

These models vary in their performance and computational cost, and choice depends on the application's accuracy and latency requirements.

## Q64. Compare general re-rankers and instruction-following re-rankers in RAG.

General re-rankers in RAG systems primarily focus on re-ranking retrieved document chunks just based on their semantic relevance to the user query.

In contrast, instruction-following re-rankers go a step further by dynamically adjusting rankings based on additional user-provided instructions such as document recency, source reliability, or metadata criteria.

## Q65. Why is the cross-encoder typically used as the re-ranker rather than the bi-encoder?

The cross-encoder is typically used as the re-ranker rather than the bi-encoder because it processes the query and candidate document chunks together, allowing it to capture intricate contextual interactions and provide more accurate relevance scores.

While bi-encoders encode queries and document chunks separately, enabling fast and scalable retrieval of broad candidate sets, they miss detailed relationships between query-document chunk pairs.

Cross-encoders, though slower and more resource-intensive, excel in precision, making them well-suited for re-ranking a small set of top candidates identified by the bi-encoder. This combined approach balances scalability with accuracy, leveraging bi-encoders for efficient candidate retrieval and cross-encoders for refined final ranking.

**Q66. A RAG system retrieves 20 candidate document chunks but can only fit 5 in the LLM's context window. Without re-ranking, how might this limitation affect response quality, and what specific problems would a re-ranker solve?**

When a RAG system retrieves 20 candidate document chunks but can only fit 5 in the LLM's context window, the limitation can cause the model to miss critical information from the discarded document chunks. Without re-ranking, the top 5 document chunks may not be the most relevant, leading to incomplete or less accurate answers.

A re-ranker solves this by analyzing and scoring all retrieved document chunks based on relevance and contextual alignment with the query, ensuring the most relevant chunks are included in the limited window.

This filtering reduces retrieval noise, enhances coherence, and maximizes the usefulness of the input for the generative model, thereby improving the overall quality of the response.

**Q67. Describe a scenario where a BM25 retrieval might return relevant chunks but in poor ranking order. How would a neural re-ranker specifically address this limitation?**

A typical scenario where BM25 retrieval yields relevant document chunks but in poor ranking order arises when the query uses synonyms or phrases that vary from those in the documents. This is because BM25's exact keyword matching may surface all relevant items, but fail to prioritize those most contextually aligned due to its lack of semantic understanding.

For instance, searching for "car maintenance" might retrieve document chunks about "vehicle upkeep" and "automobile servicing," but BM25 may rank less relevant document chunks higher if they have keyword overlaps rather than semantic closeness. Neural re-rankers explicitly address this by leveraging deep contextual and semantic signals, reordering the candidate set to prioritize document chunks that best match the query's intent and meaning.

**Q68. If your RAG system serves both simple factual queries and complex analytical questions, how would you decide when to bypass the re-ranker for efficiency while maintaining quality?**

To decide when to bypass the re-ranker in a RAG system, queries should be classified based on complexity. Simple factual queries like "What is the capital of France?" require straightforward and well-known answers. Re-ranker can be skipped for simple factual queries, as the initial retrieval is likely to yield highly relevant results.

For complex analytical questions, such as those requiring synthesis or reasoning across multiple chunks, the re-ranker should be used to ensure the most relevant chunks are prioritized.

**Q69. Describe the vector pre-computation and storage strategy in a bi-encoder + cross-encoder pipeline. Why can't cross-encoders pre-compute text representations like bi-encoders can?**

The RAG pipeline leverages bi-encoders for fast retrieval and cross-encoders for the precise reranking of top candidates.

Bi-encoders pre-compute chunk representations by encoding them into fixed-size dense vectors offline and then storing them in a vector database for efficient retrieval.

Cross-encoders, however, cannot pre-compute chunk representations because they jointly encode query-chunk pairs, capturing intricate interactions through attention mechanisms, requiring both inputs at inference time to produce a relevance score.

**Q70. Compare the noise reduction capabilities of re-rankers versus simply increasing the similarity threshold in initial retrieval. When would each approach be more appropriate?**

Increasing the similarity threshold in initial retrieval reduces noise by filtering out less similar chunks but risks missing relevant ones due to embedding limitations. Re-rankers reduce noise by prioritizing relevant chunks by deeply understanding query-chunk relevance.

The choice depends on the trade-off between computational cost and precision requirements. Re-rankers are preferred for high-stakes applications like legal or medical

searches. Increasing the similarity threshold is simpler and faster, suitable for resource-constrained environments.

### Q71. What challenges do re-rankers face regarding computational overhead and latency?

Re-rankers in RAG systems face significant challenges related to increased computational overhead and latency, as each query-chunk pair must be processed. This latency increase can hinder high-throughput environments, making re-rankers computationally expensive compared to initial vector searches and limiting scalability.

### Q72. In real-time applications with strict latency requirements, describe two specific optimization strategies you could implement to reduce re-ranking overhead while preserving most of the quality gains.

Two effective strategies to reduce re-ranking overhead while preserving quality gains in real-time RAG applications are

1) Query classifier: Deploy a query classifier to identify complex or analytical queries, invoking the re-ranker only for these while bypassing it for simple factual queries.

2) Model distillation: Train a smaller, faster re-ranking model to mimic the behavior of a larger, more accurate model, enabling quicker inference with minimal quality loss.

These approaches balance latency and quality by minimizing computational load without significantly compromising the relevance of retrieved results.

### Q73. How would you evaluate the effectiveness of a reranker in a RAG system? Which metrics (e.g., MRR, MAP, NDCG) would you prioritize and why?
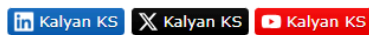
The effectiveness of a re-ranker in a RAG system is best evaluated using ranking metrics that capture how well it prioritizes relevant chunks. Mean Reciprocal Rank (MRR) is key when the focus is on how quickly the first relevant chunk appears, ideal for question-answering scenarios.

Mean Average Precision (MAP) is useful when multiple relevant chunks matter, measuring both precision and ranking quality across results. Normalized Discounted Cumulative Gain (NDCG) excels when relevance is graded rather than binary, rewarding the correct order of highly relevant chunks.

---

## LLM Survey Papers Collection



This repo is highly useful to stay updated with LLM research.

---

**Q74. Explain the difference between Precision@k and Recall@k in the context of RAG. When might you prefer one over the other?**

Precision@k focuses on accuracy of the retrieval by measuring the proportion of the top-k retrieved chunks that are relevant to the query. Recall@k focuses on completeness of the retrieval by measuring the proportion of all relevant chunks that are retrieved within the top-k results.

You might choose Precision@k when you want to ensure high-quality, relevant chunks to reduce noise. On the other hand, you might choose Recall@k when it is crucial to capture as many relevant chunks as possible.

**Q75. Why is MRR unsuitable when there are multiple relevant chunks per query, and how does MAP address this limitation?**

MRR (Mean Reciprocal Rank) considers the rank of the first relevant chunk and disregards the ranks and presence of other relevant chunks. This limitation makes MRR more appropriate for scenarios where a single chunk sufficiently answers the query.

In contrast, MAP (Mean Average Precision) addresses this by averaging the precision across all relevant ranks, accounting for the presence and order of all relevant chunks. Hence, MAP is preferred over MRR for cases where multiple relevant chunks contribute to answering a query comprehensively.

**Q76. Given a retrieval result, show how to manually calculate the MAP@5 (Mean Average Precision at 5). What does MAP reveal about the retrieval system that raw Precision does not?**

To manually calculate MAP@5, list the top 5 retrieved items for each query and note the positions where relevant items appear; then, compute precision at each relevant position (e.g., if the first relevant item appears at rank 2, precision = 1/2) and average these values to get the Average Precision (AP) for that query. Repeat this for all queries and take the mean of their APs for MAP@5.

MAP@5 reveals a retrieval system's ability to rank relevant items higher. In contrast, raw Precision only measures the proportion of relevant items retrieved, ignoring their order. This makes MAP@5 a better indicator of how well the system prioritizes relevance at the top of the result list.

**Q77. If all the relevant chunks are at the very bottom, how would this affect MRR, MAP, and NDCG metrics? Explain each.**

If all relevant chunks are at the bottom of a ranked list for a search query, MRR (Mean Reciprocal Rank) would be low, as it measures the reciprocal of the rank of the first relevant chunk. MAP (Mean Average Precision) would also be low, as it averages precision across all relevant chunks, penalizing late appearances heavily due to increasing denominators in precision calculations.

NDCG (Normalized Discounted Cumulative Gain) would similarly be low, as it discounts the relevance scores of chunks appearing later in the ranking, reducing the cumulative gain.

**Q78. Suppose your RAG retriever gets perfect Recall@10 but low Precision@10. What problems could this cause for the downstream generator?**

Perfect Recall@10 means all relevant chunks are retrieved within the top 10 results. Low Precision@10 indicates many of those retrieved chunks are irrelevant. If a RAG retriever achieves perfect Recall@10 but low Precision@10, the downstream generator receives all the relevant information mixed with much irrelevant content.

This will confuse the generator model and increase the chance of generating off-topic or inaccurate responses.

**Q79. Compare and contrast "order-aware" and "order-unaware" retrieval metrics in RAG, giving examples for each from the set (Precision, Recall, MRR, MAP, NDCG).**

Order-aware retrieval metrics consider the ranking of retrieved items, emphasizing the importance of higher-ranked relevant results. For example, Mean Reciprocal Rank (MRR) and Normalized Discounted Cumulative Gain (NDCG) are order-aware, as MRR evaluates the rank of the first relevant item and NDCG accounts for relevance scores and ranking positions.

Order-unaware metrics focus solely on whether relevant items are retrieved and ignore the order.  Precision and Recall are order-unaware, measuring the proportion of

relevant items retrieved (Precision) and the proportion of relevant items found out of all relevant items (Recall), without considering their order.

## Q80. How would the value of NDCG@k change if all relevant chunks are retrieved but in the reverse order (least to most relevant)?

NDCG@k rewards placing highly relevant chunks at earlier ranks and applies a logarithmic discount to relevance scores at lower positions. So reversing the order pushes the most relevant chunks further down the list—making them less valuable in the NDCG calculation.

While all relevant items are present, their suboptimal positions reduce the overall score since NDCG is sensitive to both the presence and order of relevant items in the top k. The value of NDCG@k will decrease compared to the ideal ranking, but will remain higher than a ranking with irrelevant chunks at the top.

## Q81. What is the significance of Context Precision@K in evaluating a RAG retriever, and how does it differ from standard Precision@k in traditional information retrieval?

The standard Precision@k in traditional information retrieval just measures the proportion of relevant items among the top-k results and ignores the order. Unlike standard Precision@k,  Context Precision@K not only checks whether relevant chunks are retrieved, but also whether they appear at higher ranks in the context.

Context Precision@K ensures useful information is prioritized in the retrieved context which greatly impacts the quality of the generated answer.

## Q82. Why does Context Precision@K use a weighted sum approach with relevance indicators, and how does this better reflect RAG retriever performance?

Context Precision is computed as the weighted sum of Precision@k, normalized by the number of relevant chunks.  Here the weighted sum accounts for both the presence and the rank of relevant chunks in the retrieved context. By multiplying Precision@k with the relevance indicator at each position, the metric rewards cases where relevant information appears earlier, reflecting the importance of ranking quality.

This approach better evaluates RAG retrievers, since in generative settings, not just retrieving relevant chunks but placing them at higher ranks significantly impacts the model's ability to produce accurate answers.

**Q83. Given a retrieval result where relevant chunks appear at positions 2, 4, 6, and 8 out of 10 total chunks, manually calculate the Context Precision@10. What does this score tell us about the retriever's ranking ability?**

To calculate Context Precision@10 with relevant chunks at positions 2, 4, 6, and 8, first assign relevance indicators $v_k=1$ at these ranks and 0 elsewhere. Calculate Precision@k at each relevant rank: at 2, Precision@2 = 1/2 = 0.5; at 4, Precision@4 = 2/4 = 0.5; at 6, Precision@6 = 3/6 = 0.5; at 8, Precision@8 = 4/8 = 0.5. The weighted sum is 0.5+0.5+0.5+0.5=2.

Dividing the weighted sum by the total number of relevant chunks (4) gives Context Precision@10 = 0.5. This score indicates the retriever has moderate ranking ability, retrieving relevant chunks but not consistently ranking them at the very top, thus limiting optimal prioritization of useful context.

**Q84. A RAG system achieves Context Precision@5 = 0.8. What are the possible scenarios that could lead to this score?**

A Context Precision@5 score of 0.8 in a RAG system indicates that not all of the top five retrieved chunks are relevant to the ground truth. This could occur if, for instance, four out of five chunks are relevant ($v_k = 1$) and one is irrelevant ($v_k = 0$), leading to a lower weighted sum of Precision@k when normalized by the total number of relevant chunks.

Another scenario might involve three relevant chunks and two irrelevant ones, with the relevant chunks ranked higher but still resulting in a score less than 1 due to the presence of irrelevant chunks.

**Q85. Explain the possible reasons for a RAG retrieval system with consistently low context precision.**

Context Precision is computed as the weighted sum of Precision@k, normalized by the number of relevant chunks. So low Context Precision@k scores reflect the presence of

high proportion of irrelevant chunks or poor ranking of relevant chunks within the top K results.

This could stem from ineffective query understanding, where the system misinterprets the user's intent, or a poorly designed retrieval algorithm that fails to prioritize chunks matching the ground truth. Additionally, a noisy or low-quality document corpus might contain few relevant chunks, causing irrelevant ones to dominate the retrieved set.

## Q86. Compare Context Recall with traditional information retrieval recall. Why is Context Recall computed using "ground truth claims" rather than simply counting relevant documents?

Context Recall in RAG retrieval differs from traditional information retrieval recall by focusing on the completeness of information through ground truth claims rather than merely counting relevant documents.

While traditional recall counts how many relevant documents are retrieved, Context Recall decomposes the reference answer into individual claims and checks if these specific claims are found in the retrieved context.

This approach ensures a more fine-grained evaluation of whether all necessary pieces of information required to answer the query are present in the context or not.

## Q87. What does context precision measure in a RAG retriever, and how does it differ from context recall?

Context Precision in a RAG retriever measures how well the system ranks relevant chunks of information higher than irrelevant ones within the retrieved context, emphasizing the prioritization of useful data. In contrast, Context Recall assesses the completeness of the retrieved context, evaluating whether all the relevant pieces of information necessary to answer the query are present.

Together, they provide complementary insights: context precision ensures useful information is prioritized, whereas context recall ensures that no important information is missed.

## Q88. In a RAG pipeline, how might context recall impact the completeness of generated answers? Describe a scenario illustrating this relationship.

Context Recall in a RAG pipeline directly impacts the completeness of generated answers by measuring how well the retriever gathers all relevant pieces of information required to answer the user query.

For instance, if a question about a historical event requires multiple claims or facts, a low Context Recall score indicates that some key facts were missed in the retrieved context, leading to incomplete answers.

A high Context Recall ensures the generator (LLM) has access to all necessary information to produce a complete and well-informed response.

## Q89. If your retriever achieves high context precision but low context recall, what types of user queries would likely suffer most?

If a retriever in a RAG pipeline achieves high context precision but low context recall, user queries that require multiple distinct pieces of information or comprehensive coverage are likely to suffer most.

Queries, like complex multi-fact questions or those needing extensive context to answer fully, will suffer because despite the retrieved chunks being relevant (high precision), many essential relevant chunks are missing overall (low recall).

This results in incomplete answers, as important claims or facts are absent, limiting the model's ability to generate a thorough response.

## Q90. In what situations would you prioritize Context Precision over Context Recall in a RAG retriever, and how would this impact the generator's performance?

In situations where precision is critical, such as in high-risk domains like healthcare, finance, or legal applications, prioritizing Context Precision over Context Recall in a RAG retriever is essential. This ensures that only the most relevant and trustworthy information is retrieved, minimizing the risk of including irrelevant or misleading content that could negatively impact the generator's response.

While this may limit the breadth of information (lower recall), it improves the quality and reliability of the generated answers by reducing noise.

**Q91. Describe a scenario where a RAG system might achieve high Context Recall but still produce poor answers. What complementary metrics would you use alongside Context Recall to get a complete picture of retriever performance?**

A RAG system might achieve high Context Recall by retrieving most or all relevant information pieces but still produce poor answers if the retrieved context contains noisy or irrelevant data that confuses the generator.

To get a complete picture of retriever performance, complementary metrics like Context Precision should be used alongside Context Recall. Context Recall along with Context Precision ensures retrieved content is not only comprehensive but is also relevant and well-ranked.

**Q92. If your RAG retriever consistently shows Context Recall scores below 0.6, what are the three potential root causes?**

Context Recall scores below 0.6 mean that the retriever is missing a significant portion of the relevant information required to answer user queries.

The three potential root causes are
(1) an incomplete or outdated knowledge base lacking necessary information,
(2) ineffective embedding model or ranking algorithms causing semantically relevant chunks to be missed, and
(3) poor chunking strategy leading to loss of key information.

**Q93. Why is it important for RAG systems to optimize both context precision and context recall simultaneously? What trade-offs might occur?**

It is important for RAG systems to optimize both context precision and context recall simultaneously. This is because context precision ensures that the retrieved information is highly relevant and ranked appropriately. The Context Recall metric ensures that all necessary information is included in the retrieved context so that the generator can output a complete answer.

The trade-off often arises because increasing recall by retrieving more chunks may introduce irrelevant chunks, lowering precision. At the same time, focusing solely on precision might omit important information, leading to incomplete responses.

Balancing these metrics helps create a RAG system that retrieves relevant content efficiently while covering the query comprehensively, resulting in accurate and thorough generated answers.

## Q94. Explain why Context Relevancy is considered a "reference-free" metric while Context Precision and Context Recall are "reference-dependent." When would you prefer using Context Relevancy over the other two metrics?

Context Relevancy is considered a "reference-free" metric because it evaluates how relevant the retrieved context is to the user's query without needing a reference answer. It measures the proportion of statements in the retrieved context that are relevant to the query.

In contrast, Context Precision and Context Recall are "reference-dependent" as they require a reference answer to determine relevance and completeness of retrieval.

Context Relevancy is preferred when reference answers are unavailable. The Context Relevancy metric offers a way to assess retrieval quality based solely on the query and retrieved context itself. This is useful for real-time scenarios where ground truth may not exist.

## Q95. Describe a scenario where a RAG retriever achieves high Context Relevancy but low Context Precision. What does this imply about the retriever's performance?

A RAG retriever achieves high Context Relevancy but low Context Precision when it retrieves a context where most statements are relevant to the user's query, but the relevant chunks are ranked lower in the retrieved list, overshadowed by irrelevant ones.

For example, if a query about "machine learning algorithms" retrieves a context with many relevant statements but places them after less relevant or noisy chunks, Context Relevancy is high (most statements are query-related), but Context Precision@K is low due to poor ranking of relevant chunks.

This implies the retriever is effective at fetching relevant content but struggles to prioritize relevant chunks over irrelevant ones.

**Q96. Suppose a RAG retriever retrieves all relevant chunks but includes many irrelevant ones, leading to low Context Relevancy. How would you improve the retriever to address this issue?**

A RAG retriever retrieving all relevant chunks along with many irrelevant ones results in low context relevancy scores. This can be addressed by improving the retriever by refining its filtering and ranking mechanisms. Techniques such as enhancing embedding model quality, applying stricter similarity thresholds, or integrating a re-ranking model can help prioritize highly relevant chunks and suppress noise.

Additionally, improving the chunking strategy to create more precise and semantically coherent chunks can reduce irrelevant retrievals. These optimizations ensure retrieved context is both comprehensive and focused on the most relevant information.

**Q97. How does the Faithfulness metric assess the quality of a RAG generator?**

The Faithfulness metric assesses the quality of a RAG generator by measuring how factually consistent the generated response is with the retrieved context. It is computed as the ratio of claims in the response that are supported by the retrieved context to the total number of claims.

A score of 1 indicates all claims are fully supported, reflecting high factual accuracy. A score of 0 shows no claims are supported, indicating complete factual inconsistency. This metric ensures the RAG system generates reliable and contextually grounded responses.

**Q98. Distinguish between Faithfulness and Context Precision metrics in RAG evaluation. Why might a system have high Context Precision but low Faithfulness, and what would this indicate about your pipeline?**

Faithfulness measures how factually consistent a RAG generator's response is with the retrieved context. This metric is computed as the ratio of supported claims to total claims in the response. The Context Precision metric focuses on the prioritization of relevant information by evaluating how well a retriever ranks relevant chunks within the top K.

A system might have high Context Precision but low Faithfulness if the retriever effectively ranks relevant chunks highly, but the generator introduces unsupported or contradictory claims not grounded in the context. This indicates a strong retrieval stage but a flawed generation stage, where the model fails to accurately interpret or utilize the retrieved information.

### Q99. A RAG system has high context precision but low faithfulness. How would you address this?

A RAG system with high context precision and low faithfulness happens when the retriever is selecting relevant chunks accurately, but the generator is producing responses with unsupported claims. To address this, one should focus on improving the generator's grounding and claim verification processes.

Use stronger cross-checking mechanisms like natural language inference models or fact-checking modules against the retrieved context. Additionally, tuning the generation prompts to encourage reliance on the context can help increase faithfulness.

### Q100. Why might a RAG system with perfect Context Recall still fail to produce accurate responses? How does the Faithfulness metric help diagnose this issue?

Context recall evaluates the completeness of the retrieved context in a RAG pipeline. Perfect Context Recall means the retrieved context includes all the ground truth claims. A RAG system with perfect Context Recall may still fail to produce accurate responses. This happens when the generator hallucinates and includes claims unsupported by the retrieved context.

The Faithfulness metric helps diagnose this issue by measuring how many claims in the generated response are factually supported by the retrieved context. A low or moderate faithfulness metric score indicates a less accurate response, i.e., the response includes unsupported claims.

### Q101. Explain how hallucinations in LLMs specifically impact the Faithfulness metric. What techniques could you implement to improve the Faithfulness metric score?

The faithfulness metric measures the proportion of claims in the response that are backed up by context. Hallucinations reduce the faithfulness metric score by

introducing unsupported claims in the generated response. These unsupported claims either contradict or have no basis in the provided context.

To improve Faithfulness scores, techniques such as incorporating natural language inference (NLI) or fact-checking models to verify claims, using prompt engineering to discourage unsupported generation, etc., can be used.

## Q102. How does Response Relevancy differ from Context Relevancy, and why do you need both metrics to properly evaluate a RAG system?

Response Relevancy measures how well a RAG system's generated response aligns with the user's query by calculating the ratio of relevant statements in the response to the total statements. Context Relevancy evaluates the relevance of retrieved context to the query by measuring the proportion of relevant statements in the context.

Both metrics are essential because Context Relevancy ensures the retriever fetches relevant context, while Response Relevancy verifies that the generator produces an answer directly addressing the query.

A RAG system could retrieve relevant context but generate an off-topic response, or vice versa. Hence, evaluating both ensures the entire RAG pipeline—retrieval and generation—performs effectively.

## Q103. The generator's response mentions facts not present in the retrieved context. Describe how faithfulness and response relevancy metrics would be impacted.

The faithfulness metric measures the proportion of claims in the response that are backed up by context. Therefore, the score will decrease if the LLM-generated answer contains unsupported claims.

In the case of the Response Relevancy metric, the score will decrease only if the unsupported facts are irrelevant to the user query. Otherwise, the score will remain high.

This underscores a key difference between these two metrics: Faithfulness metric looks for answer's factual consistency with the context, while Response Relevancy assesses answer's relevancy with the query.

**Q104. How does the Response Relevancy metric help evaluate whether a RAG generator is addressing the user's query effectively?**

The Response Relevancy metric is computed as the ratio of relevant statements in the response to the total number of statements. So, this metric checks the effectiveness of the RAG generator by measuring how well the response aligns with the user query.

A score close to 1 means the answer directly addresses the query with little to no irrelevant content. A score close to 0 means that the answer contains information that is not related to the question.

**Q105. When evaluating RAG generator output, what are the risks of relying solely on response relevancy? How can including the faithfulness metric improve reliability?**

The Response Relevancy metric tells you how relevant the answer is to the user query. But this metric doesn't check if the answer is based on the retrieved context, so it misses factual errors.

The faithfulness metric is the number of supported claims divided by the total number of claims. Adding the faithfulness metric makes the system more reliable by making sure that the claims in the LLM-generated response are supported by the context.

This dual evaluation ensures the RAG system delivers both relevant and factual responses, reducing the risk of misleading outputs.