# Multimodal RAG systems

See beyond words

AtlasIA أطلسيـــــة

MoroccoAI

# Learning Outcomes

You will:

- Understand what is a RAG, when it's used and its core components (embeddings, vector stores, …)
- Understand multimodality
- Know a few methods to build a multimodal RAG (focus on image and text)

What we won't cover:

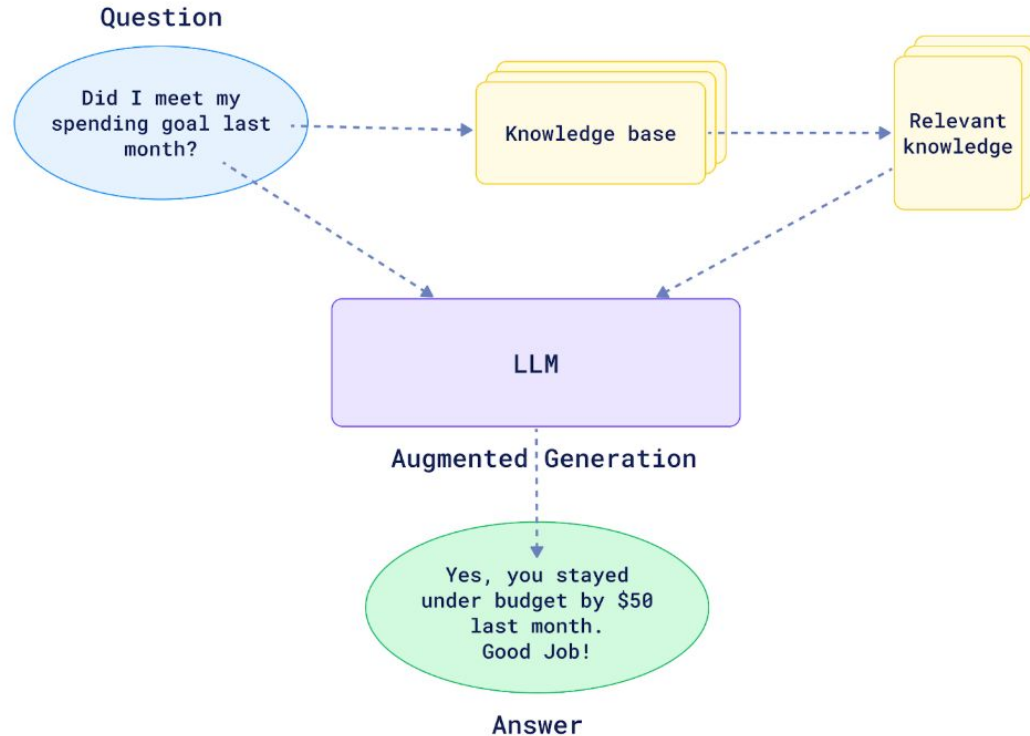- Transformers 😥
- Maths behind components

# Agenda

- Why use RAG?
- Architecture and core components
- What is multimodality?
- Multimodal RAG design approaches

# Retrieval Augmented Generation
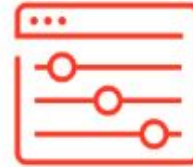
- Why and When do we need it?

# Why use RAG ?

Question

Did I meet my spending goal last month?

Knowledge base

Relevant knowledge

LLM

Augmented Generation

Yes, you stayed under budget by $50 last month. Good Job!

Answer

# Why use RAG ?



Prompt engineering · Retrieval augmented generation (RAG) · Fine-tuning · Pre-train from scratch

Complexity/Compute-intensiveness

# RAG in a nutshell

# RAG components

**Embeddings**

Information from documents is stored as vector embeddings. This format supports efficient similarity searches to retrieve relevant data for your query.

**Vector Search**

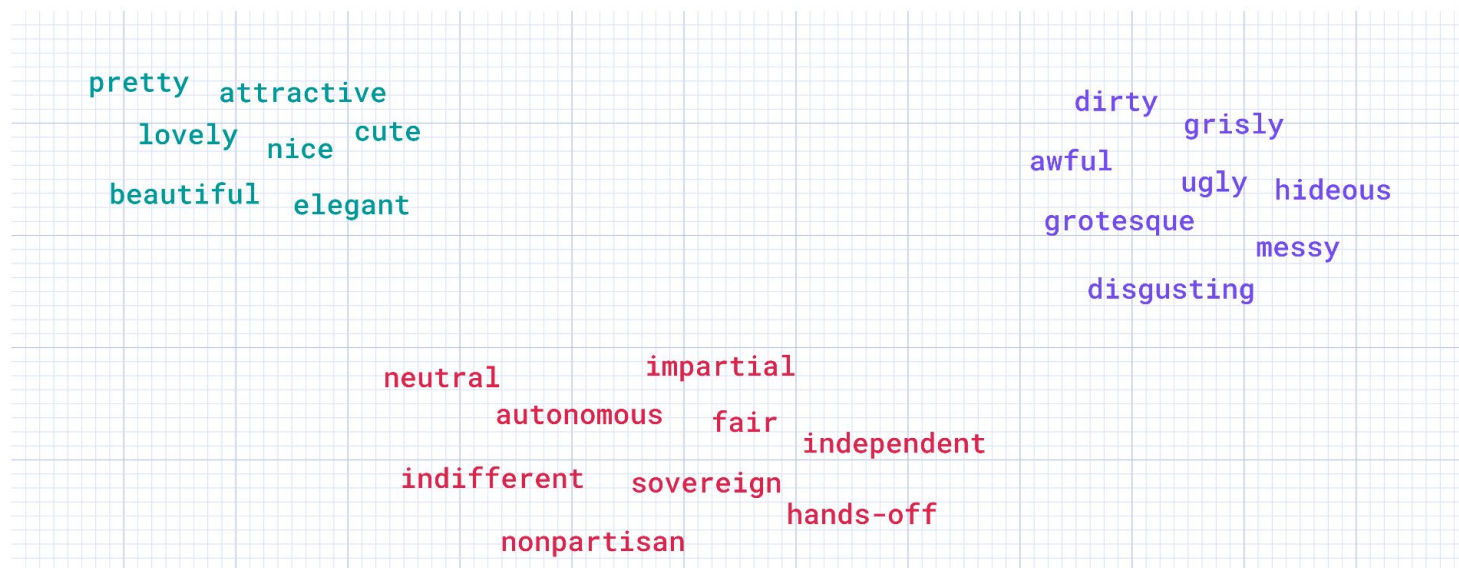Similarity search is applied on the vector database to retrieve the most relevant documents to the query

**LLM**

Takes the query augmented by the relevant context and generates the final answer

# Embeddings

# Embeddings

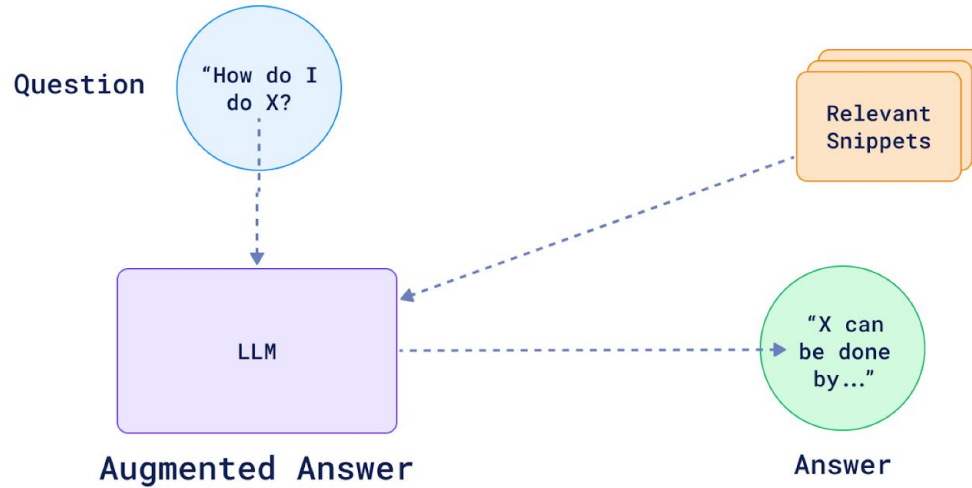Similar objects are nearby while different objects are distant from each other in the vector space.

pretty attractive
lovely nice cute
beautiful elegant

dirty grisly
awful ugly hideous
grotesque
messy
disgusting

neutral impartial
autonomous fair
independent
indifferent sovereign
hands-off
nonpartisan

# RAG in a nutshell: Indexing

# RAG in a nutshell: Retrieval

**Retrieval**

"How to do X?"

Question

Embedding Machine

Embedding

Vector Database

Relevant Snippets

# RAG in a nutshell: Generation

# RAG architecture (in real-world applications)

# Summary (Part I)

- **RAG**: A system that combines retrieval of relevant information from external sources with generation using an LLM to create contextually enriched outputs.

- **Workflow:**
    - Index external documents (textual content) in the vector store
    - Embed the user's query
    - Search for the most similar documents
    - Concatenate the relevant documents to the query and generate an answer

# How to handle documents with more than text?

# What is multimodality 🤔?

# What is multimodality 🤔?

# How to build a RAG that handles images and text?

# 1st approach: Ground all modalities into one primary modality

# 1st approach: Ground all modalities into one primary modality

# 2nd approach: Embed all modalities into one vector space

# 3rd approach: Have separate stores per modality
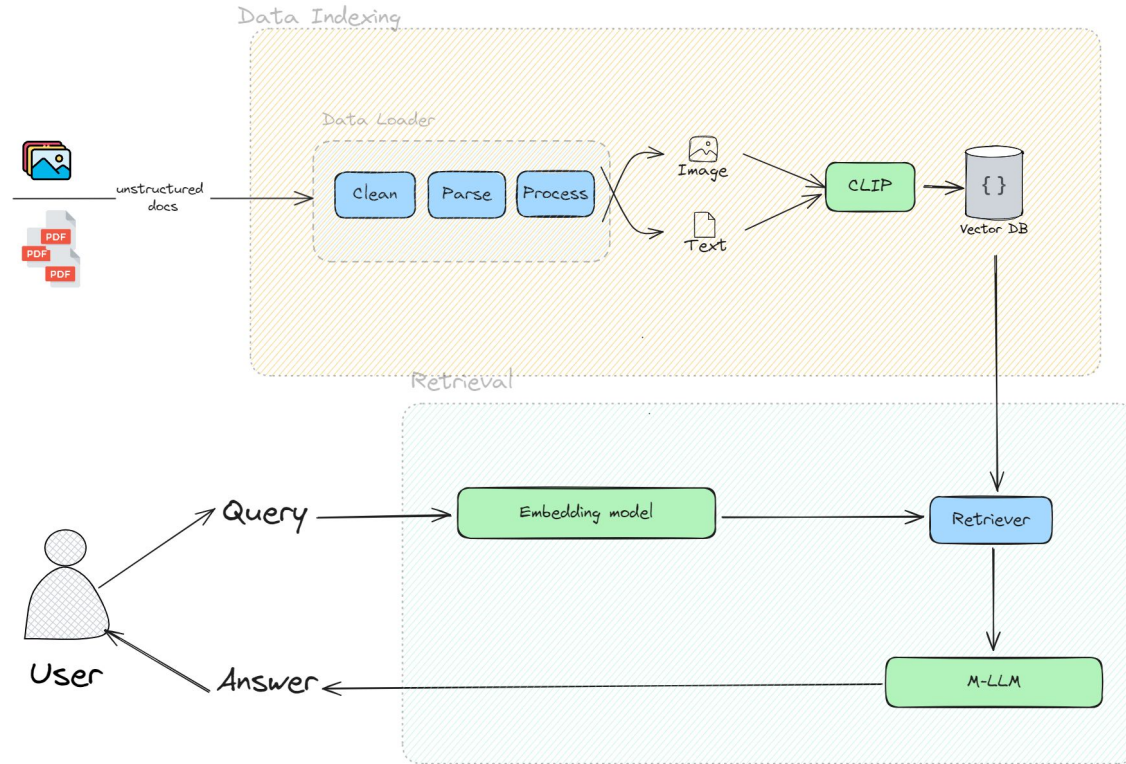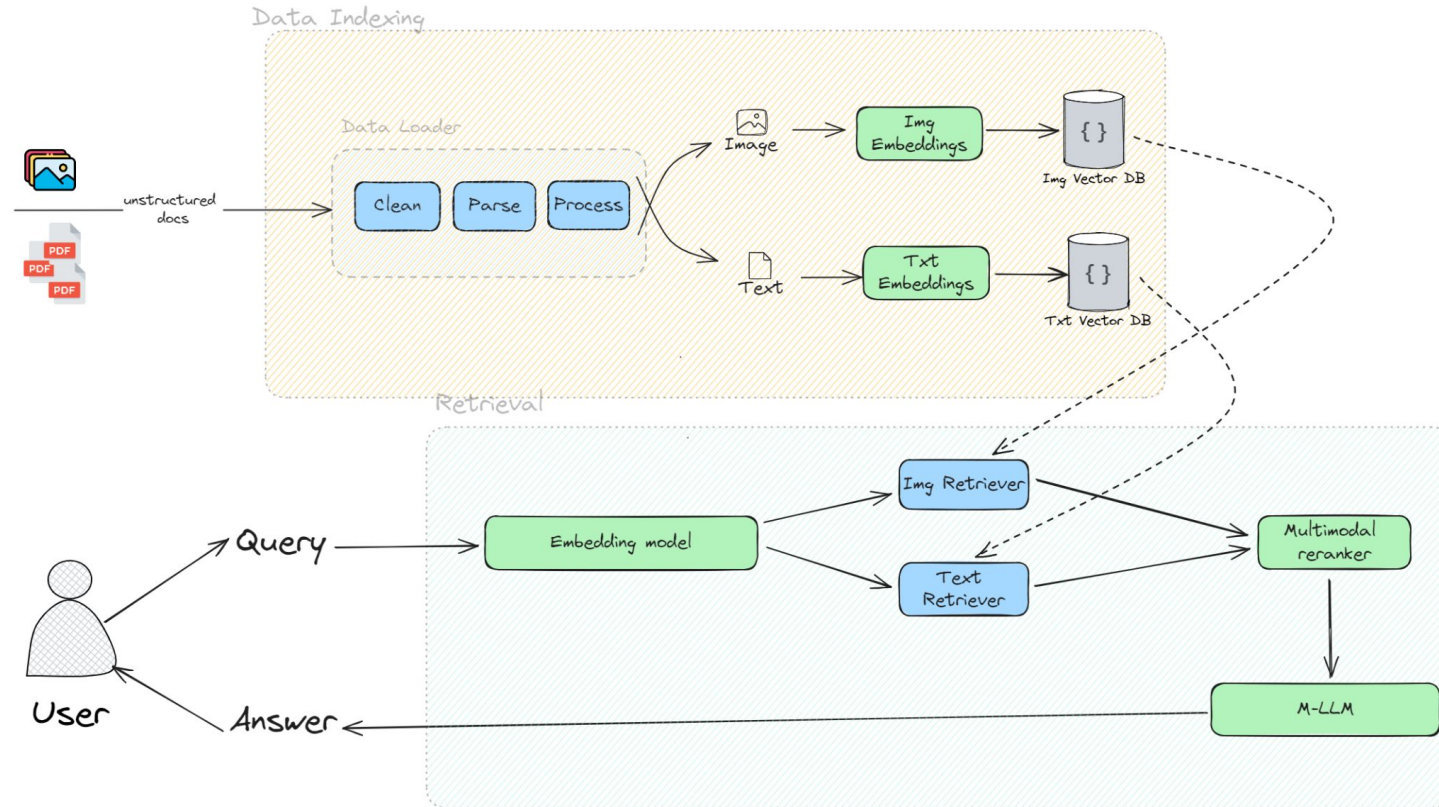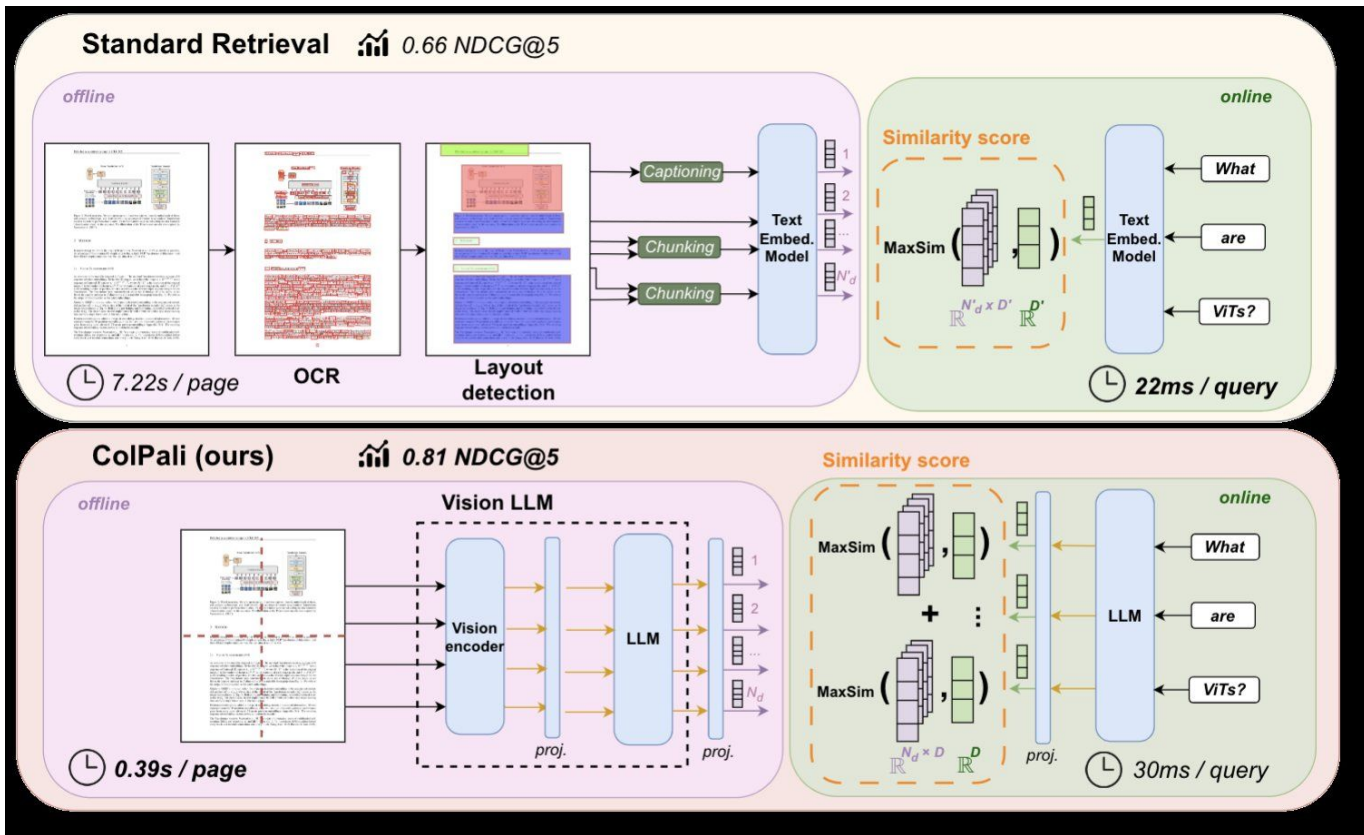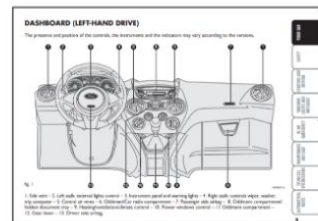
# 4th approach: ColPali

# Summary

- No one-size fits all method:
    - The choice depends on the use case, type of documents, system constraints etc
    - For eg. CLIP models typically offer only generic insights into objects and shapes without providing detailed explanations and do not allow more than 80 words.
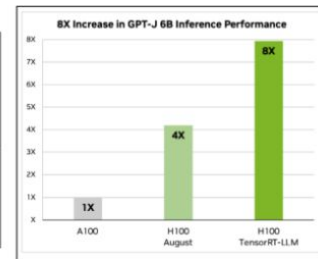
- Challenges:
    - Document parsing
    - Factuality and hallucinations
    - Splitting
    - Latency
    - …



- Hard to capture all information in text
- Very few "key points" of attention. More Focus on "General Imagery"



**DASHBOARD (LEFT-HAND DRIVE)**

- Some details can be captured in Text
- Has both "key points" of attention and "General Imagery"



8X Increase in GPT-J 6B Inference Performance

- Details can be perfectly captured in text.
- "Has key points of attention"