

# Hands on lab: Cómo crear aplicaciones de IA Generativa con LLMs

26 octubre 2023



# índice

- Ejemplo práctico implementación sistemas RAG (Retrieval Augmented Generation)

# 03 ejemplo práctico implementación sistemas RAG

# RAG

## Intro

La Retrieval-Augmented Generation se refiere a un enfoque de procesamiento de lenguaje natural (NLP) y aprendizaje automático que combina elementos de **modelos de recuperación** y **modelos de generación** para mejorar la calidad y relevancia del texto generado.

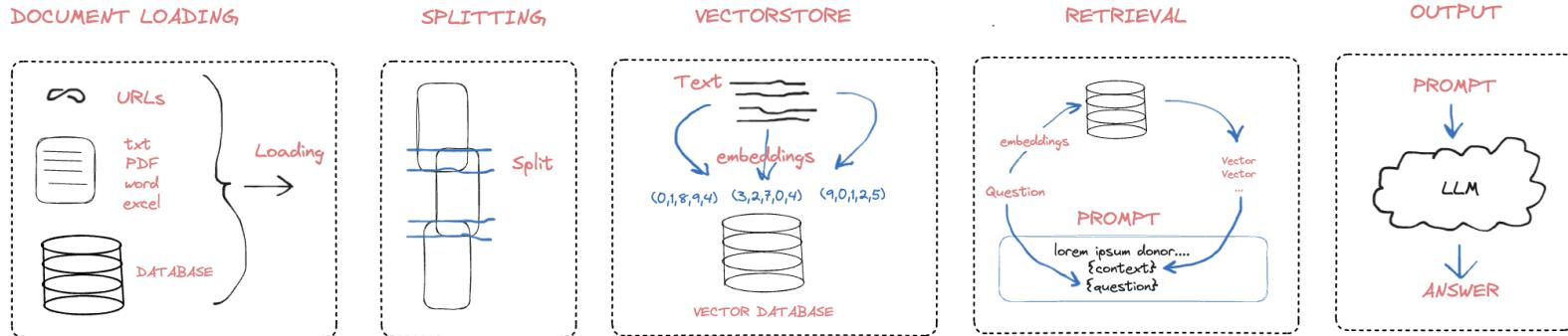
Los **modelos de recuperación** se centran en encontrar información relevante en un conjunto predefinido de documentos o fuentes de datos.

Los **modelos de generación** son capaces de generar texto parecido al humano desde cero o completar frases basadas en un contexto dado

Esta técnica se utiliza a menudo para tareas como generación de texto, respuesta a preguntas y resumen de contenido. Combina las fortalezas de los métodos basados en recuperación y los métodos basados en generación.

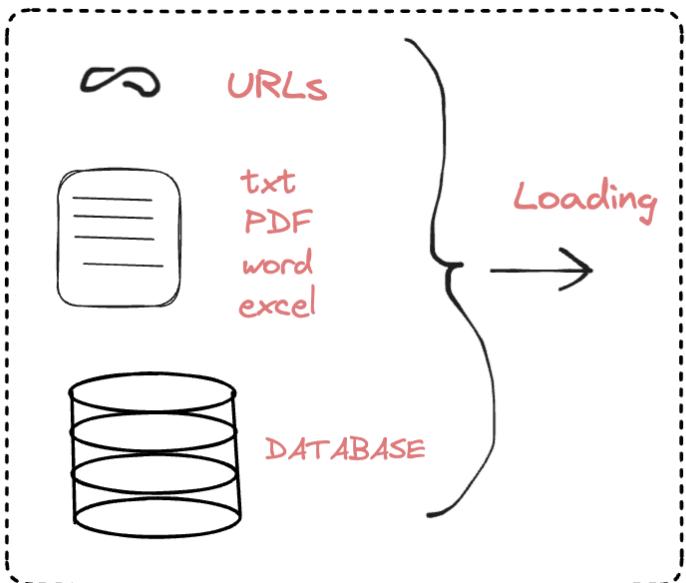
# RAG

## Steps



# RAG

## Document loading

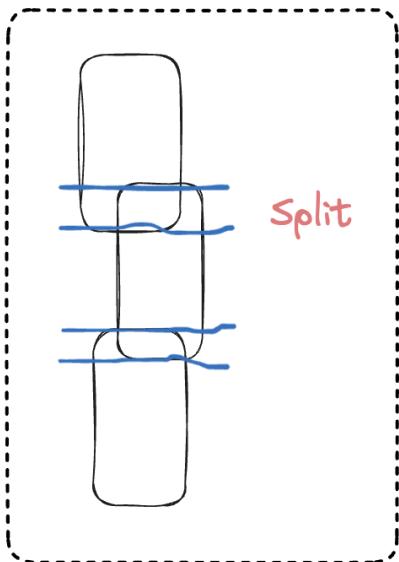


```
import { DirectoryLoader } from "langchain/document_loaders/fs/directory";
import { TextLoader } from "langchain/document_loaders/fs/text";
import { PDFLoader } from "langchain/document_loaders/fs/pdf";

const loader = new DirectoryLoader(
  'documents',
  {
    ".txt": (path) => new TextLoader(path),
    ".pdf": (path) => new PDFLoader(path)
  }
)
```

# RAG

## Splitting

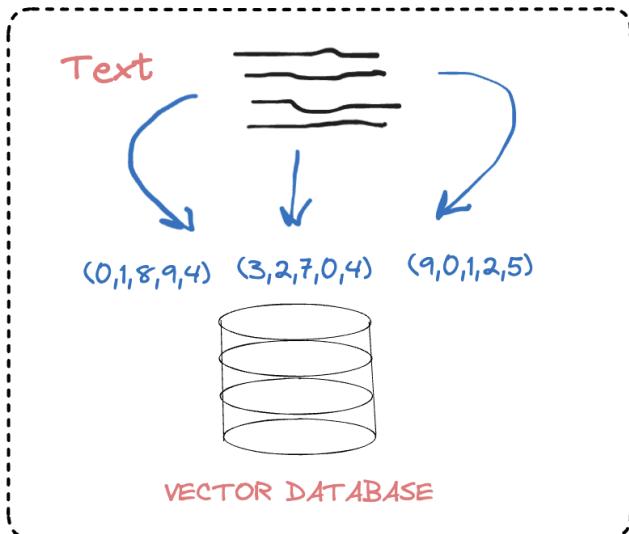


```
import { RecursiveCharacterTextSplitter } from "langchain/text_splitter";

const splitter = new RecursiveCharacterTextSplitter({
    chunkSize: 1000,
    chunkOverlap: 200,
});
```

# RAG

## Vectorstore



```
import { PineconeClient } from "@pinecone-database/pinecone";
import { PineconeStore } from "langchain/vectorstores/pinecone";
import { OpenAIEmbeddings } from "langchain/embeddings/openai";

await PineconeStore.fromDocuments(document, new OpenAIEmbeddings(), {
  pineconeIndex,
  // namespace: <namespace-name>,
})
```

# RAG

## Vectorstore

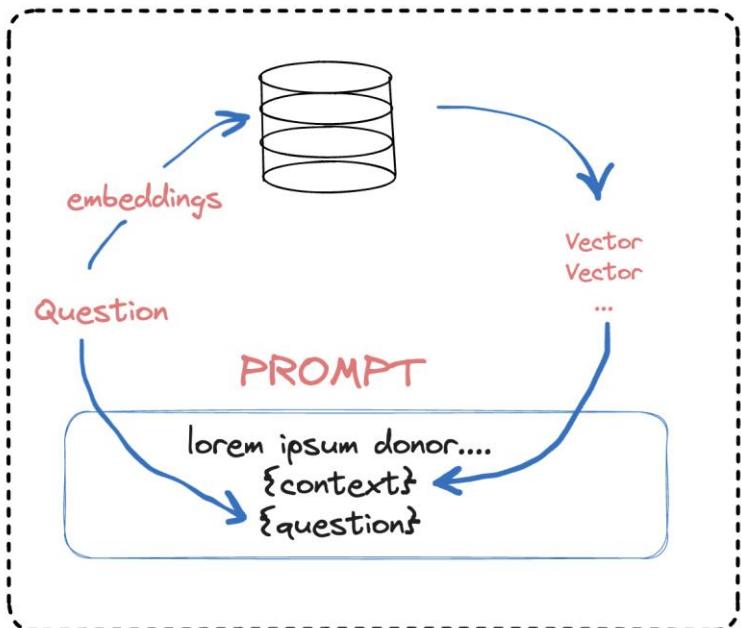
The screenshot shows the Pinecone search interface. On the left is a sidebar with navigation links: PROJECT (Search application, Indexes, API Keys, Members), Docs, Support Center, a user profile (pedro jimenez ca...), and Settings. At the bottom is an Upgrade button. The main area is titled "Search application Indexes" and has tabs for BROWSER and METRICS. It displays two search results:

ID	VALUES
d2ceca75-f7b7-4c5c-a1a5-8...	0.000747660233, -0.0141103...
80e7ae2c-1b63-46e6-9036...	-0.0220680479, -0.0115444...

Each result includes a SCORE (1 and 2 respectively) and METADATA fields. The METADATA for result 1 includes: loc.lines.from: 1, loc.lines.to: 63, loc.pageNumber: 22, pdf.info.CreationDate: "D:20231014191312+02'00)". There is a "Show 9 more" link at the bottom of each result.

# RAG

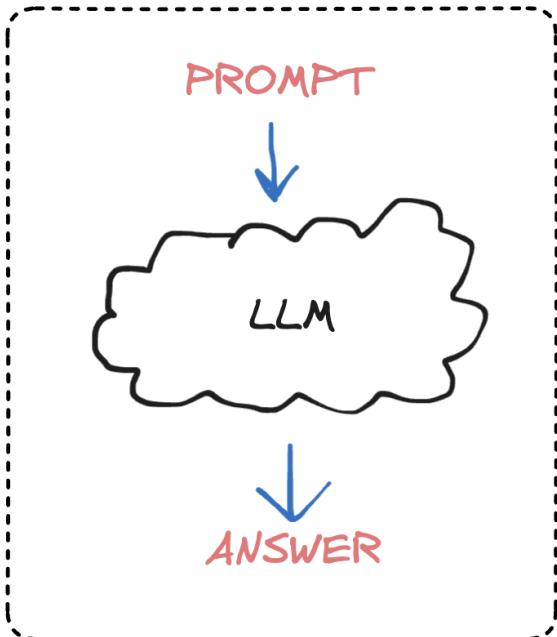
## Retrieval



```
const vectorStore = await PineconeStore.fromExistingIndex(  
  new OpenAIEmbeddings(),  
  {  
    pineconeIndex: index,  
    filter: {}  
  }  
);  
const retriever = vectorStore.asRetriever();  
const chain = await makeChain(retriever);
```

# RAG

## Output



```
const questionPrompt = ...buildQuestionPrompt();
const question = await questionPrompt.format({
  question: req.body.question,
})

const answer = await chain.call({ question, chat_history: [] });
```

# RAG

## Docs



**LangChain**

NodeJS

[Installation | !\[\]\(341b5bdc31177a6c7da7dc713da0d169\_img.jpg\) Langchain](#)

Python

[Introduction | !\[\]\(eaac180de418db4eae4b4cefebda75e8\_img.jpg\) Langchain](#)



**Pinecone**

[Pinecone](#)



**OpenAI**

[Introduction - OpenAI API](#)

# ¡Gracias!

diverger