




# Information Retrieval and Synthesis Workflow with Gen AI

This file is meant for personal use by [venkhatbalaji@gmail.com](mailto:venkhatbalaji@gmail.com) only.  
Sharing or publishing the contents in part or full is liable for legal action.

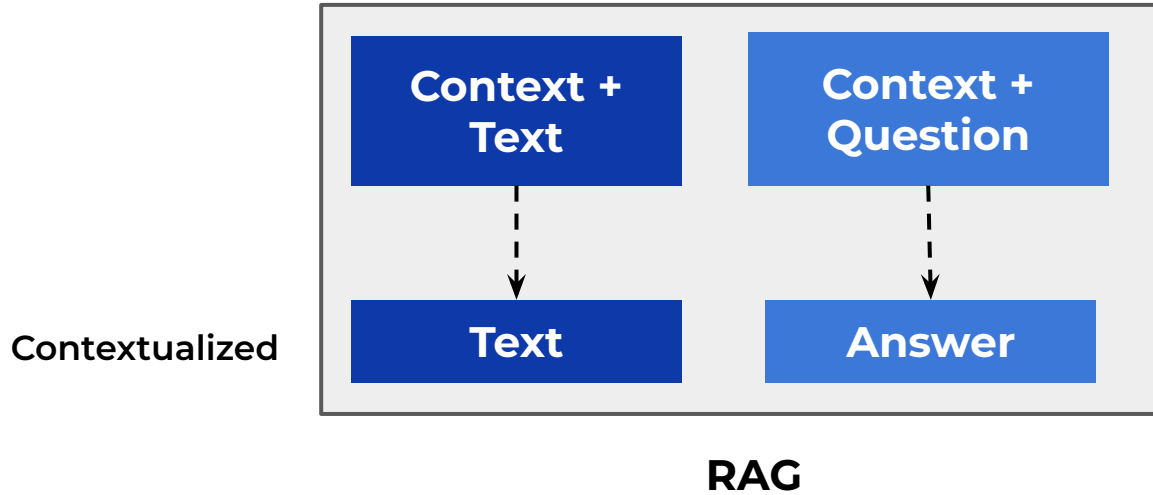


# Agenda

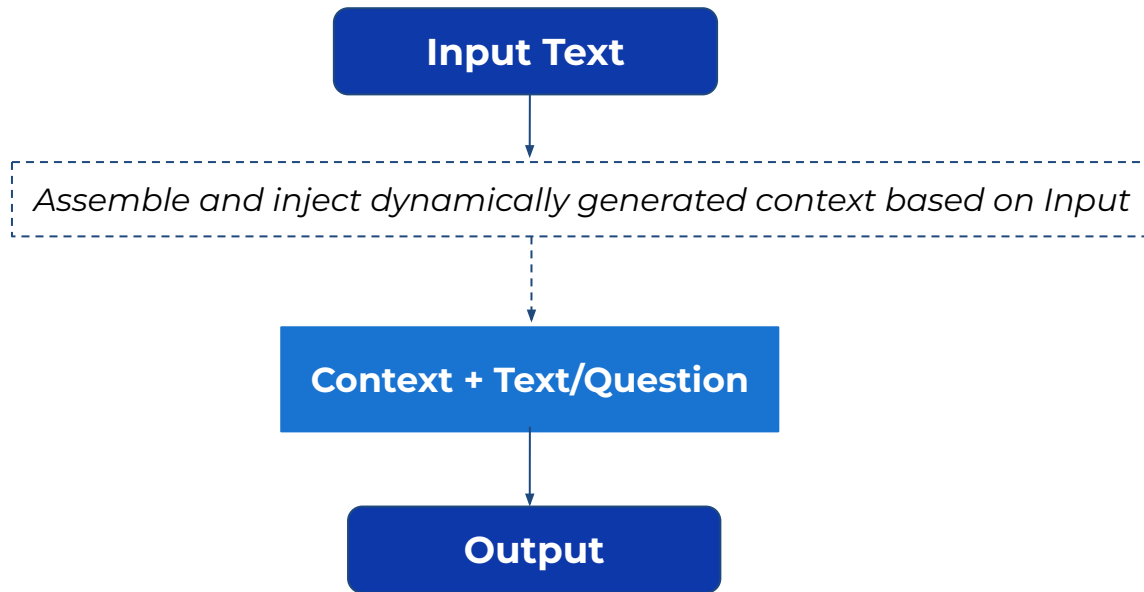
In this session, we will discuss :

- Overview of Retrieval Augmented Generation (RAG) and its Working
- Building Blocks of RAG
- Data Preparation Process with respect to RAG
- Devising and Evaluating Prompts with respect to RAG

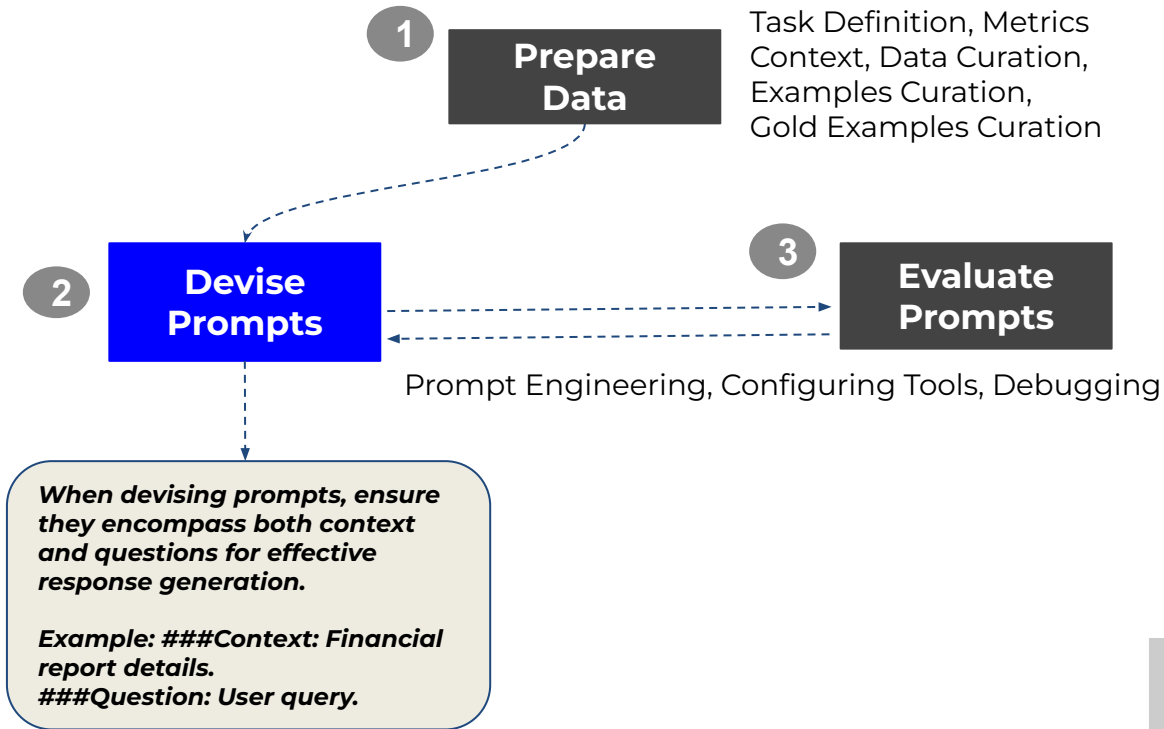
# Retrieval Augmented Generation (RAG)



# Working of RAG



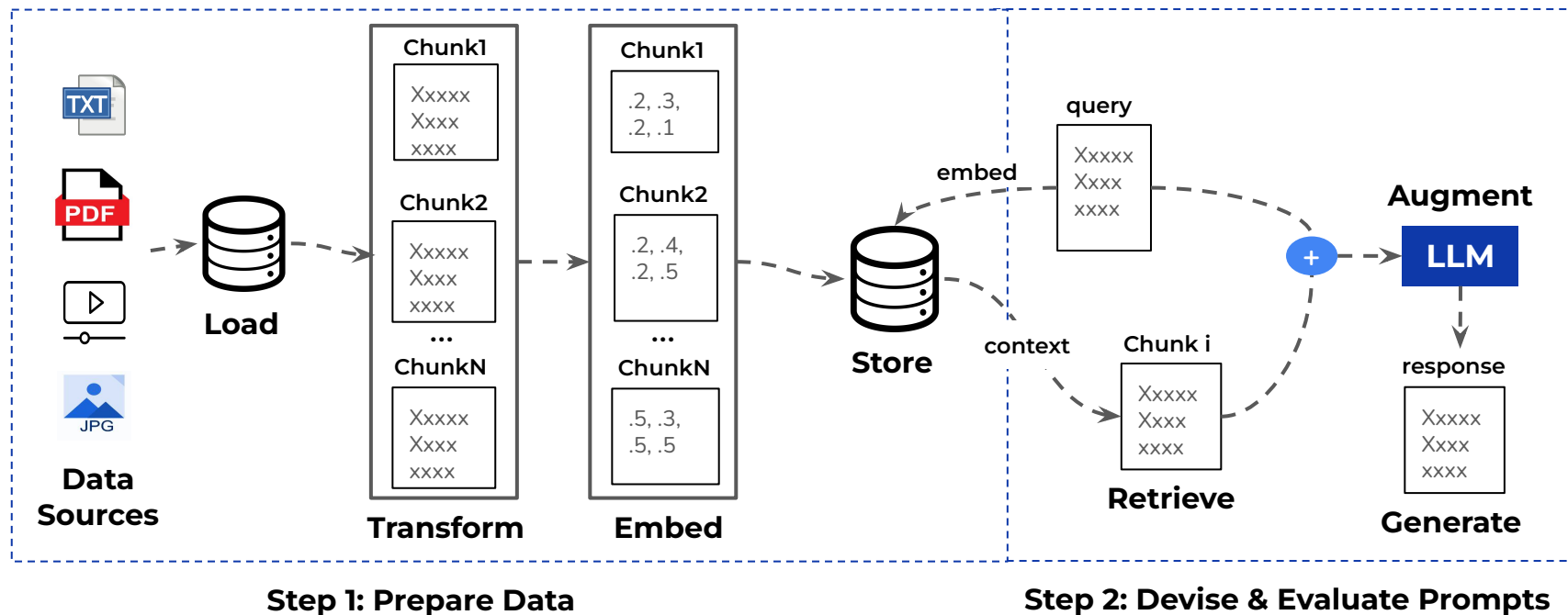
# Structure of RAG



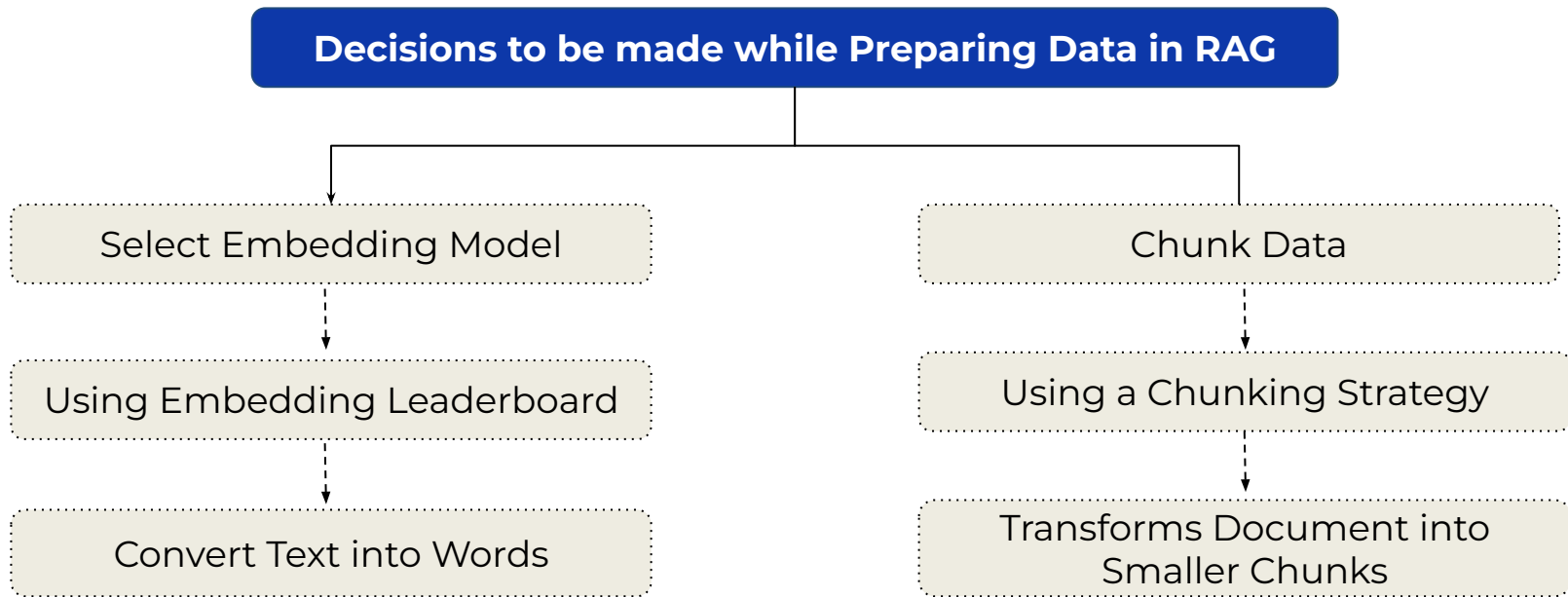
# Building Blocks of RAG

This file is meant for personal use by venkhatbalaji@gmail.com only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Building Blocks of RAG

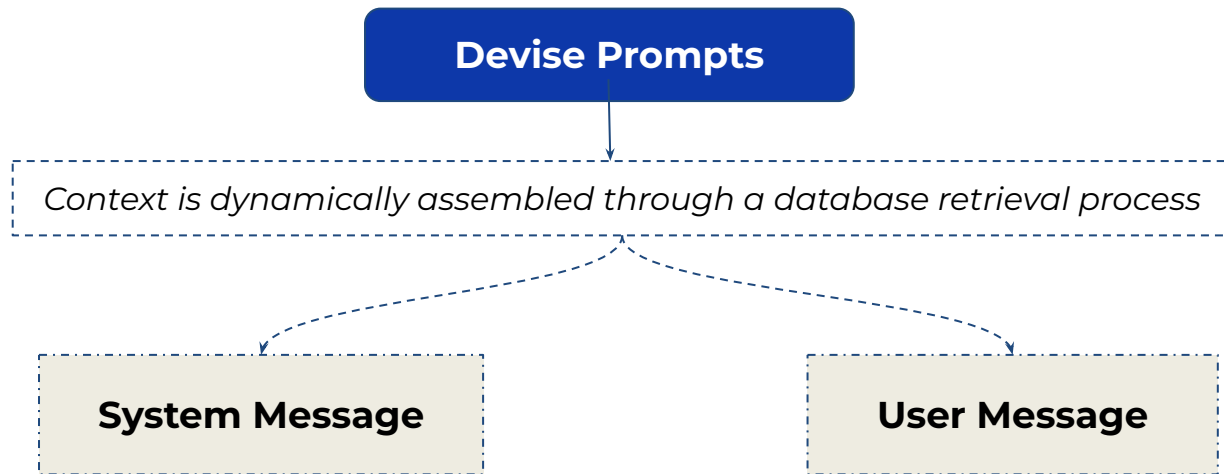


# Prepare Data in RAG





# Devise Prompts in RAG



# Evaluate Prompts in RAG

Evaluate Prompts

Accuracy

**Assess the effectiveness of prompts used in RAG tasks.**

**Factors:**

**Clarity:** *How clear is the prompt in conveying the task?*

**Relevance:** *Is the response relevant to the query posed by the user?*

**Faithfulness to the context:** *Is the context used correctly to create the response?*

**Ensure prompts facilitate accurate and meaningful model predictions.**

# Data Preparation Process

This file is meant for personal use by [venkhatbalaji@gmail.com](mailto:venkhatbalaji@gmail.com) only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Select Embedding Model

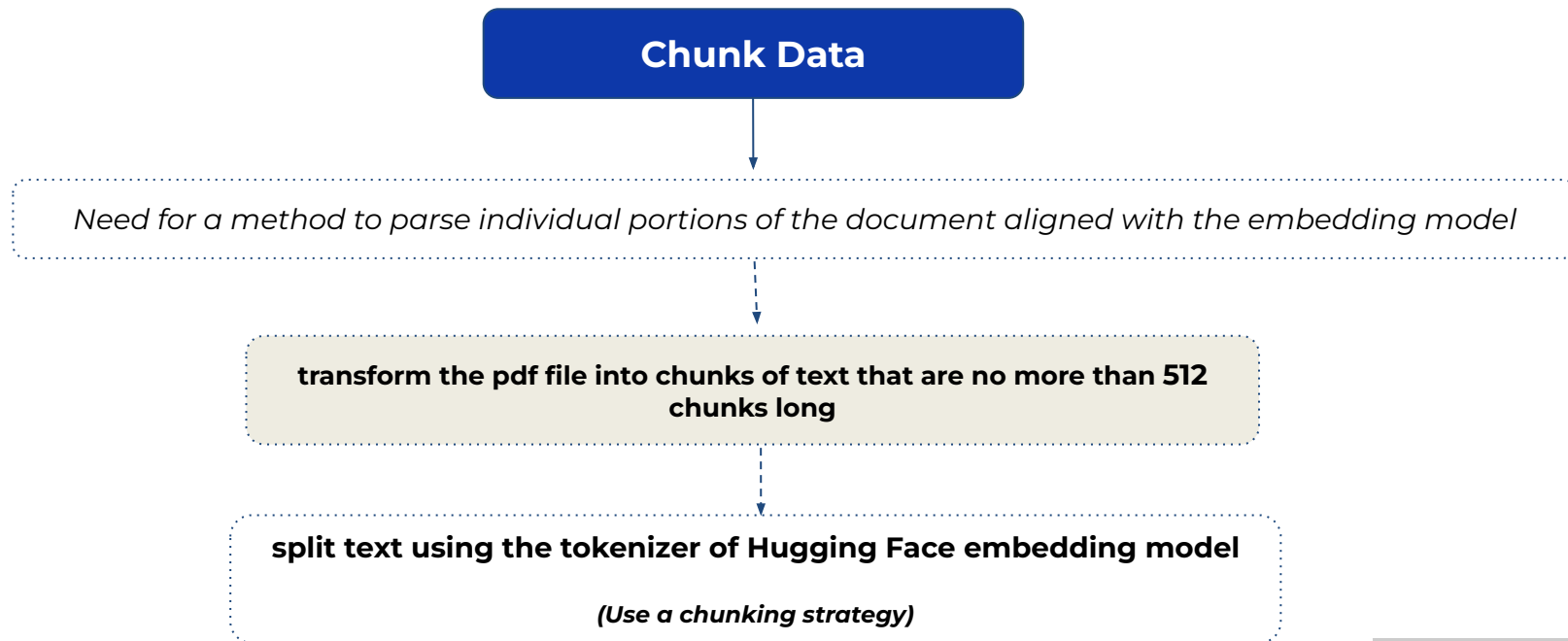
## Embedding Model

*Encodes text into vector representations that act as good features for LLM retrieval tasks*

**Selecting an open source model from Embedding Leaderboard**  
*(To make this choice, look at the task to solve and then choose the embedding model close to Open AI `text-embedding-ada-002` on the leaderboard)*

**create a vectorized representation of the user\_input by using the `embed_query` method**

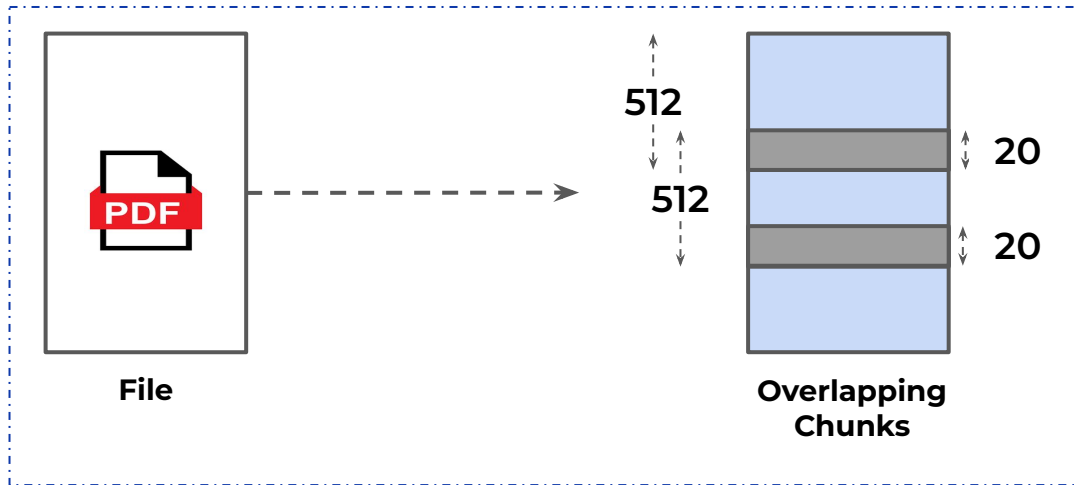
# Chunk Data



# Chunking Strategy: Example

`chunk_overlap = 20`

*ensures that the chunks are related to each other (i.e., there is some continuity between the chunks)*



This file is meant for personal use by [venkhatbalaji@gmail.com](mailto:venkhatbalaji@gmail.com) only.  
Sharing or publishing the contents in part or full is liable for legal action.

# Create Vector Database

**Create vector database**

*Generate a vector for each chunk and save this chunk along with the vector representation*


**To add embeddings data to the database, create an index and push the embeddings by chunk to the index**

**Important components of the index to be specified during creation**

**Dimension of the embedding  
generated by the embedding  
model**


**Metric used to define the similarity  
between a pair of documents**

*(E.g., - Cosine similarity for indexing text)*



# Devising and Evaluating Prompts

This file is meant for personal use by [venkhatbalaji@gmail.com](mailto:venkhatbalaji@gmail.com) only.  
Sharing or publishing the contents in part or full is liable for legal action.





# Devising Prompts

## Prompt Design

*Context is dynamically assembled through a database retrieval process*

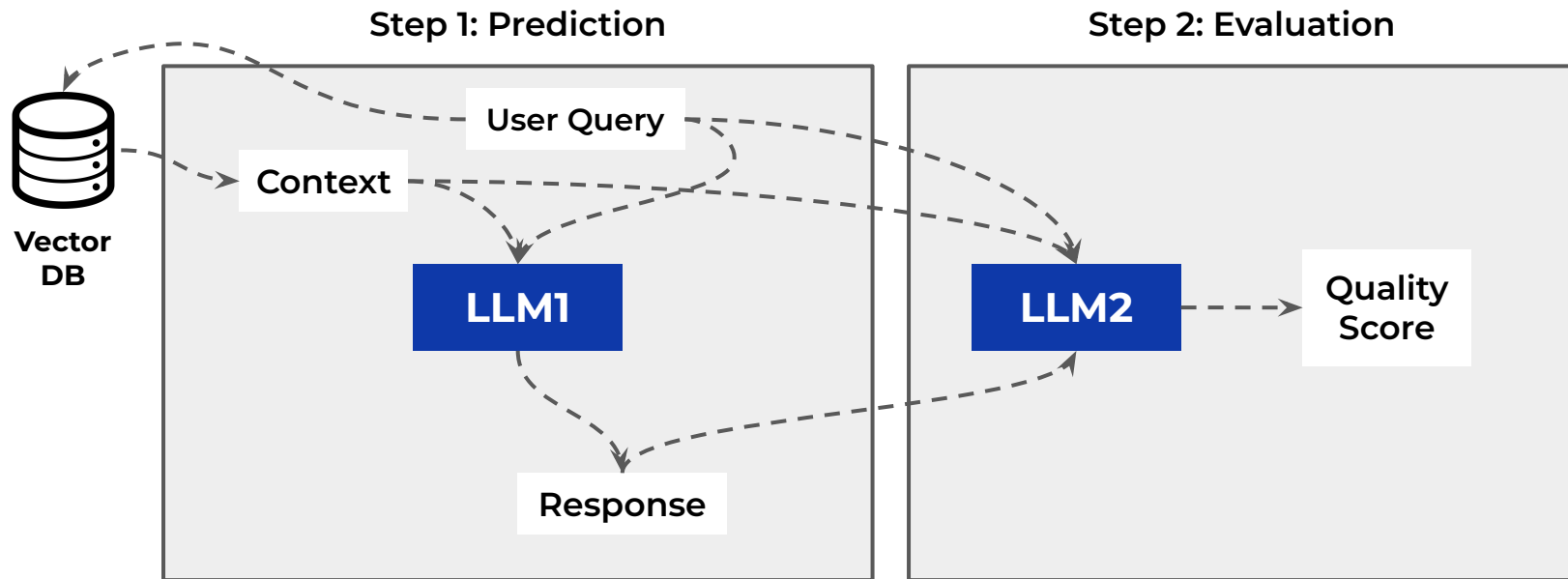
### System Message

*Here we provide a distinct set of instructions regarding the task*

### User Message

*Here we clearly define the sections where the context will be inserted and where the user input will be injected*

# Evaluation Process in RAG



# Summary

