

## RAG - Game Changer in GenAI Space

**Shekhar Jha**

**Chief Architect | Global Architect Leader | AVP | Integration, Data & Cloud COE,  
USA**

<https://www.linkedin.com/in/shekhar78jha/>

### ***Problem***

Generative AI Challenge - AI hallucinations

It is challenging to rely on Generative AI (GenAI) to provide users and customers with custom, correct and real-time answers to their questions.

Reason - GenAI uses Large Language Models (LLMs) trained on large amounts of publicly available Internet data. Major LLM vendors (OpenAI, Google, Meta) do not train their models on a regular basis as this takes time and money. This means that the data from which the LLMs are based, is never new. Second, LLMs lack access to a massive pool of value-added private information in organizations – the information that enables the individual responses users want.

This is why real-time data scarcity remains one of the biggest challenges for the commercial adoption of GenAI-based tools. For questions to which they do not know the answer, LLMs often provide misleading or irrelevant responses, known as hallucinations, because they rely on the statistical skew of their training data.

### ***Solution***

“By Adopting RAG technique”

RAG technique addresses this challenge. It also helps to avoid AI hallucinations — in which AI generates false or misleading data — and to ensure that answers are based on accurate, current knowledge.

RAG (Retrieval-Augmented Generation) combines the capabilities of LLMs with the ability to extract and integrate knowledge from the external environment.

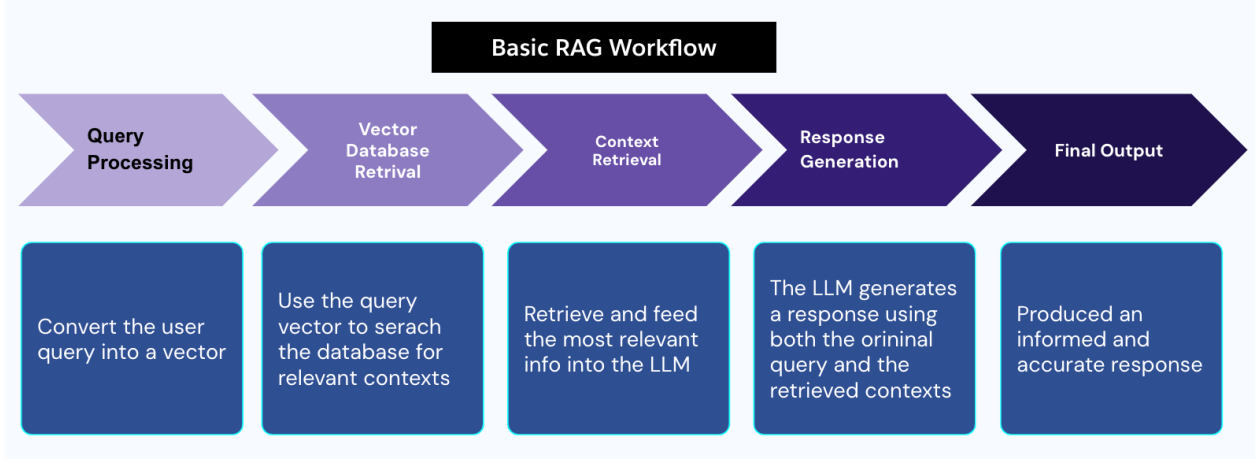
Purpose - It revolutionizes generative AI by giving LLMs the ability to combine private trusted enterprise data (source attribution) with public data, pushing the interface beyond users. This model improves the capacity of LLMs to deliver more meaningful, relevant and accurate answers to user questions. Sources may be included in the output. Even users can search the source documents themselves in case they want further clarification or explanation. This will build more confidence and faith in your generative AI algorithm.

Retrieval: When the user asks a question, RAG first extracts relevant knowledge from outside knowledge sources like databases, websites or company documents.

Augmentation: The information extracted is combined with the existing knowledge and processing resources of the LLM.

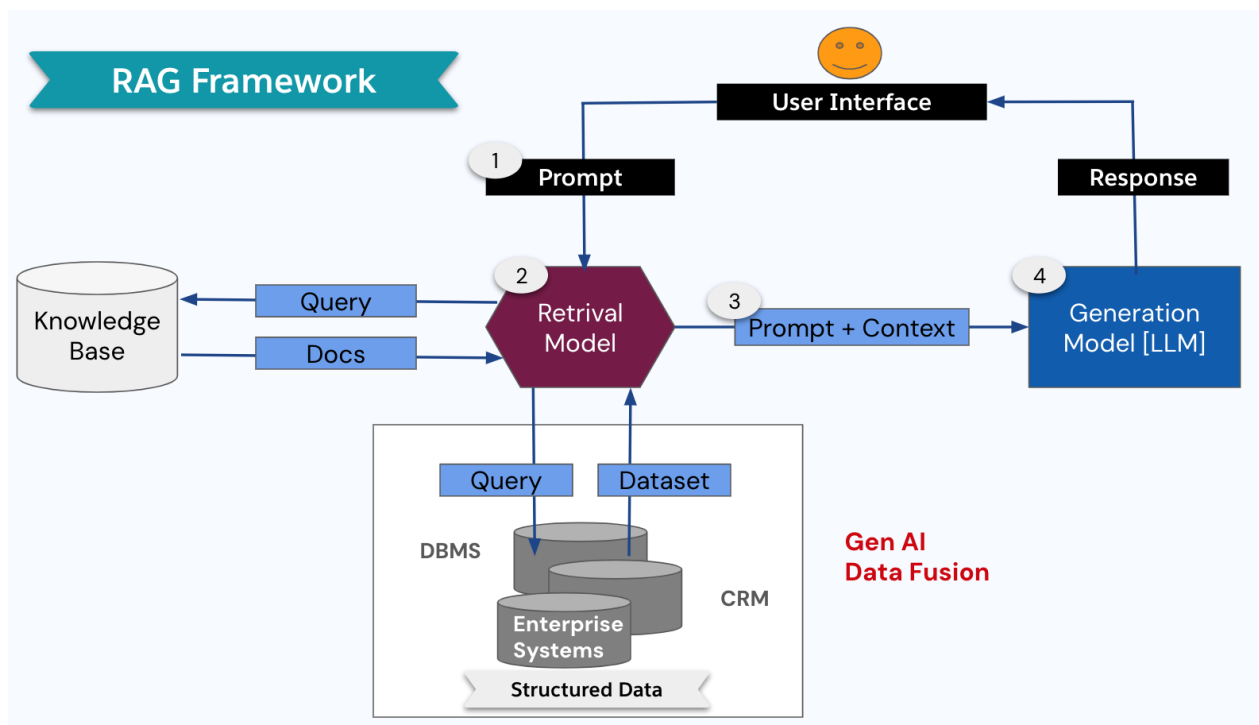
Generation: LLM takes advantage of this enriched knowledge to generate a fully-fledged and informative response to the user's query.

### Basic RAG Workflow



### Understanding the RAG-GenAI Connection

To see it, in the RAG diagram below (source: Gartner), the retrieval model searches, queries, and pre-programs the most relevant data from the prompt that the user provides, converts it into an enhanced contextual prompt, and feeds it into the LLM for analysis. The LLM, in turn, delivers a more precise and specific message to the user.



### The Benefits of RAG for GenAI

RAG-AI connection lets organizations jumpstart their GenAI apps with:

Swift implementation [FAST]

Training for an LLM is a long, costly process. RAG enables faster and cheaper delivery of new, validated data to an LLM, making GenAI easier for businesses to onboard.

**Accuracy [Enhanced Accuracy and Relevance]**

RAG enables LLMs to read and update relevant information in real time, thus ensuring the reliability and timeliness of their responses.

**Up-to-date information [Improved Factual Grounding]**

As responses are based on external knowledge, RAG avoids producing facts that are false or misleading.

**Enhanced Reasoning [Expanded Knowledge Base]**

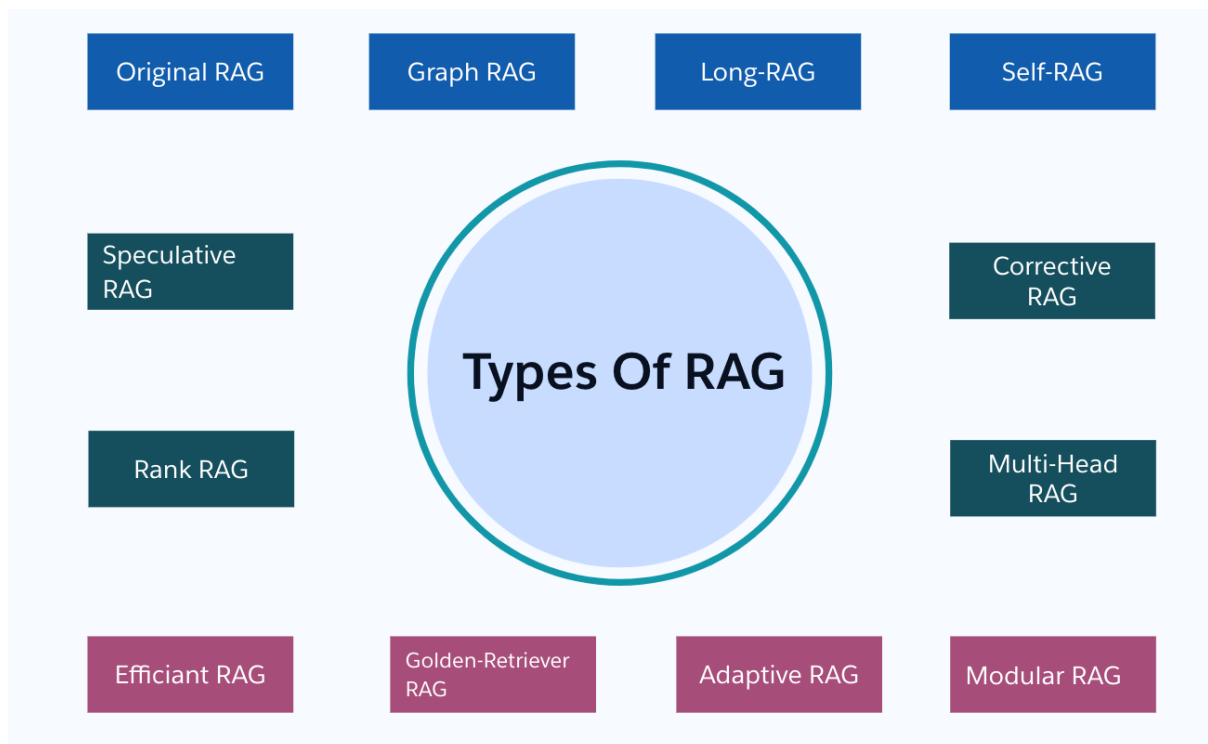
RAG allows LLMs to use much more data than their original training data, allowing them to gain knowledge and abilities.

**Business Specific [Domain-Specific Expertise ]**

RAG is field- or industry-specific, so LLMs can provide domain-focused expert responses to niche-specific questions.

**Customizable [Personalized Responses]**

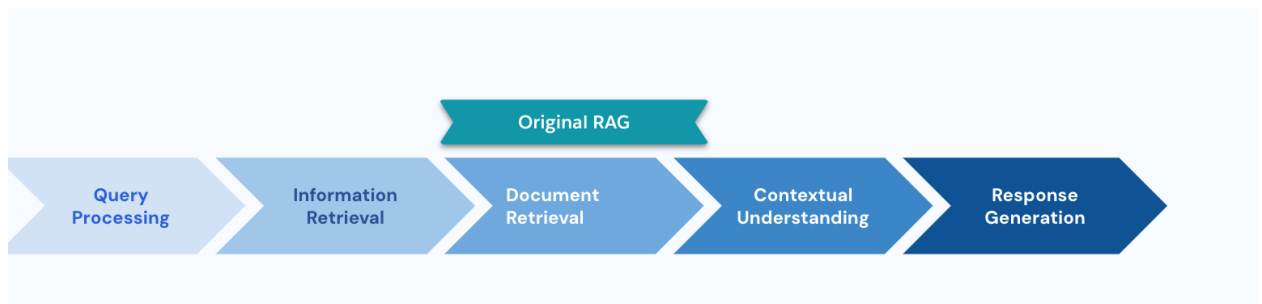
RAG can customize responses according to user preferences, past interactions, etc.

**Exploring Types of Retrieval-Augmented Generation (RAG) Models****1. Original RAG**

Original RAG allows language models to learn from and draw on external information in order to give deeper and more informative answers to difficult questions.

Summary -

- Strengths: Solid foundation, flexible in scope.
- Limitations: Complex computations, simpler capabilities than specialized versions.
- Use Case: Multipurpose question and answer platforms or simple summary text.



**Query Processing:** The user submits a query using natural language. The query is preprocessed to omit the stop words and normalize it.

**Information Retrieval:** The initialized query is turned into a vector, and the query vector is then used to query a vector database populated with vector documents. The vector database efficiently locates the relevant documents based on semantic similarity.

**Document Retrieval:** The best documents are extracted from the vector database.

**Contextual Understanding:** The documents pulled are read to obtain the information. This might involve techniques like:

- Keyword extraction
- Named entity recognition
- Summarization

**Response Generation:** Data obtained is stitched together with the original query to create a full picture. This enhanced context is passed to a language model, which returns a response that is:

- Factual and informative
- Coherent and relevant to the query

### **Benefits:**

- **Better Precision:** With factual information in play, RAG eliminates the possibility of hallucination or misremembering responses.
- **More Factuality:** RAG guarantees that the output responses are based on actual facts.
- **More Relevance:** The system delivers relevant and informative responses by taking into account the context of the query.

### **Limitations:**

- **Dependent on Data Quality:** The quality of data in the vector database directly influences the quality of generated responses.
- **Bias Risk:** If the training data or the vector database is biased, the generated output may reflect those biases.
- **Cost Per Calculation:** RAG can be computationally costly when it comes to big data and multi-question queries.

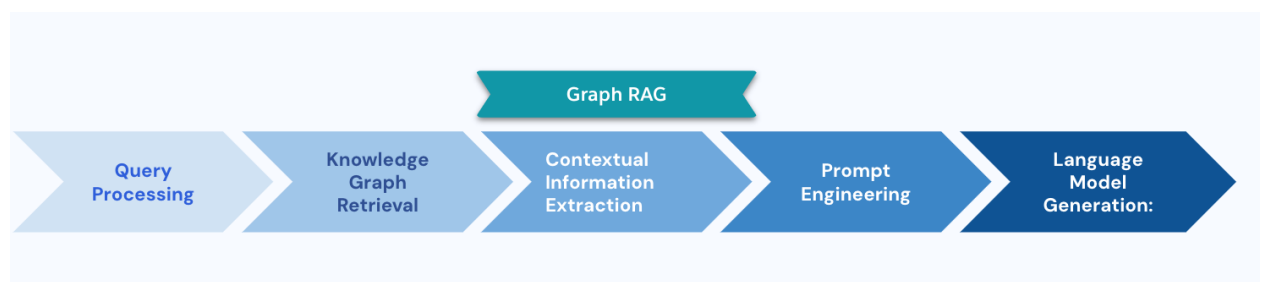
## 2. Graph RAG

GraphRAG is an innovative approach that bridges the gap between RAG and Knowledge Graphs (KGs). It allows Large Language Models (LLMs) to think in terms of more detailed information and give accurate and more meaningful responses.

Combining the efficacy of LLMs with the knowledge structure of KGs, GraphRAG opens up a new window for AI-driven applications.

### Summary -

- Advantages: Improved contextualization via entity relations, ideal for structured data and knowledge graph generation.
- Low-lights: Hard to set up, difficult to scale for big graphs or random data.
- Use Case: Identifying new drugs by connecting chemical molecules, proteins, and diseases.



**Query Processing:** The query is read by the user to determine the key entities and relations.

**Knowledge Graph Retrieval:** The retrieval engine queries the KG for matching nodes and edges. It can hop over a number of hops in search of higher associations.

**Contextual Information Extraction:** The contextual information, such as entities, relationships, and text descriptions are extracted from the KG.

**Prompt Engineering:** The extracted data is utilized to build a comprehensive prompt for the language model. The prompt can include:

- The original query
- Relevant entities and their descriptions
- Relationships between entities
- Specific questions to be answered

**Language Model Generation:** The language model analyzes the command and returns a response. It uses the organized knowledge to answer questions more accurately and in depth.

### Benefits:

- **Factual Precision:** With responses anchored to formal knowledge, GraphRAG mitigates the possibility of hallucinations.
- **Enhanced Reasoning:** GraphRAG allows LLMs to reason through large information by tracing relationships in the KG.
- **More Contextual Understanding:** Rich contextual knowledge provided by the KG enables more thoughtful and sophisticated responses.
- **GraphRAG can explain its answers** by reasoning their way through the KG.

**Limitations:**

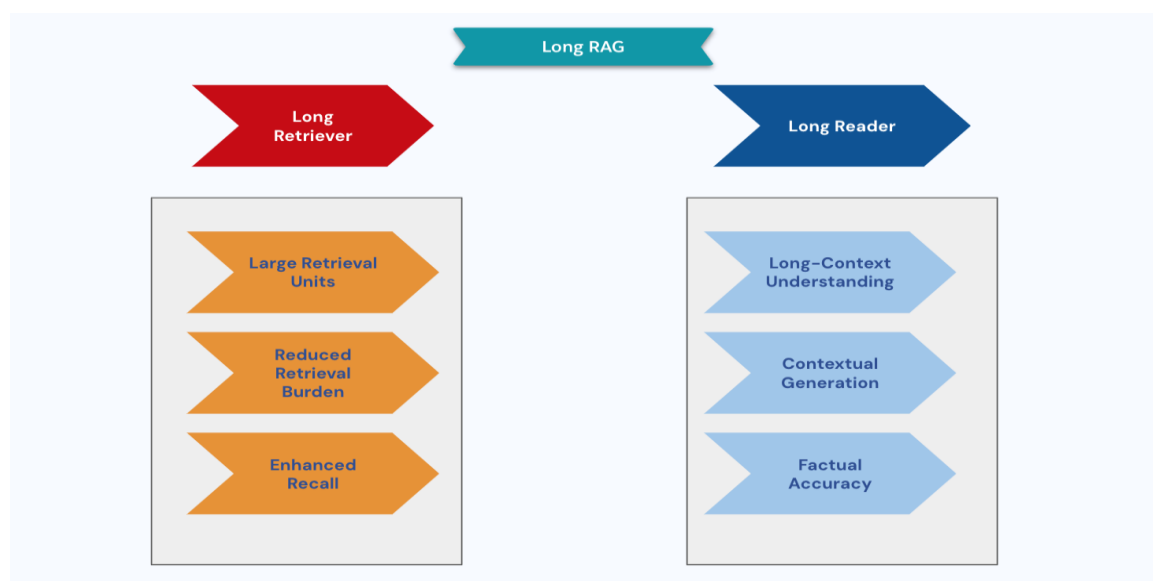
- Requirements for Knowledge Graphs: Designing and maintaining accurate, comprehensive and accurate knowledge graphs is an arduous and resource-intensive process.
- Complexity and Ambiguity of queries: Graph RAG may have trouble handling ambiguous or complex queries that involve extensive thought and exploration of the knowledge graph.
- Scalability: When knowledge graphs get bigger and more detailed, navigating and reasoning through them is computationally expensive, which constrains scalability.
- Factuality and Hallucination: Though Graph RAG could increase factuality by exploiting the power of structured knowledge, it may also generate false or misleading information, if you're dealing with partial or contradictory data.

**3. LongRAG**

LongRAG is a retrieval-augmented generation (RAG) framework that employs LLMs to create text based on data returned from a knowledge reservoir. It's built to tackle longer context questions and tasks that are challenging for conventional RAG models. LongRAG was able to perform a wide range of long-context tasks such as question answering, summarizing and translation. Also, it has been demonstrated to be more accurate than standard RAG models for such functions. LongRAG has two main components: a long retriever and a long reader. The long retriever is extracting relevant information from a knowledge base, and the long reader is reading the information returned. The long retriever can be used to access data from almost any type of source, such as text, code and data. The long reader is built to produce text that's correct and consistent despite the partial or incongruous information it gets.

Summary -

- Strengths: Takes care of bigger contexts, saves a lot of memory.
- Weaknesses: Risk of overloading, slower performance for longer documents.
- Use Case: Constructing summaries of lengthy legal or scientific documents.



## Long Retriever

- **Large Retrieval Units:** LongRAG splits the whole Wikipedia into 4K-token units, which is far larger than RAG designs. It enables the access of richer and context-specific data.
- **Reduced Retrieval Cost:** By increasing retrieval unit size, LongRAG reduces the amount of units that the retriever has to retrieve. This significantly decreases the number of CPUs and makes the retrieval more efficient.
- **High Recall:** The long retriever prioritizes recall, retrieving as much relevant context as possible. This puts the weight of retrieving onto the reader, who is better able to extract precise information from long sequences.

## Long Reader

- **Long-Context Completion:** Long-context LLMs can be programmed to parse long strings of text. This allows the reader to use all of the data that has been extracted from the long retriever.
- **Contextual Generation:** The reader creates the content based on the query input and the long context returned. This makes them possible to produce better and more consistent answers, even if the information returned is partial or inconsistent.
- **Factual Accuracy:** The dual information view from LongRAG helps to attain factual accuracy by accounting for the retrieval context as well as the original query. This is especially relevant to high-level questions involving multiple jumps of reasoning.

## Benefits:

- **Durability:** LongRAG is robust against noisy and incomplete data, making it ideal for use cases.
- **Scalability:** LongRAG easily scales to accommodate larger data and more advanced queries.
- **Improved Retrieval:** The long retriever is faster and better at finding the right information, increasing recall scores.
- **Superior Writing:** The long reader can write much more precise and coherent text using the full data of the long retriever.

## Limitations:

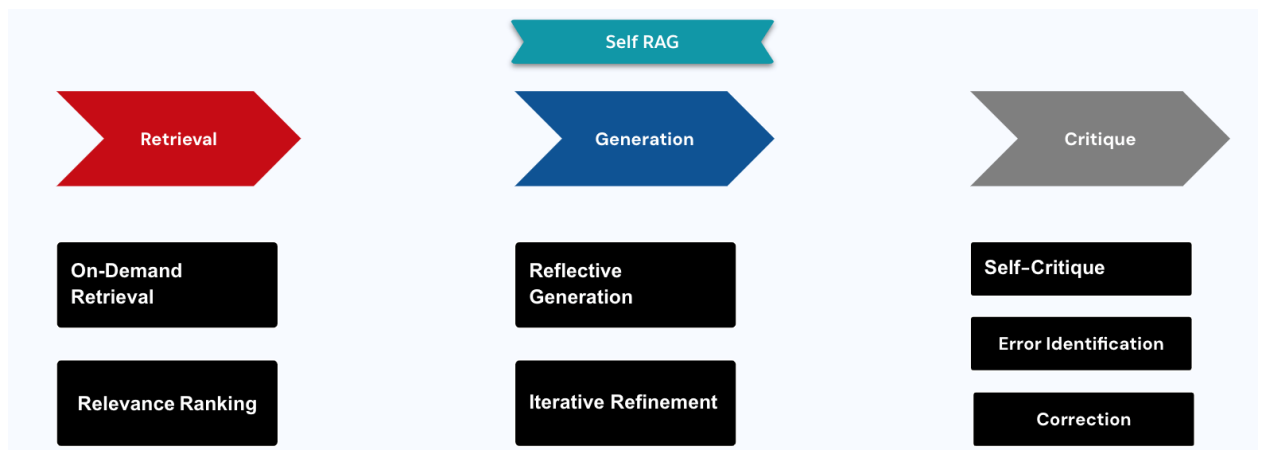
- **Limitations of the Context Window:** LongRAG does have some limitations despite the bigger context window. Long documents or queries can still overflow the model's information capacity.
- **Archiving** — Long, contextual sentences can be harder to locate than shorter ones.
- **Factuality and Hallucination:** While LongRAG alleviates some of the hallucinations of regular RAG, it can still produce inaccurate or deceptive data if you provide unclear or contradictory data.
- **Model Complexity and Cost:** LongRAG models demand a high computing overhead to train and infer, so they're hard to use for small teams or people with a limited budget.

#### 4. Self-RAG

Self-RAG is a new framework which teaches an LM how to retrieve, generate and critique its own generations, giving it more truthfulness and quality but less range. Self-RAG – In contrast to the traditional RAG method, self-RAG pulls information on demand, enabling adaptive retrieval algorithms. Moreover, Self-RAG will self-examine its own generation from multiple fine-grained dimensions by predicting reflection tokens as part of generation. This allows the model to detect any possible errors or mismatches in its output.

Summary -

- Strengths: Very little supervision, continuous improvement via self-teaching.
- Weaknesses: Unpredictable results, complex training architecture.
- Use Case: Product descriptions for large sets of unseen product data.



#### Retrieval

- On-Demand Retrieval: Unlike other RAG frameworks that use precompiled knowledge bases, Self-RAG retrieves data on demand. It enables adaptive and less rigid retrieval methods in relation to the task and environment.
- Relevance Ranking: The model automatically learns to prioritize information retrieved from the source as far as the task is relevant, so the most useful information is used.

#### Generation

- Reflective Generation: In generation, the model anticipates reflection tokens to determine the output. Such tokens could be used to spot possible inconsistencies, mistakes, or biases in the text being generated.
- Iterative Refinement: The model can iteratively optimize its outputs on the basis of its predicted reflection tokens to create a more accurate and logical writing.

#### Critique

- Self-Critique: The model learns to judge its own generations on several micro-levels, including accuracy, consistency, and relevance.
- Identification of Errors: By examining the reflection tokens, the model is able to detect errors or inconsistencies in its results.



- **Correction:** Now the model can correct these errors or deviations by getting some new information or changing its generation.

**Benefits:**

- **Better Factuality and Quality:** Self-RAG can generate more reliable, informative outputs if it is trained to extract the data it needs and self-evaluate its generations.
- **Greater Flexibility:** Self-RAG is flexible to a wide variety of tasks and interests because it can be programmed to create reflection tokens that configure its behaviour.
- **Greater Efficiency:** Self-RAG's adaptive retrieval mechanism fetches information only when required, thus minimizing computing effort.

**Limitations:**

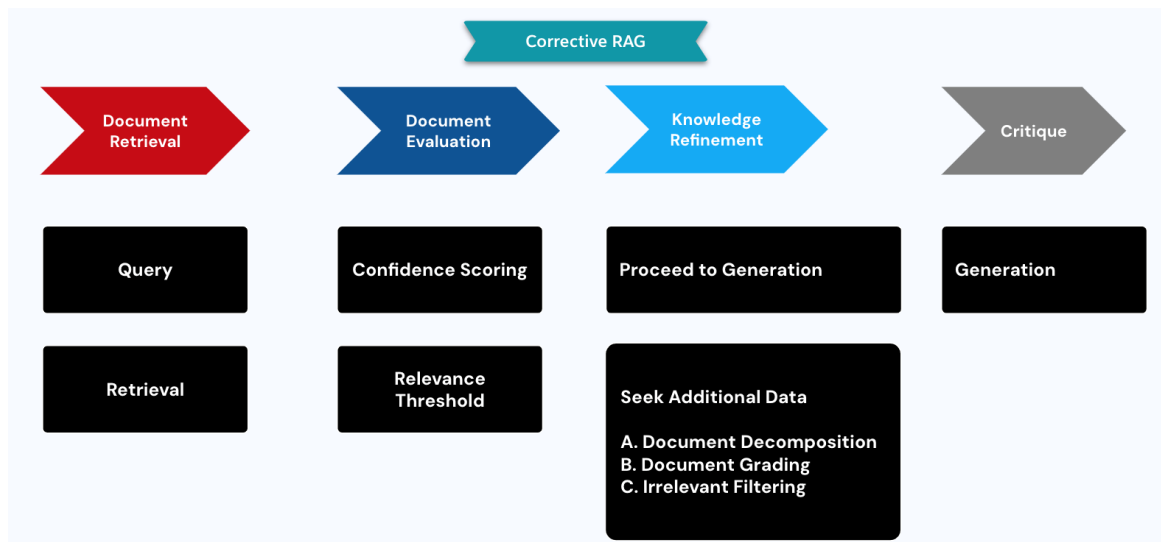
- **Self-Evaluation:** The model relies on self-evaluation to evaluate its own performance and learn where it can improve. This is somewhat limiting, especially when the job is particularly complex or sophisticated. The model might not always pick up on its own weaknesses.
- **Potential Overfitting:** Too much self-regard or too much dependence on a particular feedback loop can cause overfitting where the model gets too specific to a particular data or task. This can slow its generalisation to new contexts.
- **Computational Cost:** Self-RAG is computationally expensive, due to its generation, comparison, and finalization iterations. This might impact its scalability for large scale applications.

### 5. Corrective RAG

Corrective RAG (CRAG) is a method for increasing the robustness of RAG models. It achieves this by incorporating a self-check feature that lets the model check the validity of the extracted documents and correct them if necessary.

**Summary -**

- **Strengths:** Accurate (with correction) and self-teaching.
- **Weaknesses:** More complex model, performance impact.
- **Use Case:** Healthcare diagnostic tools with a strict focus on precision.



### Document Retrieval

- Prompt/Query: The requestor sends a query.
- Retrieve: A standard RAG system fetches relevant documents from a knowledge base.

### Document Evaluation

- Confidence Scoring: A lightweight assessor computes a confidence score for each returned document.
- Relevance threshold: A parameter is defined that specifies whether a document has sufficient relevance to generate.

### Knowledge Refinement

- Move to Generation: If at least one document crosses the relevance threshold, then the model proceeds to generation.
- Ask for Extra Data: If all documents do not meet or the grader does not know what he is looking for, the system requests extra data, such as searching on the web, to supplement the search.
  - Document Separation: The returned documents are cut up into "strips of knowledge."
  - Document Grading: Each knowledge strip is rated in terms of its relevance and veracity.
  - Null Filtering: Relevant or low-quality knowledge strips are eliminated.

### Generation

- Text Generation: The model produces text from the enhanced information, perhaps adding insights from the evaluation.

### Benefits:

- **Better Precision:** CRAG prevents errors and hallucinations by enabling the model to be based on current and correct information.
- **Greater Reliability:** Through a narrowing of the knowledge, CRAG increases the reliability of the output text.
- **Greater Flexibility:** CRAG can take different situational decisions by gathering more information in emergencies.

### Limitations:

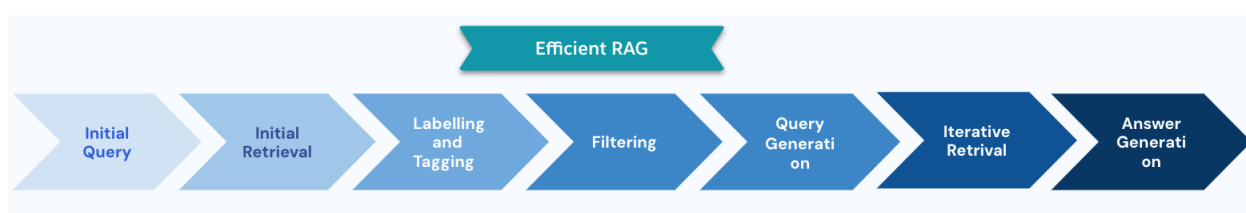
- **Dependence on Human Feedback:** The level of the corrections from human experts has a huge impact on how the model learns. Negative or inaccurate feedback can cause the model to entrain on incorrect patterns.
- **Poor Context Processing:** Corrective RAG can be unable to discern the larger context of a question or the subtleties of language, especially with difficult or opaque prompts.
- **Potential Overfitting:** The model trained on a small number of corrections can overfit to certain patterns and fail to generalize.
- **Computational Cost:** The process of constantly generating, evaluating, and cleaning up the responses can be computationally costly, especially on large systems.

### 6. Efficient RAG

EfficientRAG is a multi-hop query retrieval method that repeatedly creates new queries without large language models. It is a plug and play model which automatically searches for useful information through many fetch rounds to increase relevance in the information and eliminate unnecessary information, then optimize the quality and accuracy of answers.

#### Summary -

- **Strengths:** Built for speed and efficiency.
- **Weaknesses:** There is some sacrifice in precision and relevancy, not much slack for challenging tasks.
- **Use Case:** Real-time chatbots with high response rates.



#### Initial query

- The client sends a request.

#### Initial Retrieval

- The system extracts the relevant documents from a knowledge base using a standard retrieval approach.

#### Labeling and tagging

- The "Labeler & Tagger" feature goes through the extracted documents and locates relevant information and objects.
- It assigns labels and tags to pertinent sections of the text to draw attention to key data.

#### Filtering

- "Filter" element filters out the junk from returned documents.
- It is focused on the most relevant bits of the text that should be of use to the query.

#### Query generation

- Based on the filtered data and entities found, the engine produces new, optimised queries.
- Such novel questions try to probe more specific details of the original question.

#### Iterative retrieval

- The same goes for the new queries that are then executed to get more relevant data.
- Through this iterative strategy, the system learns gradually what is known about the query and finds the most relevant information.

#### Answer Generation

- Once the system has collected enough data, it can compute a complete, correct response to the original query.
- It can be as simple as summarizing the results or separating out facts or giving a detailed description.

#### **Benefits:**

- **Scalability:** Reduces large language models (LLMs) requirement for every iteration and is cost-efficient and scalable.
- **Improved Precision:** By selecting the right data and continually fine-tuning the query, EfficientRAG will come up with better, more detailed answers.
- **Multi-hop Reasoning:** Processes advanced questions that involve multiple steps of reasoning.
- **Plug and Play:** Easy integration with current RAG systems.

#### **Limitations:**

- **Retrieval Efficiency:** Indexing and searching large knowledge resources are computationally expensive, especially for real-time applications.
- **Model Complexity:** Large language models (and specifically the ones employed in Efficient RAG) require massive computing resources to train and infer.

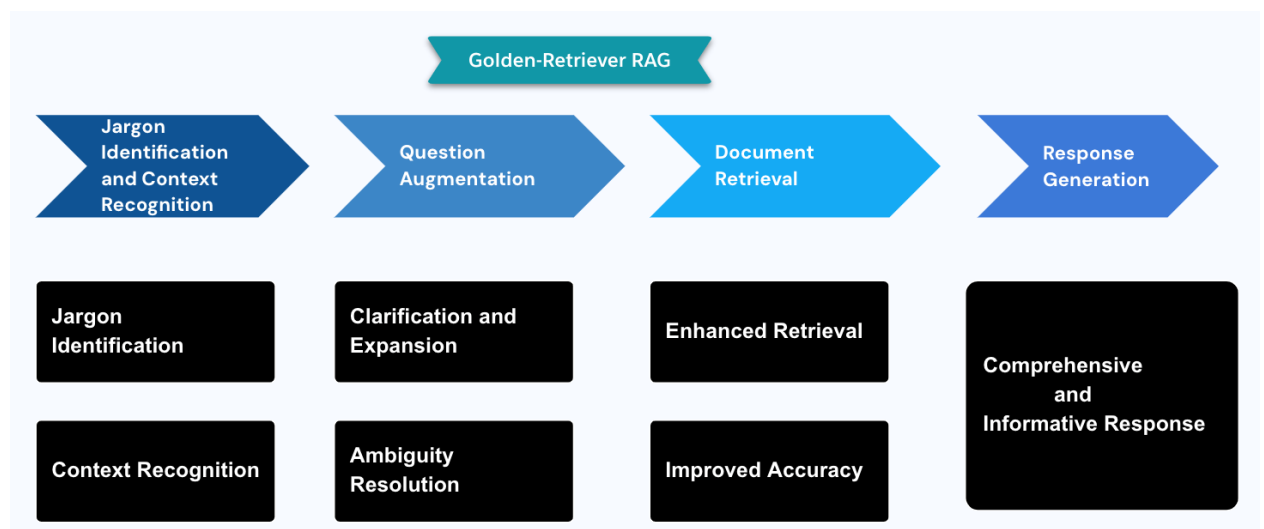
- **Data Quality and Bias:** Availability and quality of the knowledge base has huge implications for Efficient RAG performance. Errors in the data can show up in the output responses.

## 7. Golden-Retriever RAG

Golden-Retriever is an engine built to search and index vast industrial databases efficiently. It tackles the issues of old LLM fine-tuning and RAG approaches in the context of domain-specific jargon and context-decoding.

Summary -

- **Strengths:** High-accuracy retrieval, reliable results.
- **Weaknesses:** High accuracy required for narrow use cases, requiring significant investment.
- **Use Case:** Fact-checking apps requiring highly accurate information retrieval.



### Jargon Identification and Context Recognition

- **Identification of Jargon:** Golden-Retriever interprets the query input from the user in order to discover any industry-specific terms unfamiliar to a general-purpose language model.
- **Context Recognition:** it uses the existing knowledge base to detect context when these words are spoken thereby providing correct meaning.

### Question Augmentation

- **Clarification and Expansion:** The captured jargon and context are applied to supplement the original question to make it more precise and applicable to the knowledge base.
- **Ambiguity Reduction:** It explains the ambiguities and incorporates context to enhance the accuracy of the query.

### Document Retrieval:

- **Enhanced Retrieval:** The augmented query fetches the most relevant documents from the knowledge base.
- **Better Retrieval Performance:** Through enabling context and removing uncertainties, Golden-Retriever brings document retrieval accuracy to a whole new level.

#### Response Generation

- **Exhaustive and Informing Response:** The returned documents are interpreted to produce an explanatory and informative answer to the user query.

#### **Benefits:**

- **Improved Accuracy:** Overcomes the issue of language and context-specificity in the domain, providing more accurate and useful data.
- **High Efficiency:** Speedy retrieval and generation processes allows for faster and responsive use of the knowledge base.
- **Scalability:** Can be leveraged for massive industrial knowledge sources.

#### **Limitations:**

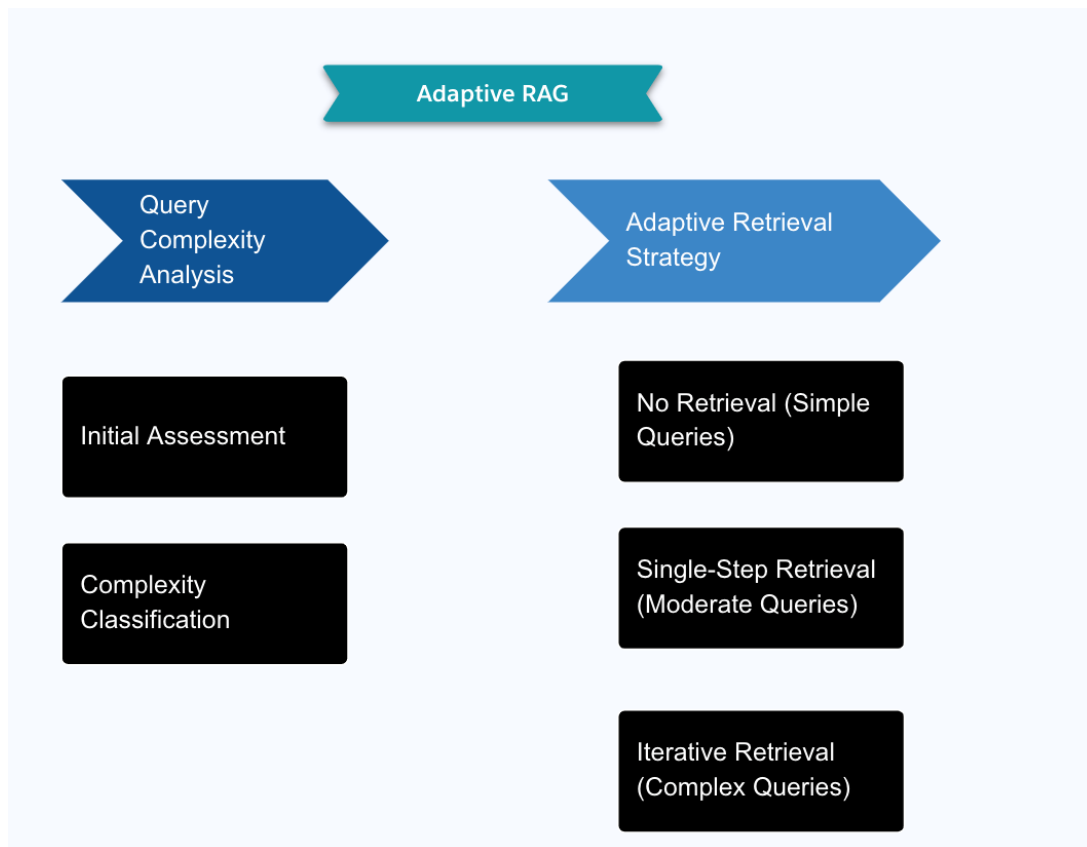
- **Jargon and Domain-Specific Terms:** Golden-Retriever might not be comfortable with specialized terminology and acronyms (especially in relation to proprietary documents). This can result in misrecognition and potentially false responses.
- **Ambiguous Queries:** If the query is unclear or contextless, Golden-Retriever may be unable to pick out the best documents. This might lead to redundant or partial responses.
- **Computational Cost:** Expanding queries and finding the correct documents is computationally costly for large knowledge sets. This can restrict Golden-Retriever scalability.

#### 8. Adaptive RAG

Adaptive RAG is a technique that can improve the speed and accuracy of RAG by automatically tuning the retrieval method depending on the complexity of the query.

Summary -

- **Strengths:** Task-based tuning, cross-domain usability.
- **Weaknesses:** Detailed tuning, risk of overfitting.
- **Use Case:** Individualised recommendation engines based on preferences.



### Query Complexity Analysis

- **Initial Assessment:** The system first assesses the query to understand the level of complexity. This is achieved by taking into account parameters such as:
  - Query length
  - Keyword density
  - Semantic complexity
  - Specific question type
- **Complexity Classification:** The query is ranked according to its complexity into one of three classes:
  - **Simple:** Simple questions which you can resolve directly and without the external information.
  - **Moderate:** Questions that require a little bit more context to get an answer right.
  - **Hard:** Searches that require extensive background knowledge about the topic, and may involve multiple rounds of information-gathering.

### Adaptive Retrieval Strategy

After finding out the query complexity, the system chooses the optimal retrieval approach:

- **No Search (For Simple Queries):** For simple queries, the LLM produces an answer without extracting additional information. This way, it is fast and does not incur computational cost.

- **Single-Step Retrieval (Moderate Queries):** For moderately complicated queries, the system retrieves a single set of relevant documents from the knowledge base and feeds them into the LLM. The LLM queries the retrieved data and returns a response.
- **Iterative Retrieval (Highly Complex Queries):** Iterative approach is taken when you are retrieving a highly complex query. It starts by initially retrieving some of these documents, parses them, and refines the query from those initial outcomes. This more specialized query is used to search a second set of documents, and so on until the LLM comes up with a fully-formed and accurate answer.

**Benefits:**

- **More Efficiency:** Without wasting time retrieving irrelevant data for simple queries, Adaptive RAG saves considerable amount of computing and time.
- **Improved Accuracy:** By customizing the retrieval strategy to the query complexity, Adaptive RAG increases the accuracy and relevance of retrieved information resulting in accurate and informative answers.
- **Scalability:** Adaptive RAG can be used across multiple query complexities and domains making it highly scalable.

**Limitations:**

- **Quality Retrieval Is Critical:** Adaptive RAG works only when the documents are of high quality. If the retrieval system doesn't catch the right data, then the model won't work properly.
- **Model Bias & Hallucination:** Despite adaptive retrieval, there are biases within the language model due to the training data. They also might produce inaccurate or misleading data, particularly when addressing complicated or uncertain questions.
- **Cost Per Unit:** Query refinement and retrieval, which happens repeatedly, are computationally costly in large applications. This can limit Adaptive RAG scaling.

**9. Modular RAG**

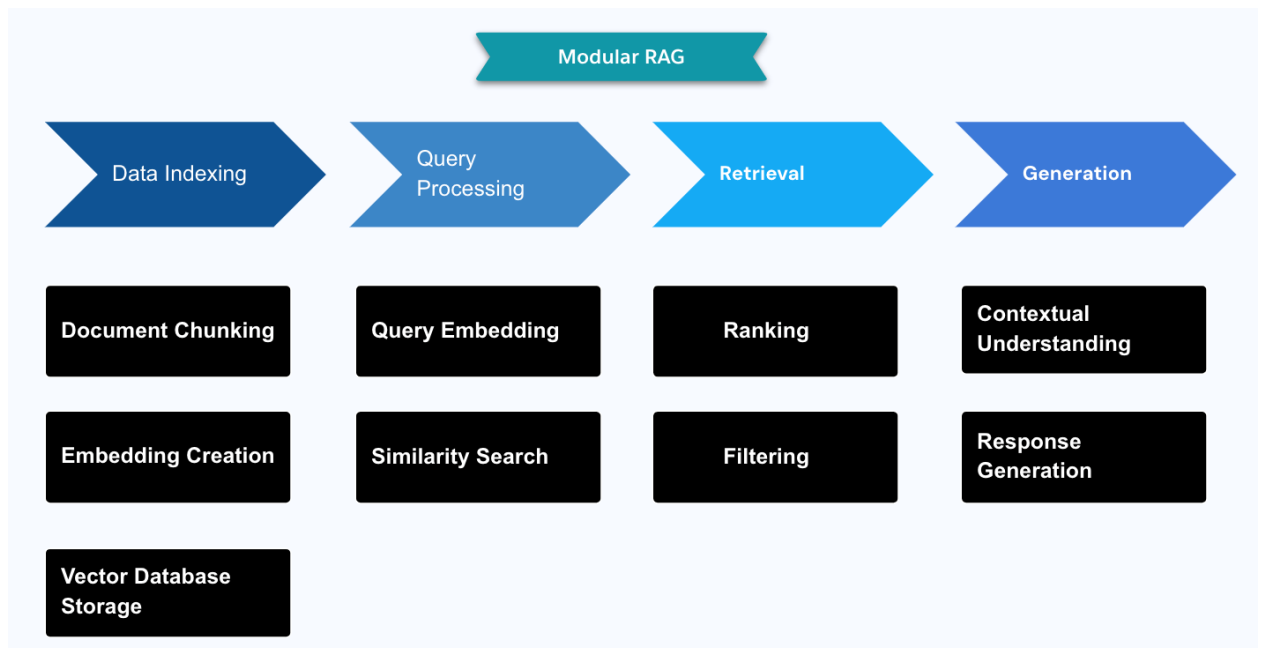
Modular RAG – Modular RAG is an Advanced form of Retrieval Augmented Generation that reduces the original RAG into smaller independent modules.

This modular architecture gives flexibility, scalability, and performance.

Summary -

- **Strengths:** Customizable components, flexible architecture.
- **Cons:** Integration problems, more frequent maintenance.
- **Use Case:** Building elaborate talk agents with interoperability modules.





### Data Indexing

- **Document Chunking:** The big documents are shredded into manageable chunks.
- **Creation by Embedding:** Every chunk is merged to a number (embedding) with the help of methods such as BERT or other language models.
- **Vector Database:** These embeddings are stored in vector database and thus are used for similarity search.

### Query Processing

- **Query Embed:** Also the user query is embedded.
- **Similarity Search:** Query embeddings are compared to the vector database embeddings to determine most relevant documents.

### Retrieval

- **Priority:** The documents returned are ranked by the relevance of the query.
- **Filtering:** Disproportionate or redundant documents can be eliminated.

### Generation

- **Contextual Understanding:** The language model processes the retrieved documents to Contextual Intelligence: The language model analysed the retrieved documents for context information.
- **Response Generation:** The model returns a response in the form of the query and relevant documents obtained from the database.

### Benefits:

- **Flexibility:** Modules can be combined and modified according to requirements.
- **Scalability:** You can simply add or remove modules and scale up or down the system.

- Speed: specialized modules can be tuned for a specific task, thereby increasing speed.
- Stability: The system can scale to new data sources and new types of queries.
- Improved Automation: By decomposing the process into modules, the RAG pipeline is easily manipulated by developers at a much finer scale.

**Limitations:**

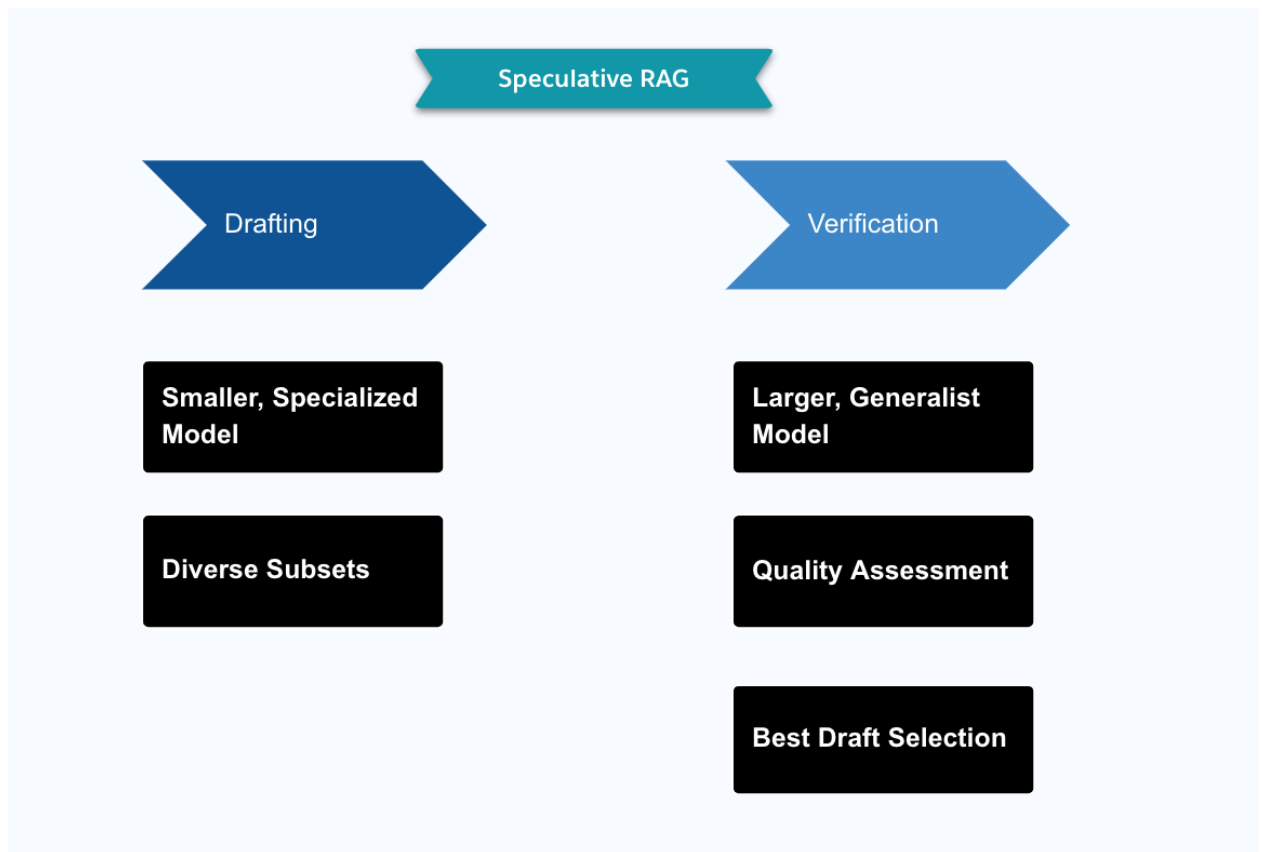
- Context Window Limitations: The maximum number of text that can be processed per RAG module is 256. This can be a limitation with deep or long-form queries (e.g. for legal or medical use cases).
- Hallucination Risk: Even with the best retrieval and generation algorithms, RAG models still generate reasonable yet incorrect information if supplied with poor quality or false data.
- Data Freshness and Quality: If the knowledge base isn't continuously updated, the model might display old or misguided data. Loud and biased data can be damaging to the model performance resulting in false positives.
- Complex Query Handling: Modular RAG might not be able to handle unclear/complex queries, especially when there is no indication of intent. There are a lot of queries which would involve multiple steps of reasoning and information search, especially when it's not explicit on what these concepts are referring to.

## 10. Speculative RAG

Speculative RAG is a revolutionary RAG that aims to increase both the efficiency and the accuracy.

Summary -

- Pros: proactive fetch for performance advantage, new method.
- Fails: The possibility of non-relevant fetches, a bigger model.
- Use Case: Text or email completion prediction based on user requirement.



It splits generation into two phases:

#### Drafting

- **Small, Dedicated Model:** A small, specialized language model, the "RAG Drafter" is used to create multiple draft answers.
- **Different Subsets:** This model is applied to different subsets of the returned documents. These subsets are chosen to give you multiple answers to the question, avoiding repetition and expanding the possibilities.

#### Verification

- **Bigger, Generalist Model:** The drafts that the RAG Drafter creates is sent to a bigger, generalist language model called "RAG Verifier".
- **Quality Assessment:** This larger model grades each draft in terms of quality, consistency, and application to the original question.
- **Best Draft:** RAG Verifier selects the best draft from the list of generated drafts.

#### Benefits:

- **Better Quality:** With the multiple drafts from multiple perspectives, Speculative RAG can come up with more accurate and complete responses.
- **Optimisation of Efficiency:** Drafting is done with a smaller, custom model that saves computational resources and accelerates the process of generation.

- **High Latency:** Multiple drafts are generated and checked in parallel, which minimizes the latency, and therefore translates into higher response times.
- **Resource Efficiency:** By assigning the drafting to a reduced model, Speculative RAG saves precious resources and is a cheaper alternative.

### Limitations:

- **Drafter Model Limitations:** The first drafts produced by the mini, superefficient drafter can be low-quality or high-quality. The verifier model won't do a very good job if the drafter's drafts are low quality.
- **Verifier Model Limitations:** Verifier model (which is generally bigger and stronger than the language model) can be a computationally expensive model when working with many complex queries and large sets of knowledge bases.
- **Training Data Requirements:** The drafter and verifier models need high-quality data, that is diverse, representative and not biased.
- **Ethical Considerations:** Misuse of Speculative RAG to create false and malicious content is an issue. If training data is biased, models may be biased or discriminatory in their output.

### 11. RankRAG

RankRAG is a new system that extends RAG models. It is about improving the accuracy and relevancy of the data on which language models compute answers.

Summary -

- **Pros:** More relevant with better ranking algorithms, more accuracy.
- **Cons:** Repetitive ranking algorithm construction, bias risk.
- **Use Case:** Search engines ranking the results with the highest relevance.



#### Retrieve

- The LLM pulls contexts from a knowledge base from the user query.
- This initial fetch can be done by keyword matching, semantic search, vector databases, etc.

#### Rerank

- RankRAG reranks returned contexts to find the best match for that query.

- This reranking process is necessary as it will let the model pick the most informative and relevant contexts.
- For each context, LLM scores relevance with the query.
- Those contexts are then ranked on relevance score, and only the best-ranked contexts are invited to the next step.

#### Generate

- The LLM pulls together the best-ranked contexts to return a detailed, informative response to the query from the user.
- The LLM draws on the knowledge of the contexts selected to come up with an answer that is correct and applicable to the question.

#### Benefits:

- Accuracy and Relevance: With the correct detection and usage of relevant data, RankRAG provides more accurate and informative answers.
- Greater Efficacy: The aggregate design of RankRAG removes the computational overhead for different context ranking and answer generation models.
- Increased Context Awareness: RankRAG is able to comprehend the context of a query in order to give you better contextually adapted answers.
- Better Generalization: RankRAG is a generalisable tool for various applications as it is applicable in new domains and new activities.

#### Limitations:

- Reliance on Initial Retrieval: If the initial retrieval step does not generate documents with the greatest relevancy, ranking and generation processes can also fail to generate meaningful and precise results.
- Computational Cost: Rank process, especially in the case of large language models, is computationally demanding and demands large hardware and software resources.
- Data Quality and Bias: RankRAG performance can be influenced in many ways by the quality of knowledge base and training data. Noisy or distorted data may cause incorrect and discriminatory answers.

## 12. Multi-Head RAG

Multi-Head Retrieval Augmented Generation (Multi-Head RAG) is a novel way of expanding the RAG model. RAG models increase the quality and relevance of answers by pulling appropriate documents from a knowledge base and integrating them into the generation process.

Summary -

- Advantages: Parallel processing to increase performance, various outputs.
- Negatives: CPU intensive, complex head-to-head coordination.
- Use Case: Answering questions with multiple viewpoints or data sources.



### Query Encoding

- The language model encoder takes the query from the user and wraps it into a high-dimensional vector.

### Multi-Head Attention

- The query embedding is pushed into a multi-head attention layer.
- Every attention head in turn responds to the query and then examines various aspects or details of the query.
- This enables the model to represent different viewpoints and return multiple query embeddings.

### Document Retrieval

- For each attention head:
  - It uses the generated query embedding to query the knowledge base.
  - A similarity score is assigned between the query embedding and the embeddings of documents in the knowledge base.
  - These highest-ranking documents are extracted based on their similarity values.
  - This operation is repeated for each attention head, so that you have multiple sets of downloaded documents.

### Contextual Understanding and Generation

- The reads from all attention heads are aggregated and injected into the decoder of the language model.
- The decoder then works through the combined context and returns a better-structured, more detailed and user-friendly answer to the user's question.

### Benefits

- **Better Retrieval Precision:** Because Multi-Head RAG takes into account multiple sides of the query, it returns more specific and diverse documents.
- **Improved Response Quality:** The rich collection of retrieved documents allows the model to produce more accurate and meaningful responses.

- **Better Processing Of More Complex Queries:** Multi-Head RAG can easily process more complex queries involving different information from different sources or different perspectives.
- **More Flexibility:** Modularity in attention heads lets you customize how the model behaves for different tasks and domains.

**Limitations:**

- **Complexity and Cost of Computation:** Since MRAG has multiple retrieval/generation heads, the model architecture can be very complex.
- **Data Quality and Bias:** Knowledge base and training data quality can affect MRAG performance tremendously. Noisy or biased data can lead to false and discriminatory conclusions.
- **Hallucinations and Factual Misrepresentations:** Even with multiple heads, MRAG might still generate plausible but inaccurate knowledge on complex or sophisticated subjects.
- **Sensitivity to Prompt Engineering:** The quality of the generated response depends heavily on prompt's specificity and explicitness. Poor prompts will lead to useless or inaccurate results.

**Summary and Conclusion**

Overall, the method Retrieval-Augmented Generation (RAG) is a robust technique to use LLMs to create answers from predefined knowledge. Combining information retrieval with LLM creation allows RAG systems to make outputs accurate, relevant and true. It's a very useful strategy in all sorts of use cases such as chatbots, customer care, personalized recommendations.