

# High-Performance RAG

Llamaindex  
edition!!

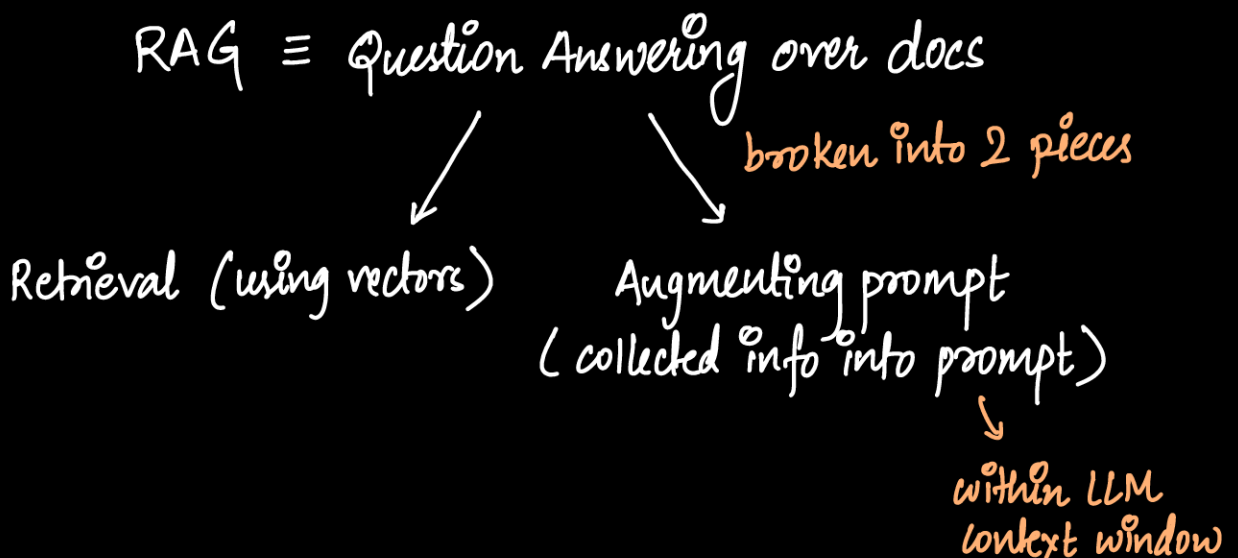
\* taken from AI Makerspace presentation by Dr Greg Loughrane  
& Chris Alexiuk

## WHY RAG?

Hallucinations lead to confident BUT FALSE responses.  
Better to have LLMs fact check against your documents.  
(especially in business context)


- Retrieval - find reference in docs
- Augmented - add reference to prompts
- Generation - improves answer to questions

Helpful in specialized domains → legal, health, finance,  
(with lots of jargons) gov, research, insurance




## RETRIEVAL (major emphasis in LlamaIndex)




Database used is usually VECTOR database  can also be called as INDEX

## BUILDING INDEX

- ① Split the docs into chunks
- ② Create embeddings for each chunk  using some embedding model
- ③ Store these embeddings in vector store index

## RETRIEVERS

- ① Ask a question
- ② Convert it into a vector  also called as retrieved context
- ③ Look for similar vectors in vector store index
- ④ Return retrieved context to put into prompt

## DENSE VECTOR RETRIEVAL



## RAG OVERVIEW

Query → Embedding Model → Similar vector  
in vector store

dense vector retrieval

Pass to LLM ← Set up prompt ← Rank order results  
(in context learning)

## LLAMA INDEX OVERVIEW

Tool to build VERY POWERFUL index and doing retrieval

“ A data framework for LLM apps to ingest, structure & access private or domain specific data ”

(Big companies → sizeable data → diff types of docs)

## LLAMA INDEX TERMINOLOGY

Sentence-level information → Nodes

PDF-level information → Documents

## NODE

(first class citizen in LlamaIndex)

- chunk of (source) doc
  - inherit metadata
  - can track where they came from
  - nodes are what get stored in vector store
  - Asking query returns "k"-most similar nodes
  - these nodes are passed through response synthesis
- ↓  
to create nodes

## NODE PARSERS

list of docs → chunks into node objects  
(different chunking strategies)

## QUERY ENGINE

(magic of LlamaIndex)

- Generic interface allowing Question Answering over your own data
- as IMPORTANT as "chain" in LangChain

## DOMAIN SPECIFIC EXAMPLE - Camelids family (veterinary domain)

- ① Data : research paper
- ② Index building : load data, chunk nodes, create embeddings, store nodes in index

## IMPROVING RETRIEVAL

Simple to Advanced RAG

check previous post

Table  
Stakes      Adv  
Retrieval      Fine-  
tuning



- ① Fine-tuning Embeddings

Reason → Specialized vocabulary in veterinary research context  
Handle mumbo-jumbo in docs

### Process

- ① Need training-validation-testing sets of (Question, Retrieved context) pairs
- ② Loss function → HF Sentence Transformers loss  
Takes positive pairs and automatically augments dataset with negative pairs (Question, wrong context) built-in
- ③ Embedding Model → bge-small-en  
cost efficient open source model from BAAI

\* CODE EXPLANATIONS WILL BE COVERED IN NEXT POST \*

## SIMPLE-TO-BIG RETRIEVAL

We also look around the area of context returned

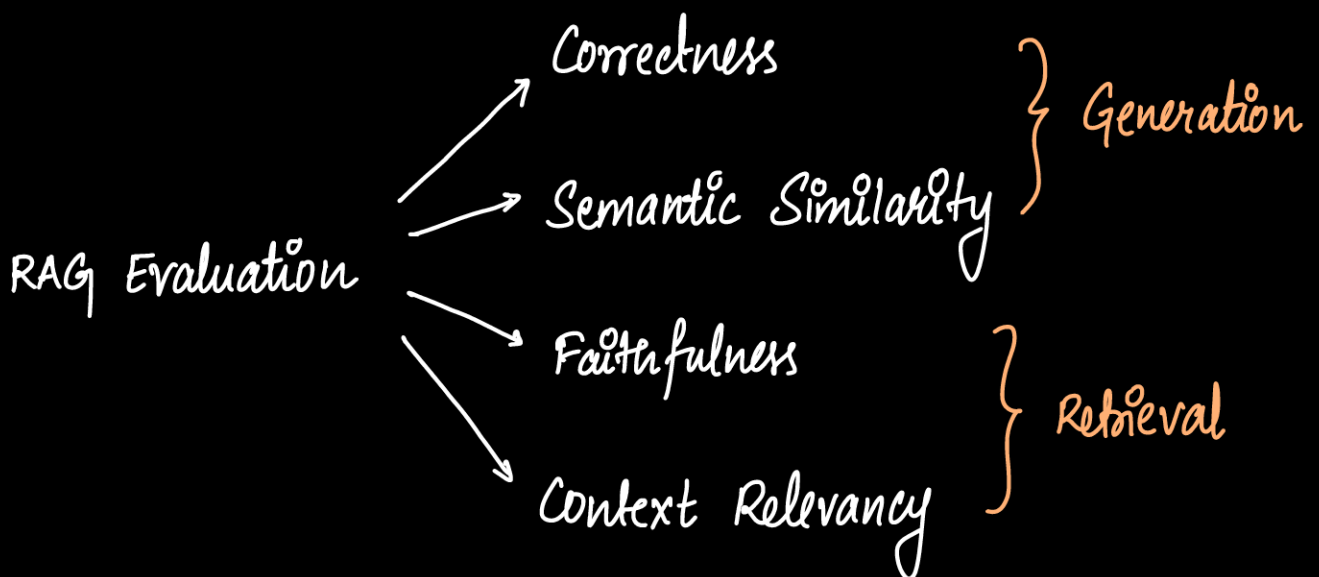
### Sentence Window Node Parser

Splits docs into nodes  $\rightarrow$  each node  $\equiv$  sentence

Window around each node  $\rightarrow$  sentences either side of node  
(look for relevant context)

## EVALUATING RETRIEVAL

Relevant context is MURKY TERM. Thus we utilize some libraries to evaluate returned context  
(built-in LlamaIndex)



## Generation Side

### ① Correctness

Comparing generated response to reference answer

created using  
GPT-4

### ② Semantic similarity

Generated response is SIMILAR to reference

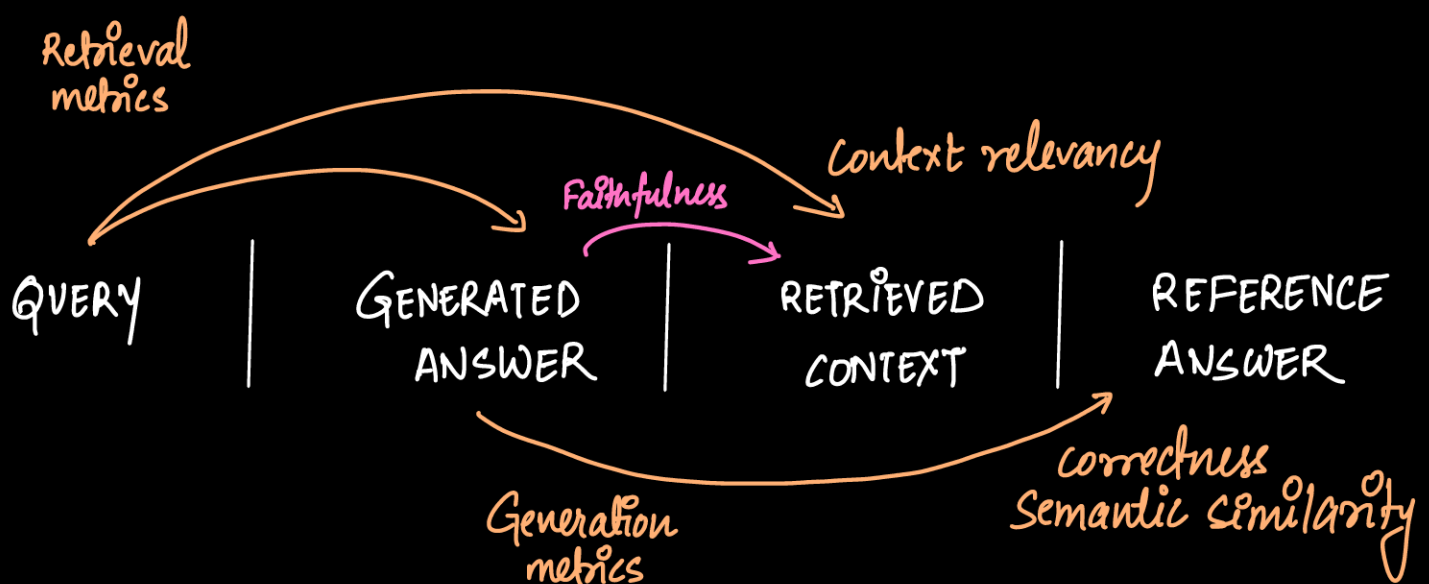
## Retrieval Side

### ③ Faithfulness

Are we seeing lots of hallucinations? Are our answers relevant to our context or not?

### ④ Context Relevancy

Are retrieved context AND answer relevant to query?





## NOTE

- When you analyze retrieval, then downstream during generation these generations also improve.
- No matter what metric you use, the results should be better compared to base RAG

## OBSERVATIONS

- Fine-tuning embedd ⊕ small-to-big retrieval improves RAG performance across all evaluation metrics
- Biggest gain are in Faithfulness & Relevancy  
massive decrease in hallucinations
- With more docs, RAG system seems to get mastery over domain and also its mumbo-jumbo
- Complex docs are next steps in more advanced retrieval methods

## CONCLUSION

- Many ways to enhance retrieval. Fine-tuning embeddings is really recommended for specialized vocabulary.
- Generation metrics also improve with Retrieval
- Evaluation of RAG systems is easier than ever  
GPT-4 still used as Ground Truth  
(it is really good at self assessment)