# RAG (Retrieval Augmented Generation) Cheatsheet

## Stages in RAG:

1. **Loading:**
   - Import your data (text files, PDFs, databases, APIs) using LlamaHub's extensive range of connectors.
2. **Indexing:**
   - Create searchable data structures, primarily through vector embeddings and metadata strategies, enabling efficient context retrieval.
3. **Storing:**
   - Securely store your indexed data and metadata for quick access without the need to re-index.
4. **Querying:**
   - Utilize LLMs and LlamaIndex data structures for diverse querying techniques, including sub-queries and hybrid strategies.
5. **Evaluation:**
   - Continuously assess the effectiveness of your pipeline to ensure accuracy, faithfulness, and response speed.

## Key Concepts:

1. **Nodes and Documents:**
   Fundamental units in LlamaIndex, where Documents encapsulate data sources and Nodes represent data "chunks" with associated metadata.
1. **Connectors:**
   Bridge various data sources into the RAG framework, transforming them into Nodes and Documents.
1. **Indexes:**
   The backbone of RAG, enabling the storage of vector embeddings in a vector store along with crucial metadata.
1. **Embeddings:**
   Numerical representations of data, facilitating the relevance filtering process.
1. **Retrievers:**
   Define efficient retrieval strategies, ensuring the relevancy and efficiency of data retrieval.
1. **Routers:**
   Manage the selection of appropriate retrievers based on query specifics and metadata.
1. **Node Postprocessors:**
   Apply transformations or re-ranking logic to refine the set of retrieved nodes.
1. **Response Synthesizers:**
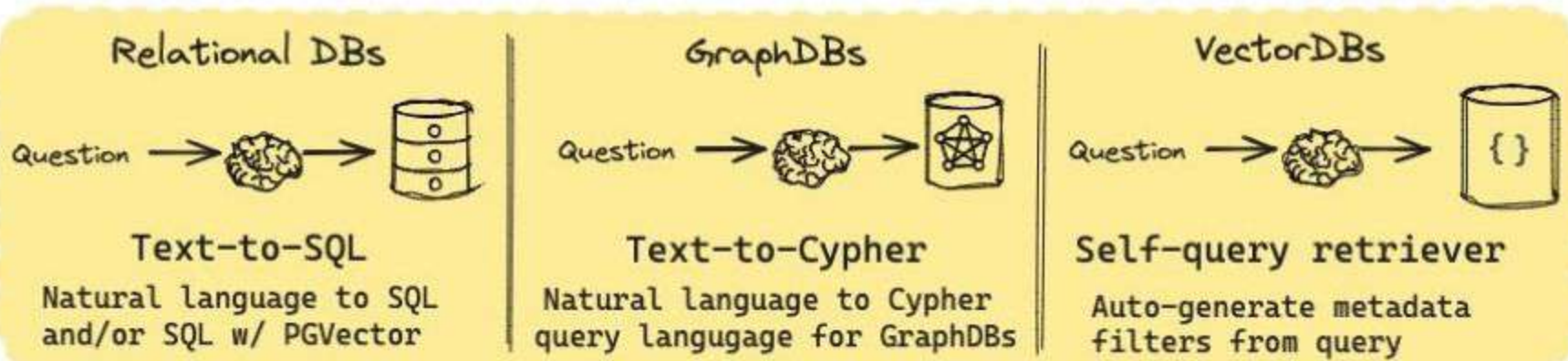   Craft responses from the LLM, utilizing user queries and retrieved text chunks for enriched answers.

## Application Types:

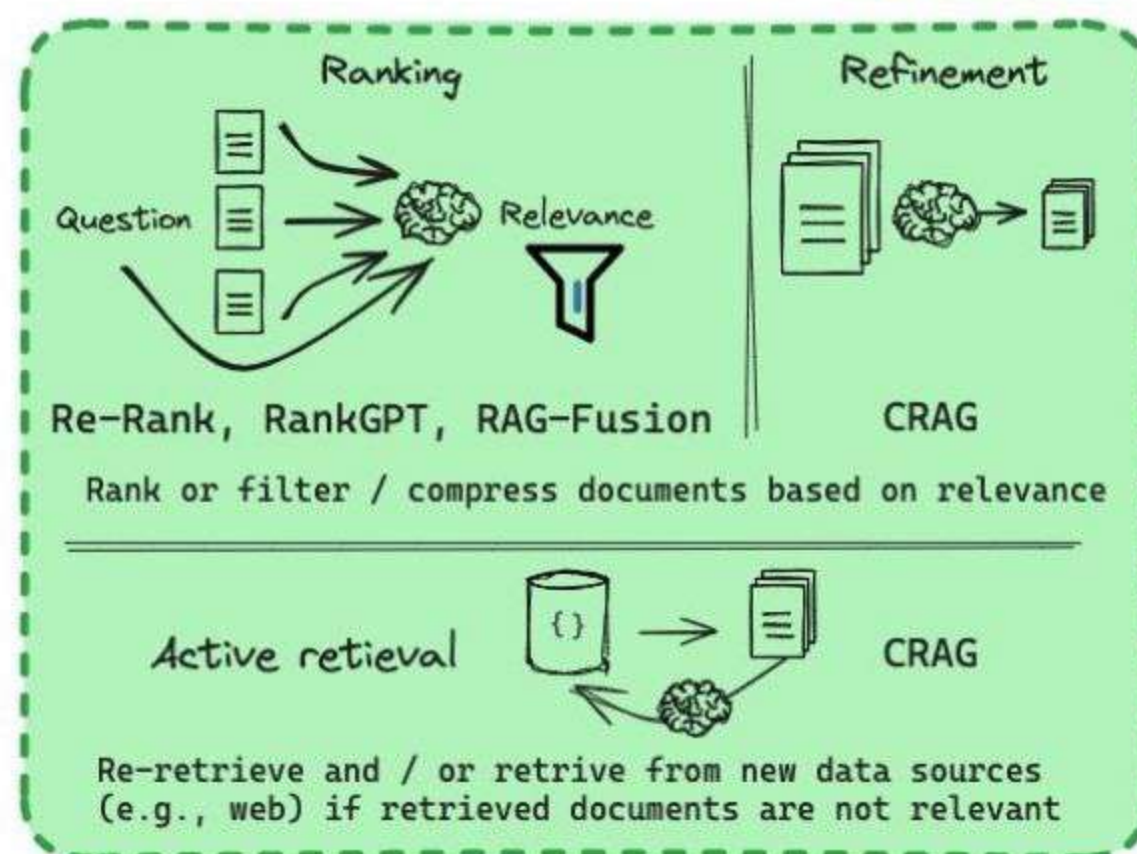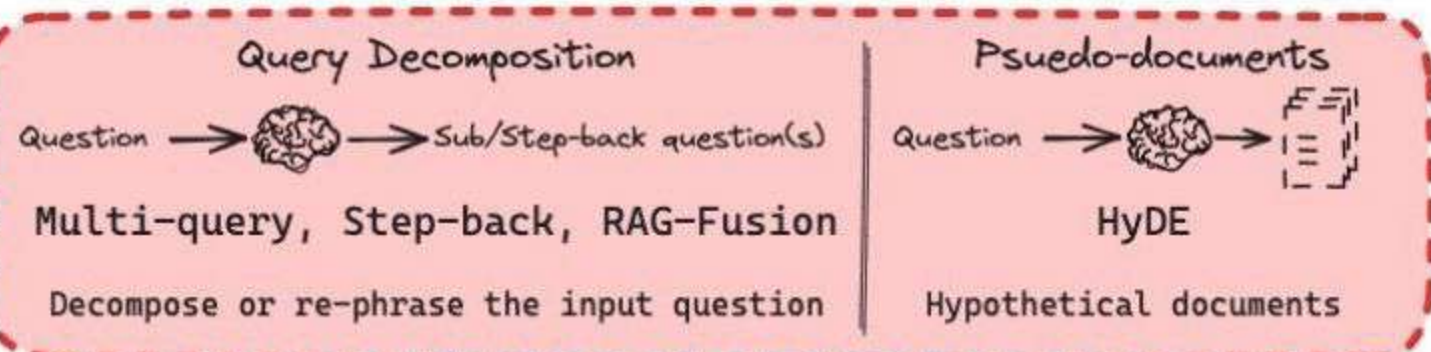1. **Query Engines:**
   - For direct question-answering over your data.
2. **Chat Engines:**
   - Enables conversations with your data for an interactive experience.
3. **Agents:**
   - Automated decision-makers that interact with external tools, adaptable for complex tasks.

**LlamaIndex**  **SingleStore**
**LangChain**  **OpenAI**

## Query Construction

| Relational DBs | GraphDBs | VectorDBs |
|---|---|---|
| Question → 🧠 → 🗄 | Question → 🧠 → graph | Question → 🧠 → {} |
| **Text-to-SQL** | **Text-to-Cypher** | **Self-query retriever** |
| Natural language to SQL and/or SQL w/ PGVector | Natural language to Cypher query language for GraphDBs | Auto-generate metadata filters from query |

## Query Translation

| Query Decomposition | Psuedo-documents |
|---|---|
| Question → 🧠 → Sub/Step-back question(s) | Question → 🧠 → docs |
| **Multi-query, Step-back, RAG-Fusion** | **HyDE** |
| Decompose or re-phrase the input question | Hypothetical documents |

## Retrieval

| Ranking | Refinement |
|---|---|
| Question → 🧠 Relevance → filter | docs → 🧠 → doc |
| **Re-Rank, RankGPT, RAG-Fusion** | **CRAG** |

Rank or filter / compress documents based on relevance

**Active retrieval**
{} → 🧠 → docs  **CRAG**

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

**Diagram credit Langchain**

## Routing

| Logical routing | Semantic routing |
|---|---|
| 🧠 → DBs | 🧠 Embed → Prompt #1 / Prompt #2 |
| Let LLM choose DB based on the question | Embed question and choose prompt based on similarity |

Question → (pink) → (orange) → Graph DB / Relational DB / Vectorstore → Documents → filter → 🧠 → Answer

## Indexing

| Chunk Optimization | Multi-representation indexing | Specialized Embeddings | Heirachical Indexing |
|---|---|---|---|
| Split → Charecters / Sections / Semantic / Delimiters | doc → Summary → {}/🗄 | doc → 🧠 → [0.1, ...] | Splits → Summaries → Cluser → ... / Cluser |
| **Semantic Splitter** | **Parent Document, Dense X** | **Fine-tuning, ColBERT** | **RAPTOR** |
| Optimize chunk size used for embedding | Convert documents into compact retrieval units (e.g., a summary) | Domain-specific and / or advanced embedding models | Tree of document summarization at various abstraction levels |

## Generation

**Active retrieval**

{} → doc → 🧠 → Answer

**Self-RAG, RRR**

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

Steve Nouri