

* RAGIs (Retrieval Augmented Generation):
LLMs face significant limitations especially in domain-specific or knowledge-intensive tasks, notably producing "hallucinations".

- RAGIs overcome these challenges by retrieving relevant document chunks from external knowledge base through semantic similarity calculation.
- It is a key technology in advancing chatbots and enhancing the suitability of LLMs for real world applications.

- RAGI has three core parts:
 - Retrieval
 - Generation
 - Augmentation

- Based on development of RAGI, it is of 3 types: Naive RAGI, Advanced RAGI & Modular RAGI.

- LLMs are great for generating text, however that does not guarantee that the text is also factually accurate, thus RAGs come in.

* Naive RAG : (Retrieval Read)

→ INDEXING :

cleaning & extracting raw data (from HTML, PDF, Word, Markdown) → convert to plain text format → these chunks are then encoded to vector representations using embeddings & database
• crucial to enable similarity searches.

→ RETRIEVAL :

gets input user query → converts to vector rep. using same encoder of indexing → computes similarity score → prioritizes & retrieves top K chunks that demonstrate greatest similarity

→ GENERATION :

posed query & selected documents → synthesised into a coherent prompt to which LLM is tasked

* Difficulties / Disadvantages of Naive RAG:

- Retrieval : struggle with precision & recall
- Generation : may still face hallucination
- Augmentation : integration may result in disjointed or inconsistent outputs
- Moreover, concern that generation models heavily rely on retrieved info, thus not adding insightful or synthesised information.

- * Advanced RAGI: (Retrieve Read)
- Specific improvements to overcome the limitations of Naive RAGI.
 - Enhancing RETRIEVAL : optimise indexing
 - pre - retrieval (structure & log)
 - post - retrieval (crucial to integrate retrieval with query)
 - Enhancing INDEXING :
 - sliding window approach
 - fine-grained segmentation
 - incorporation of metadata

- * Modular RAGI: (Rewrite Retrieve Read)
 Many modifications that include =

- 1) New Modules :
 - Search Modules : enable direct search across various data sources like search engines, databases & knowledge graphs .
 - RAGI-fusion addresses search limitation by employing a multi query strategy that expands user query into diverse perspectives , utilising parallel vector searches and intelligent re-ranking
 - Memory Module : leverages LM's memory to guide retrieval , thus aligning it closely with data distribution through iterative self enhancement .

- Predict Module aims to reduce the redundancy & noise by generating content directly through the LLM.
- Task Adapter : tailors RAG to various downstream tasks, automating prompt retrieval for zero shot inputs
↳ tasks applied on a trained model .
- The distinguishing property of Modular RAG is that it allows module substitution or reconfiguration to address specific challenges .
 - leveraging LLM's capability of refined retrieval
 - RAG system can more easily integrate with other technologies (such as fine tuning or reinforcement learning).

* RACRs vs Fine Tuning vs Prompt Engg

like providing a ~~less~~ comparable to model with a student tailored textbook internalising for information knowledge overtime, retrieval good for scenarios involving replication.

leverages model's with min necessity for external knowledge.

- not mutually exclusive, can complement each other
- RACRs outperform FT on various knowledge intensive tasks across various topics .
- alone classification based off : external knowledge & model adoption requirements .

A * Retrieval

1. Retrieval Source

1. Data Structure:

- Text → Unstructured data (text) Wikipedia
- Table → Semi-structured data (pdf)
- Table → Structured data (knowledge graphs)
- Graph neural network (GNN), LMs & RAG integration.

2. Retrieval Granularity:

- Coarse Grained: can provide more relevant info for the problem but may also contain redundant content which could distract retriever & language models in downstream tasks
- Fine Grained: increases the burden of retrieval & does not guarantee semantic integrity & meeting req. knowledge

Token, Phrase, Sentence, Proposition, Chunk, Document

Fine to Coarse

Entity, Triplet, subgraph

Granularity of text & knowledge graph

B * INDEX/NC Optimisation

- determines whether correct context can be obtained.

1. Chunking Strategy

- leads to truncation of sentences, prompting the optimisation of a recursive splits & sliding window methods, then enabling layered retrieval by merging globally related info.
- still, not the perfect balance.
- long chunks → more data more noise
- small chunks → less data less noise.

2. Meta Data Attachments :

- enriching with metadata info such as page no, file name, author, category, timestamp, limiting scope of retrieval.
- can be artificially constructed, for eg: adding summaries of paragraph.

3. Structural Index:

- hierarchical structure for documents
 - ↳ parent-child relationships.
- knowledge graph Index:
 - Use KB where nodes (text) and edges (show relation of any sort among the text)
- This, maintains consistency, and describes the connection b/w different entities & concepts, markedly reducing potential illusions

C * Query Optimisation :

1. Query Expansion
2. Multi Query
3. Sub Query
4. Chain of Verification (CoVe)

enriches
content
of query
& enable to
fully
answer the asked qs.

2. Query Transformation:

- to extreme chunks based on transferred query over user query as they may not always be optimal.
- Thus, we can either prompt the LLM to rewrite or use another smaller model.

3. Query Routing:

- Metadata routing : extracting keyword & apply filter on metadata
- Semantic Router : leveraging semantic info of the query.
- A hybrid approach can also be employed.

D * EMBEDDING

- 1. Mix / Hybrid Retrieval
- 2. Fine - Tuning Embedding Model
 - (especially when used for specific domains such as health, legal, etc).
 - capture different relevance features and can benefit from each other by leveraging complementary relevance info.

E * Adapter:

Use pretrained models to help bridge the gap which could not be achieved by finetuning.

IV GENERATION

A* Content Curation

1. Re Ranking (Assigning weightage to retrieved)
2. Content Selection / compression (small Language Model (SLM) & LLM can be used in tandem to remove unnecessary tokens).

B* LLM Finetuning:

- allows the model to adjust to data & learn more.
- Finetuning of LLM & Retriever can also be coordinated.

V AUGMENTATION PROCESS

In order to present the best generated text some augmentation is done.

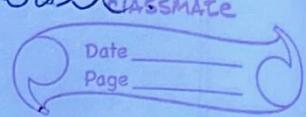
- Iterative Retrieval
- Recursive Retrieval
- Adaptive Retrieval

Iterative: enhance robustness of subsequent answer generation by offering additional contextual references through multiple iterations. After each round, retrieved info is used to refine search in next round.

Recursive: To improve depth & relevance of search results. Retrieves info, processes & then uses this output as new input, creating a layered/hierarchical process. (nest nested loops)

LVLM → Large Vision Language Model
(generate captions)

Whisper → audio to text



Adaptive :: adjusts its strategy based on the input or task. Dynamically changes retrieval approach, such as how it queries the database or which sources it uses.

* RAGs Use Cases:

- customer support & A chatbot
- Email chain analysis
- company internal documentation
- textbook & A

* Why RUN RAG locally?

→ Privacy, speed, cost

MULTIMODAL RAG

A multimodal RAG system should utilize both visual & language content in a video including:

1. Utilize the written text in the video.
2. Utilize pure visual info in the video.
3. Utilize the info being spoken about in the video.
4. Be able to hold multi-turn conversation w.r.t content of the video.

For similarity search: COSINE SIMILARITY, as it is computationally cheap & used widely

Multimodal

- ① Embeddings \Rightarrow Bridge Transformer.
- ② Video Processing \Rightarrow
 - ① Transformer (audio+image)
(no transcript all.)
 - ② Whisper model
(only image)
 - ③ LVL M.
- ③ Database for Embeddings:
 \rightarrow storing the multimodal embeddings using Lance DB (data can be private),
- ④ Retrieval \Rightarrow use similarity search
- ⑤ LVL M \Rightarrow take as input both vision & text
to generate textual outputs.