

# 20240530Meeting

312707003 黃鈺婷

# Retrieval-Augmented Generation(RAG)

- 檢索式模型+生成式模型
- 利用外部知識來增強生成的文本，使生成的文本更加準確
- 不需要重新訓練模型
- 適用於有大量資料，但多數資料未分類或標記
- 開放領域的問答任務、資訊檢索

# RAG vs. Fine-tuning

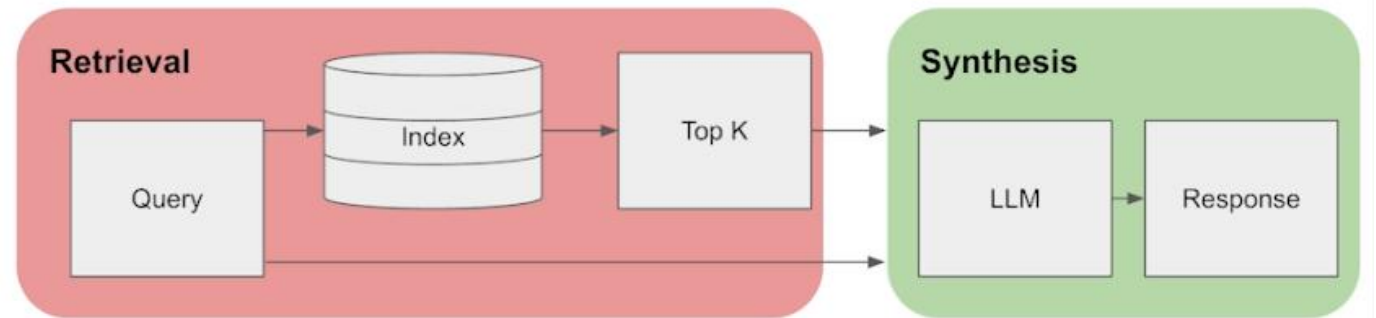
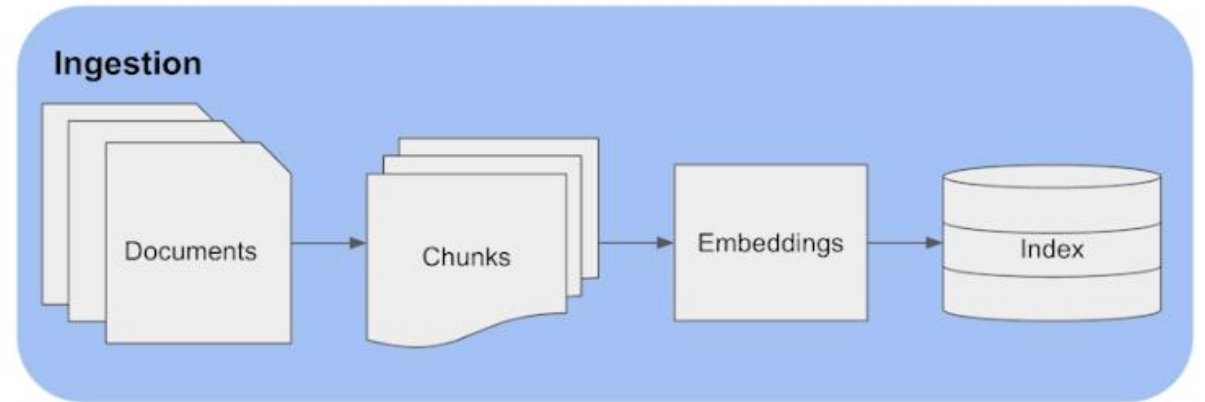
RAG	Fine tuning
直接更新檢索知識庫	需要重新訓練
需要少的數據處理	有限的資料集無法有顯著效能的提升
答案可以追溯到特定的資料來源	像黑盒子，無法知道模型為何以某種方式做出反應
無法自訂模型的行為或風格	允許根據特定語氣或術語調整模型的風格

# Neural Retrieval

- Encode the query and documents into dense vector representations
  - compute cosine similarity
  - determine which documents are most relevant to a query
- Advantage:
  - adept at dealing with long and complex queries
- Challenge:
  - performance depends on the data they are trained on

# Process of RAG

1. Vector Database Creation
2. User Input
3. Information Retrieval
4. Combining Data
5. Generating Text



# Embedding

- 高維離散的特徵映射到相對低維的連續向量空間中的表示方式
- 將原始數據轉換成一種特別的數據格式，以便 AI 或機器學習演算法能夠處理這些數據，並加入了距離的概念
- Sparse embedding
  - TF-IDF
    - lexical matching the prompt with the documents
- Semantic embedding
  - BERT
    - ◆ 擷取document和query中上下文的細微差別
  - SentenceBERT

# Sentence Embedding vs Token-Level Embedding

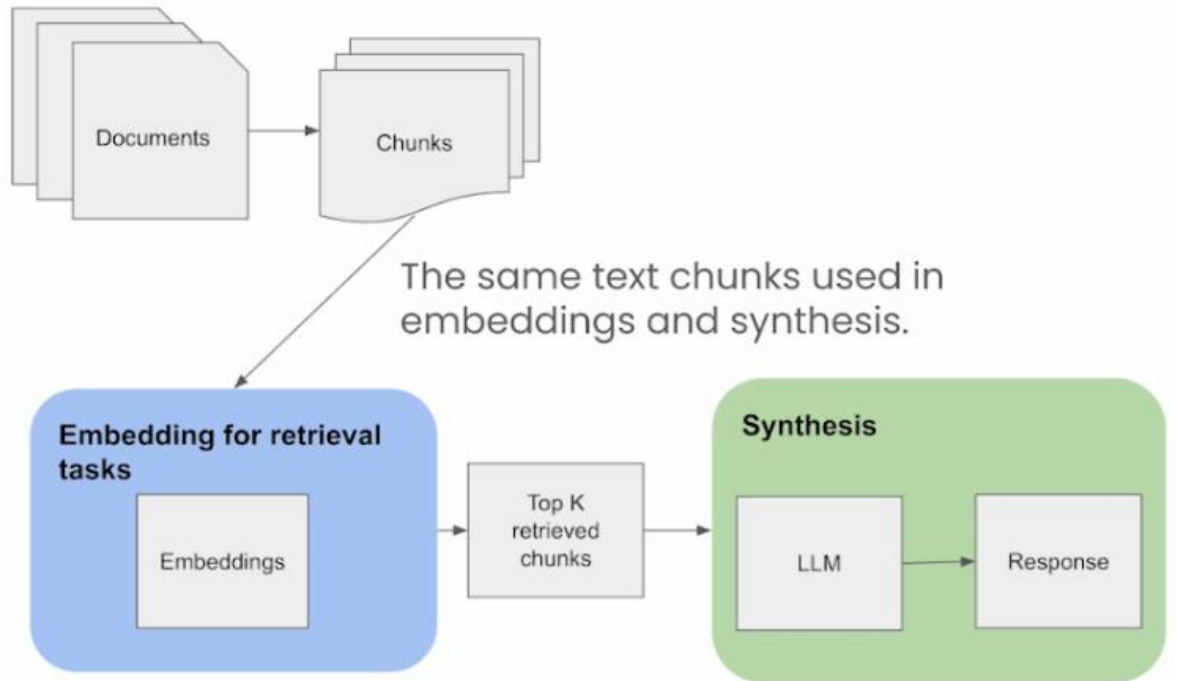
- Modification of the traditional BERT model
- Are trained specifically to understand the meaning of entire sentences
- Generate embeddings where sentences with similar meanings are close in the embedding space
- Provide a single embedding for the entire sentence
- Are more suited for tasks that rely on sentence-level understanding (like semantic search, sentence similarity)

# Retrieval

- Standard/Naive Approach
- Sentence-Window Retrieval
- Auto-merging Retriever

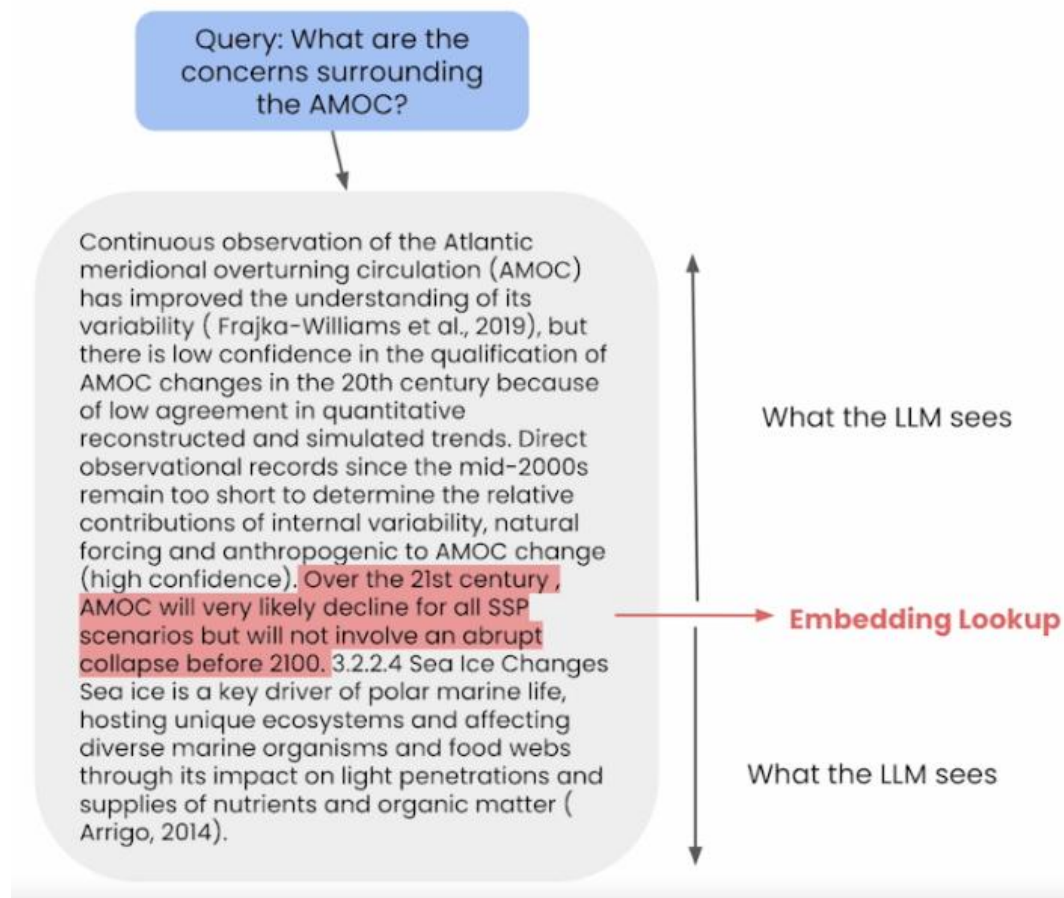
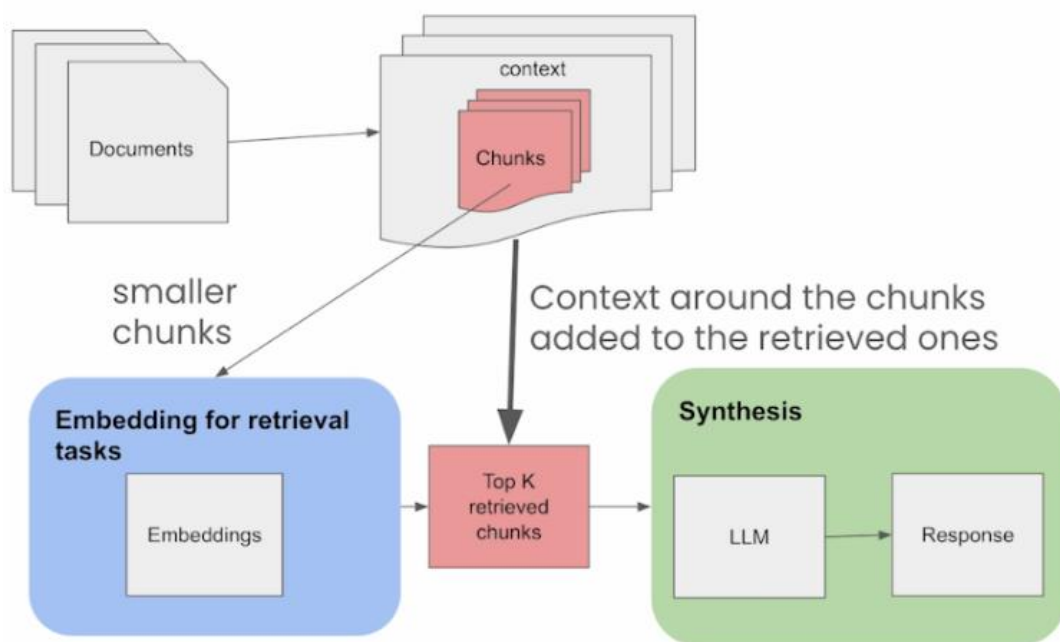
# Standard/Naive Approach

- Using the same text chunk for both embedding and synthesis, simplifying the retrieval process.
- Maintains consistency in the data used across both retrieval and synthesis phases



# Sentence-Window Retrieval

- Breaks down documents into smaller units, such as sentences or small groups of sentences



# Auto-merging Retriever

- Aims to combine information from multiple sources or segments of text to create a more comprehensive response to a query

