

DEEP LEARNING IN BUSINESS



RAG-BASED Q&A SYSTEM

PRESENTED BY GROUP 1

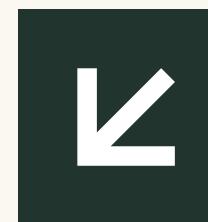


TABLE OF CONTENTS

INTRODUCTION



PROPOSED SOLUTION: RAG



SYSTEM ARCHITECTURE OVERVIEW



PHASE 1: DATA INGESTION & EMBEDDING



PHASE 2: QUERY & ANSWER GENERATION



PHASE 3: QUERY & ANSWER GENERATION



RESULTS & EVALUATION



Q&A

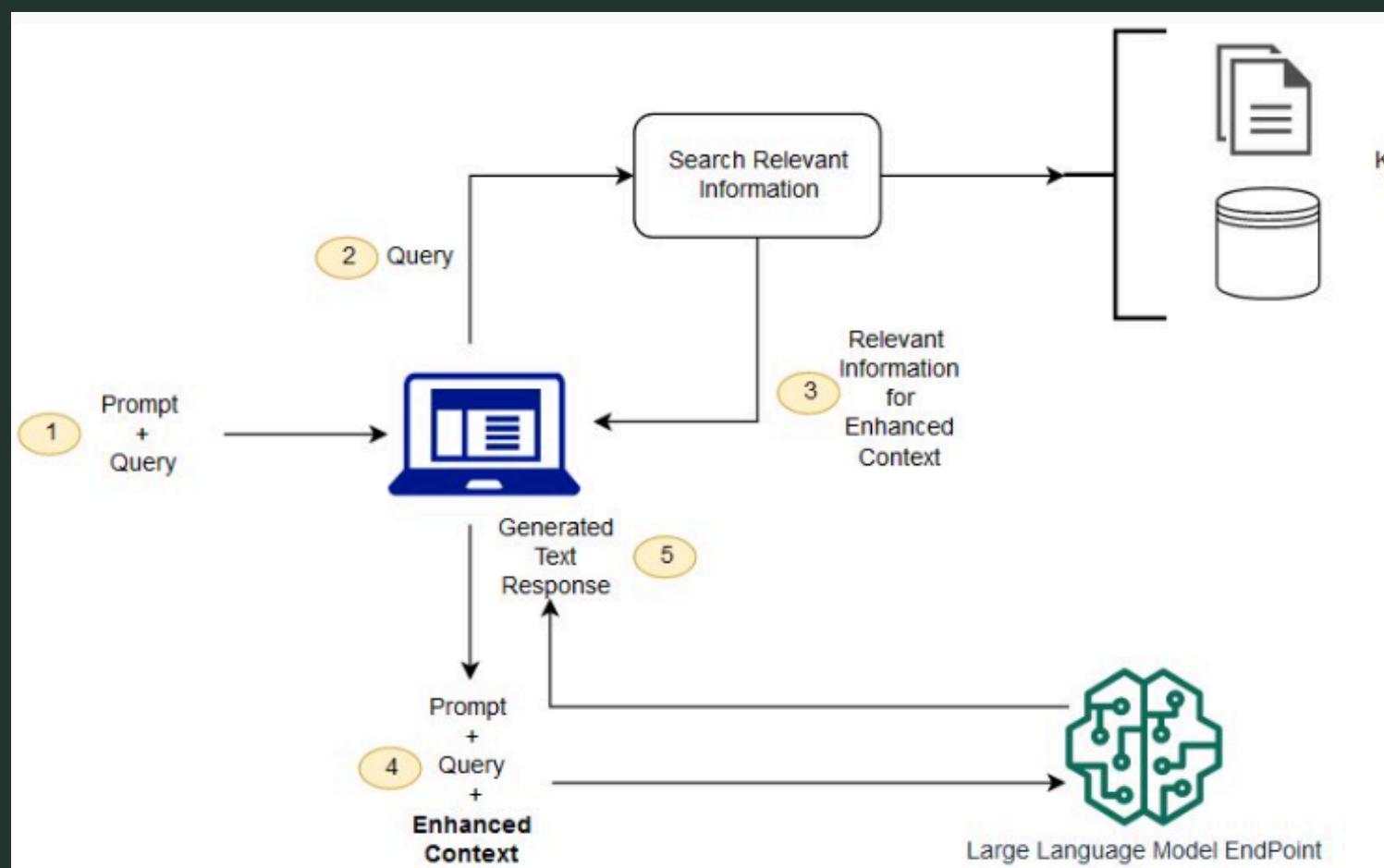


INTRODUCTION

- The Challenge: Information overload from large, specialized PDF documents.
 - The Limitation: Traditional keyword search (Ctrl+F) is ineffective for complex, semantic questions.
- > The Need: A conversational tool to ask questions and receive accurate, fast answers directly from the document's content.



PROPOSED SOLUTION: RAG



Solution: Using the **RAG** (Retrieval-Augmented Generation) model.

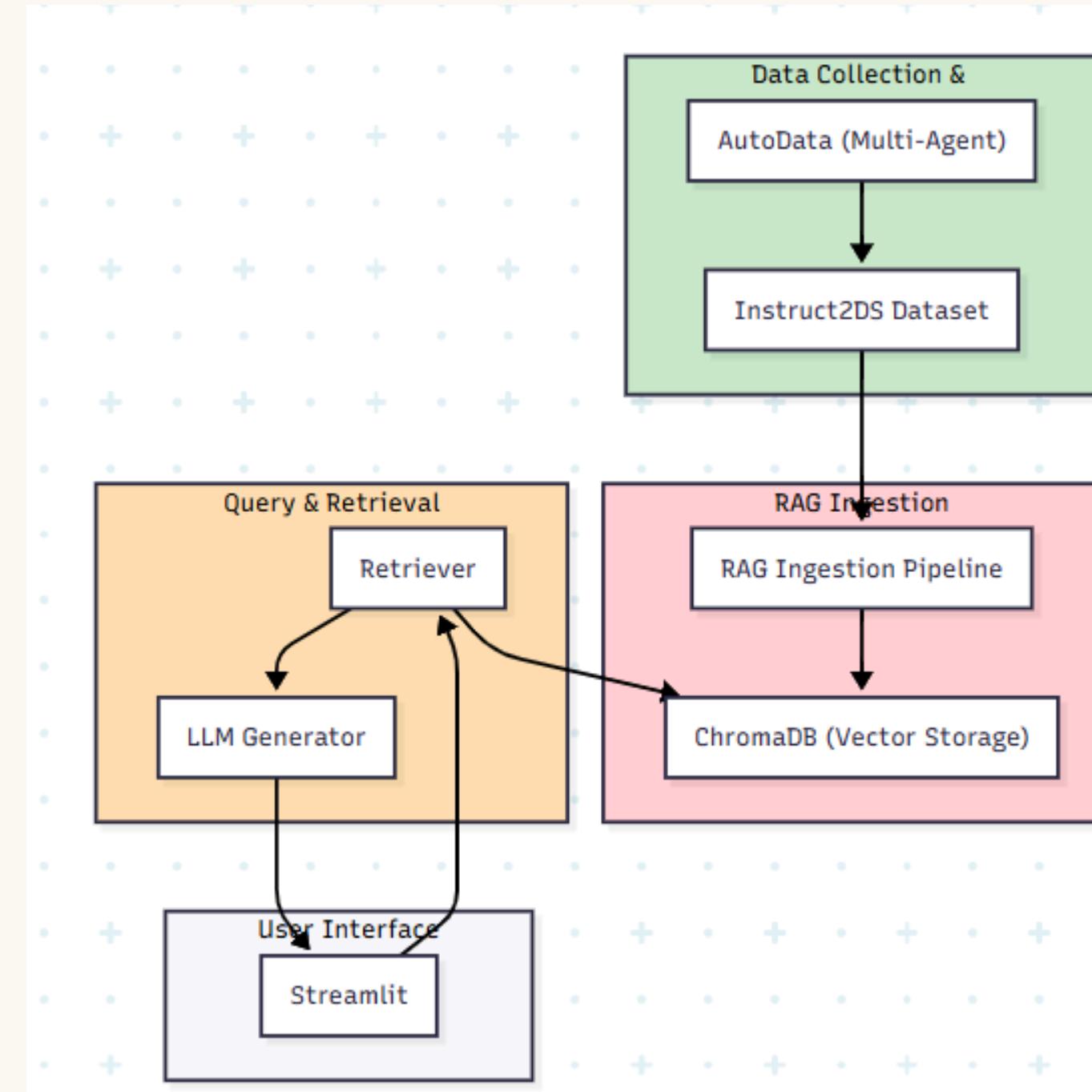
What is RAG? A technique that combines the strengths of:

- **Retrieval:** Searching and fetching the most relevant information from a database (the PDF content).
- **Generation:** Using a Large Language Model (LLM) to generate a natural answer based on the retrieved information.

Why we choose RAG?

Advantages: Ensures the answer is grounded in the document, reduces the risk of the LLM 'making up' information, and allows us to cite the source.

SYSTEM ARCHITECTURE OVERVIEW



FIRST PHASE

Data Ingestion & Embedding Phase

- Collect and prepare data
- Clean and normalize text
- Chunk documents
- Generate vector embeddings
- Store embeddings in a Vector Database
- Initialize the RAG pipeline

SECOND PHASE

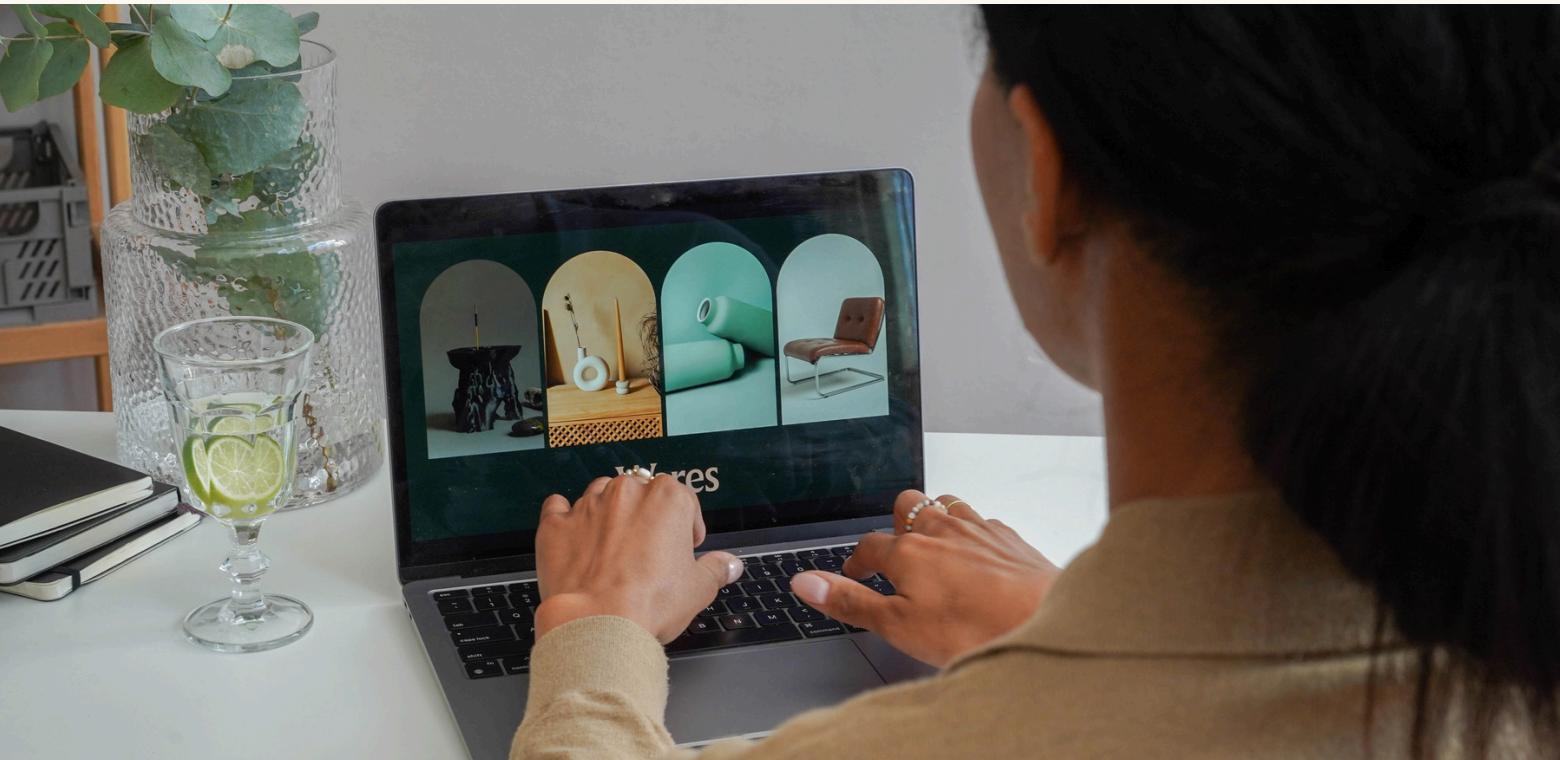
Query & Answer Generation Phase

- User submits a query
- Perform semantic retrieval
- Combine retrieved context with the LLM
- Generate a concise, accurate answer

PHASE 1:

Data Ingestion & Management Phase

Objective: Build a Knowledge Base from PDFs.



Load & Clean

Load, clean, and normalize data (from Instruct2DS).

metadata.json

pdf downloader

pdf parser



Text Chunking

Split the text into meaningful "chunks".

text_chunker



Embedding

Convert each chunk into a numerical vector (using an embedding model).

model gte-Qwen2-1.5B-instruct



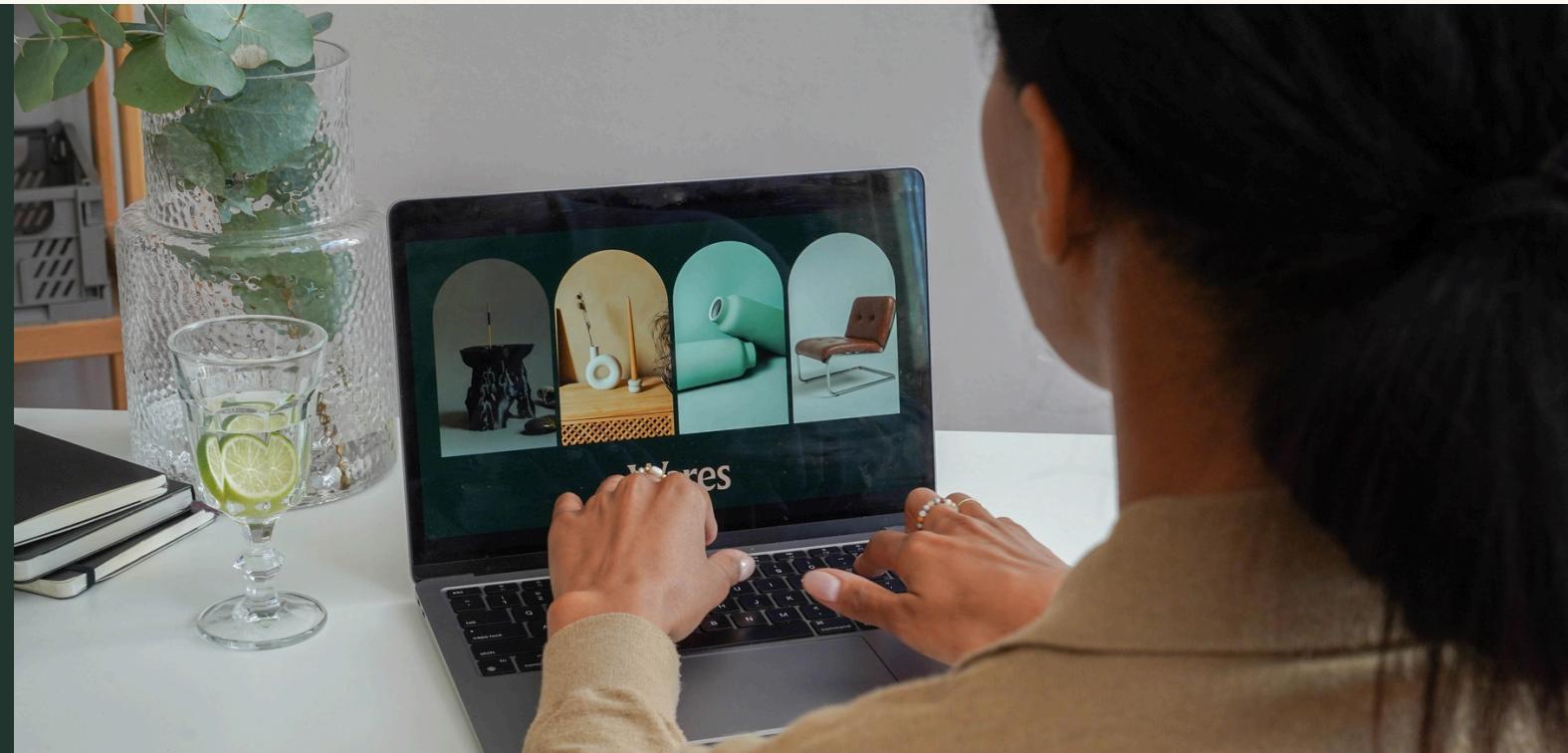
Vector Store

Store all vectors in ChromaDB.

PHASE 2: Query & Answer Generation

Objective: Find relevant text chunks from ChromaDB

→ Use those chunks to generate the final answer.



User query

Embed the user query

Query vector

Search ChromaDB (top-k)

Retrieve chunks

Return context chunks

Rag prompt

Build RAG prompt
(context + query)

Final answer

Call LLM → generate
final answer

PHASE 3: Evaluation and Fine-tuning

Objective: measure the retrieval quality and question answering ability of the RAG system using the current embedding..

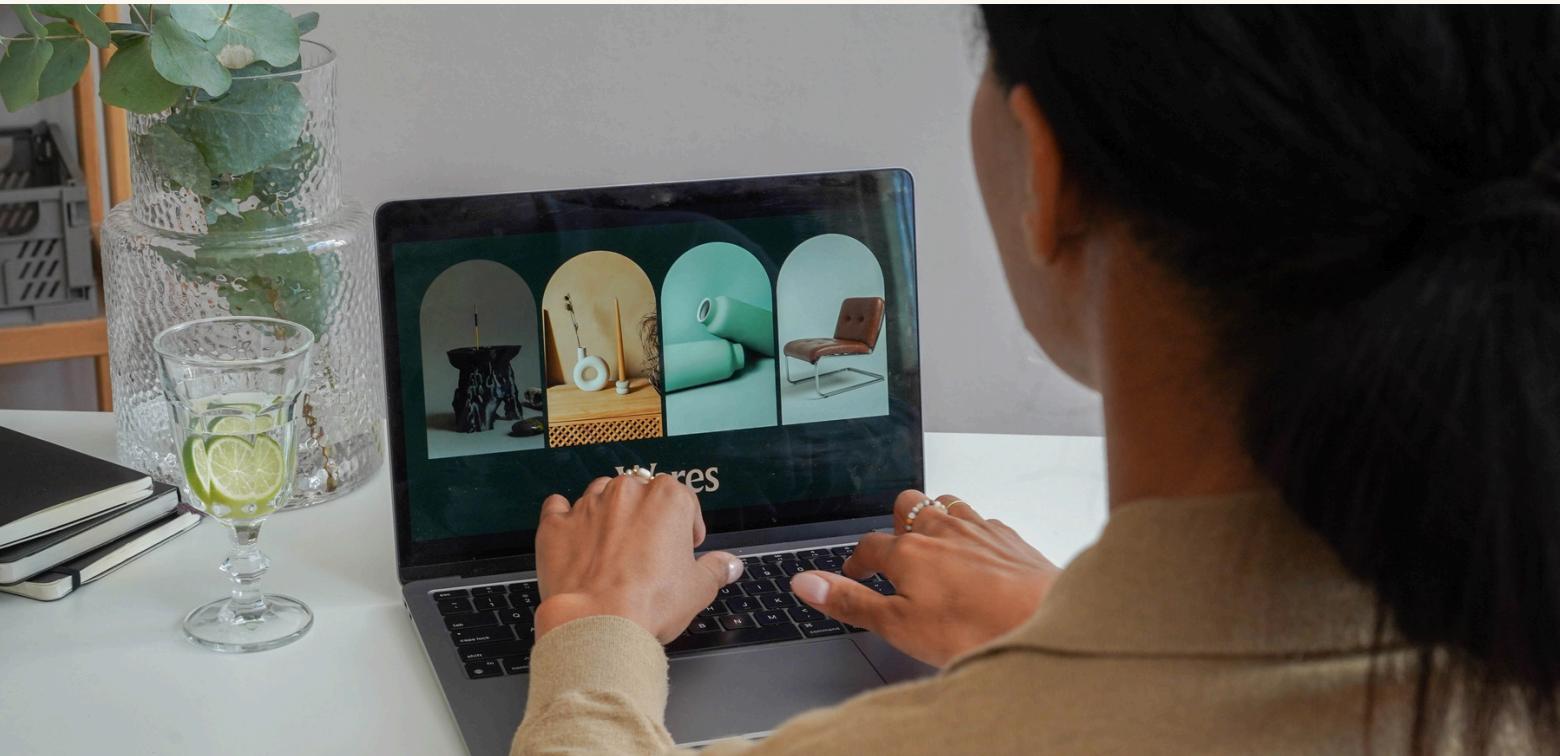


Evaluation



PHASE 3: Evaluation and Fine-tuning

Objective: measure the retrieval quality and question answering ability of the RAG system using the current embedding..



Fine-tuning



RESULT

The screenshot shows a user interface for a RAG-Based Question Answering system. On the left, there is a sidebar titled "Thiết lập truy vấn" (Query Setup) with sliders for "Top-k" (set to 5), "Max token" (set to 4000), and "Temperature" (set to 0.10). There is also a checkbox for "Show detail trace". Below this, under "Quy trình" (Process), a numbered list details the workflow:

- Câu hỏi được nhúng bằng mô hình [Alibaba-NLP/gte-Qwen2-1.5B-instruct](#).
- ChromaDB truy xuất top-k đoạn văn dựa trên cosine similarity.
- Các đoạn được nén để phù hợp giới hạn context và gửi tới [qwen2.5:7b](#) (Ollama).
- Câu trả lời được hậu xử lý với nguồn trích dẫn rõ ràng.

The main area is titled "RAG-Based Question Answering for Online PDF Documents". It includes a note about the system using Retrieval-Augmented Generation with LangGraph, ChromaDB, embedding [Alibaba-NLP/gte-Qwen2-1.5B-instruct](#), and a Qwen2.5:7b model running on Ollama. A text input field contains "hello", and a large red button labeled "Truy vấn" (Search) is visible. Below the search bar, tabs for "Kết quả", "Ngữ cảnh", "Phân tích", and "Lịch sử" are present. The result section displays the message: "Câu trả lời: KHÔNG NẮM TRONG NGỮ CẢNH".

RESULT

The screenshot shows a user interface for a RAG-based question answering system. On the left, there's a sidebar titled "Thiết lập truy vấn" (Query Setup) with sliders for "Top-k" (set to 5), "Max token" (set to 4000), and "Temperature" (set to 0.10). There's also a checkbox for "Show detail trace". Below this is a section titled "Quy trình" (Process) containing a numbered list of steps:

- Câu hỏi được nhúng bằng mô hình [Alibaba-NLP/gte-Qwen2-1.5B-instruct](#).
- ChromaDB truy xuất top-k đoạn văn dựa trên cosine similarity.
- Các đoạn được nén để phù hợp giới hạn context và gửi tới [qwen2.5:7b](#) (Ollama).
- Câu trả lời được hậu xử lý với nguồn trích dẫn rõ ràng.

The main area has a title "RAG-Based Question Answering for Online PDF Documents" and a sub-section "Transformer có được dùng trong NLP không?" (Is the Transformer used in NLP?). A large red button labeled "Truy vấn" (Search) is prominent. Below it, a result card displays the query "Câu trả lời: Có, Transformer đã được sử dụng rộng rãi trong lĩnh vực xử lý ngôn ngữ tự nhiên (NLP)." and navigation tabs "Kết quả" (Results), "Ngữ cảnh" (Context), "Phân tích" (Analysis), and "Lịch sử" (History).

THANK YOU