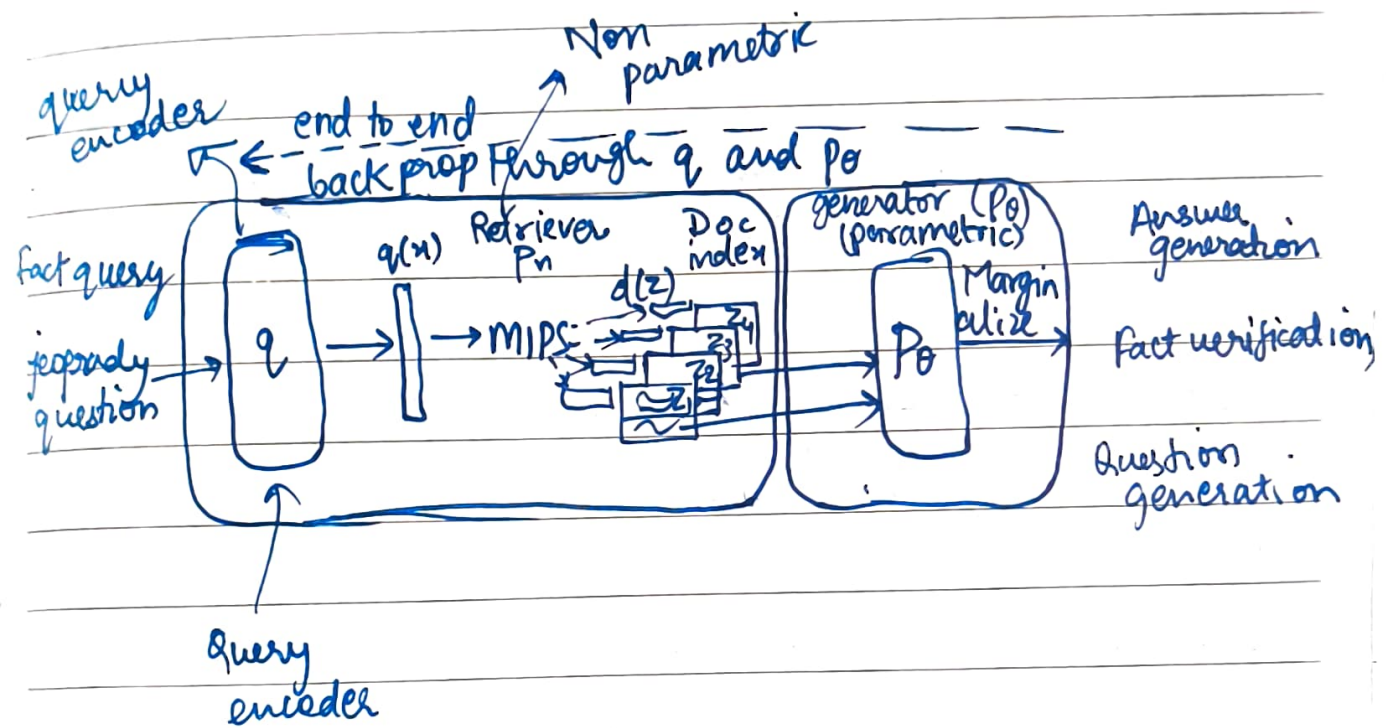


"Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks"



→ combines a parametric pre trained model & a non parametric retriever.

Architecture ☆

⇒ Input Query → denoted as 'x'. Can be a question, a fact to verify or an entity.

⇒ Query encoder → inputs 'x' → encoded in a dense vector representation generally ~~encoded in a~~ ^{based on} BERT

→ Retriever (non parametric)

↳ uses Maximum Inner Product Search (MIPS) to find top k matching documents, documents which have highest dot product with the query vector indicating relevance.

→ Here retriever uses a DPR (Dense Passage Retriever) and bi-encoder architecture (BERT)

→ Document Index → non parametric memory, storing embeddings of pre-processed & indexed documents.

→ Generator → (Parametric memory)

→ seq2seq model, takes 'x' & retrieved docs 'z'. pre trained BART model used here. Generates output response denoted as 'y'.

→ Output Generation → retrieved docs 'z' are treated as latent variable, and model marginalize over these docs to ~~create~~ produce the probability distribution for ~~the~~ generated text.

↳ the one not directly observed or explicitly chosen as a part of output but rather its possible values are considered to determine the final outcome.

paper gives two formulations of the RAG model

→ ~~RAG Sequence~~



RAG Sequence

→ same retrieved doc for prediction

→ marginalizes over top ~~K seq~~ retrieved docs to calculate prob of entire sequence

RAG-Token Model

→ different retrieved doc for prediction at each token generation step.

→ marginalizes at each token generation step. Allows generator to combine info from multiple docs when creating answer.

MIPS → method for finding top ~~K~~ ^{matching} most relevant docs.

→ finds highest dot product, in sub linear time

Parametric memory → stored within model's trainable

parameters. (A). Models that rely on parametric memory are limited by the knowledge in their training ^{data} up to its cutoff point. Lesser hallucinations

Non parametric memory → external knowledge, explicit, retrieval based memory. Knowledge can be ^{easily} updated (DPR)