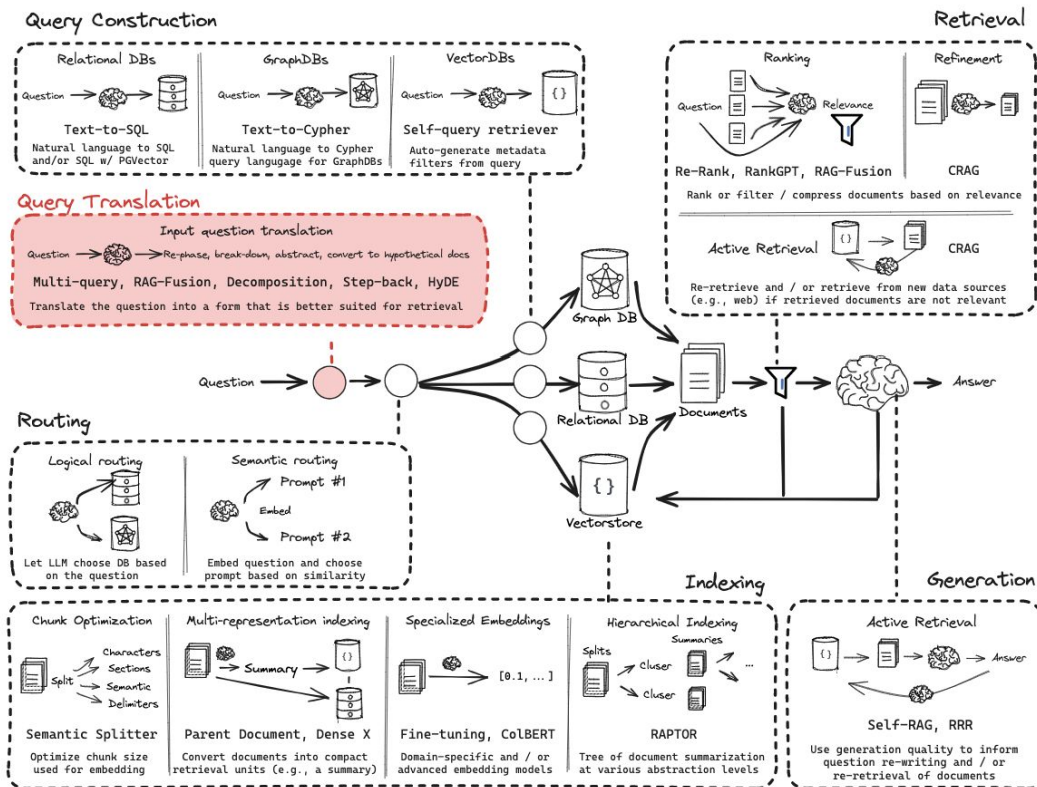# RAG from scratch: Query Translation (HyDE)

Lance Martin
Software Engineer, LangChain
@RLanceMartin

# Query Translation



## Query Construction

**Relational DBs**

Question → 🧠 → 🗄

**Text-to-SQL**

Natural language to SQL and/or SQL w/ PGVector

**GraphDBs**

Question → 🧠 → ⭐

**Text-to-Cypher**

Natural language to Cypher query language for GraphDBs

**VectorDBs**

Question → 🧠 → {}

**Self-query retriever**

Auto-generate metadata filters from query

## Query Translation

**Input question translation**

Question → 🧠 → Re-phase, break-down, abstract, convert to hypothetical docs

**Multi-query, RAG-Fusion, Decomposition, Step-back, HyDE**

Translate the question into a form that is better suited for retrieval

## Routing

**Logical routing**

🧠 → 🗄 → ⭐

Let LLM choose DB based on the question

**Semantic routing**

🧠 → Prompt #1
Embed
→ Prompt #2

Embed question and choose prompt based on similarity

## Retrieval

**Ranking**

Question → 🧠 → Relevance

**Refinement**

🧠 →

**Re-Rank, RankGPT, RAG-Fusion**          **CRAG**

Rank or filter / compress documents based on relevance

**Active Retrieval**  {} → 🧠   **CRAG**

Re-retrieve and / or retrieve from new data sources (e.g., web) if retrieved documents are not relevant

Question → ⚪ → ⚪ →
Graph DB
Relational DB  → Documents → 🧠 → Answer
Vectorstore

## Indexing

**Chunk Optimization**

Split → Characters / Sections / Semantic / Delimiters

**Semantic Splitter**

Optimize chunk size used for embedding

**Multi-representation indexing**

→ Summary → {}

**Parent Document, Dense X**

Convert documents into compact retrieval units (e.g., a summary)

**Specialized Embeddings**

🧠 → [0.1, …]

**Fine-tuning, ColBERT**

Domain-specific and / or advanced embedding models

**Hierarchical Indexing**

Splits → Cluser → Summaries → …
→ Cluser →

**RAPTOR**

Tree of document summarization at various abstraction levels

## Generation

**Active Retrieval**

{} → 🧠 → Answer

**Self-RAG, RRR**

Use generation quality to inform question re-writing and / or re-retrieval of documents

# General approaches to transform questions

## 3.1 Preliminaries

Dense retrieval models similarity between query and document with inner product similarity. Given a query $q$ and document $d$, it uses two encoder function $enc_q$ and $enc_d$ to map them into $d$ dimension vectors $\mathbf{v_q}, \mathbf{v_d}$, whose inner product is used as similarity measurement.

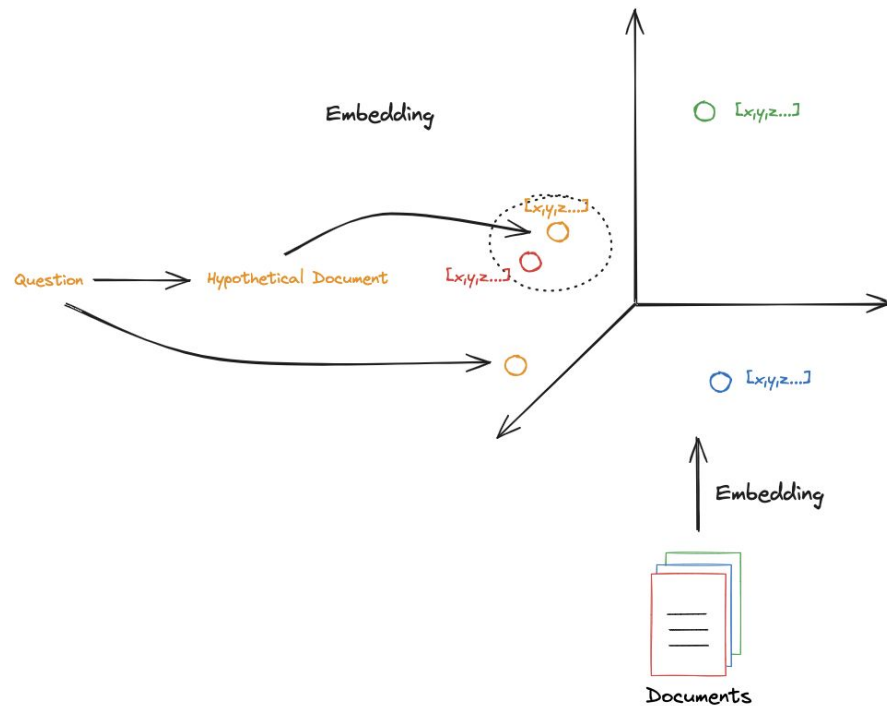$$sim(q, d) = \langle enc_q(q), enc_d(d) \rangle = \langle \mathbf{v_q}, \mathbf{v_d} \rangle \quad (1)$$

For zero-shot retrieval, we consider $L$ query sets $Q_1, Q_2, ..., Q_L$ and their corresponding search corpus, document sets $D_1, D_2, ..., D_L$. Denote the $j$-th query from $i$-th set query set $Q_i$ as $q_{ij}$. We need to fully define mapping *functions* $enc_q$ and $enc_d$ without access to any query set $Q_i$, document set $D_i$, or any relevance judgment $r_{ij}$. The difficulty of zero-shot dense retrieval lies precisely in Equation 1: it requires learning of two embedding functions (for query and document respectively) into the *same* embedding space where inner product captures *relevance*. Without relevance judgments/scores to fit, learning becomes intractable.

## 3.2 HyDE

HyDE circumvents the aforementioned learning problem by performing search in document-only embedding space that captures document-document similarity. This can be easily learned using unsupervised contrastive learning (Izacard et al., 2021; Gao et al., 2021; Gao and Callan, 2022). We set document encoder $enc_d$ directly as a contrastive encoder $enc_{con}$.

# Intuition

Embedding

$[x,y,z...]$

$[x,y,z...]$

Question → Hypothetical Document

$[x,y,z...]$

$[x,y,z...]$

Embedding

Documents

Code walk-through