

Basics of RAG

Ben Claire (Answer.AI)

RAGatouille → Lib to use COLBERT easily

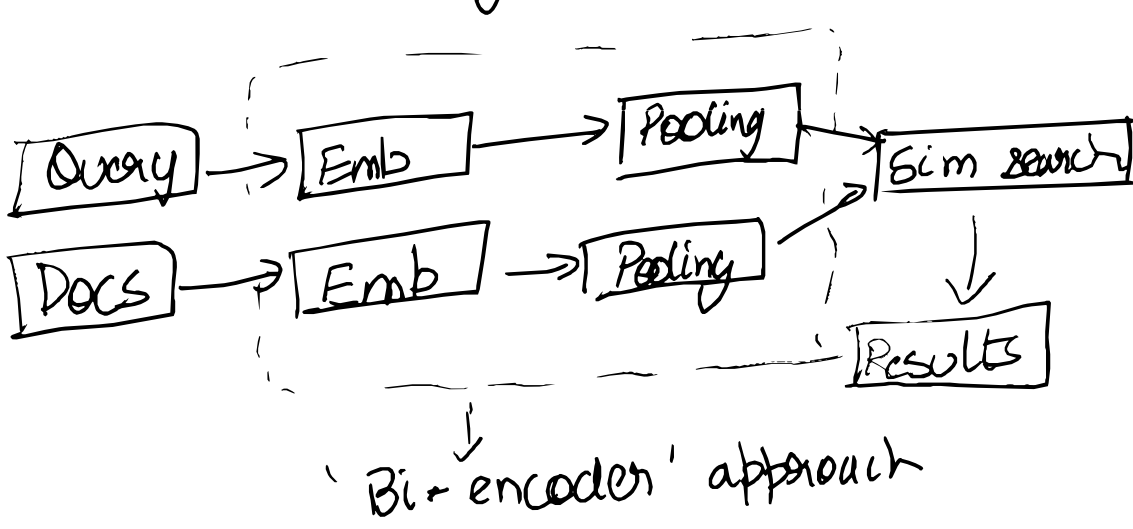
RemanReas → Lib for re-rankers

MVP → Bi-encoders single vec emb & cosine similarity are all you need.

Metadata filtering → Important

Good RAG:

- Good retrieval pipeline
- Good generative model
- Good way to link them up.



'Bi-encoder' approach

Why 'bi-encoders'

- Generally used to create 'single vector' representations. They precompute doc representations
- Docs & queries are computed separately. Not aware of each other
- Efficient but has some drawbacks i.e. retrieval perf tradeoffs

Re-ranking: The power of cross encoders

- 'Encode both query & docs at the same time
- Not computationally realistic to compute query aware document representations for every single document pair.
- It's effectively a binary classifier. The prob of positive class represents the relevance of doc to query, & is taken as the similarity score
- Re-ranking :- Use a expensive model to score only a subset of your retrieved documents
- Use Keyword Search
 - Embeddings represent info that is useful to their training queries
 - It'll never be fully representative of your queries
 - To capture signals use keyword search.
- BM25 is too strong a baseline
- Powerful when you have lot of domain-specific jargon.

Metadata Filtering

- Remove docs that are not relevant for query.
- Use Gliner to extract entities from your text & use in filtering