

LLMs APIS RAG

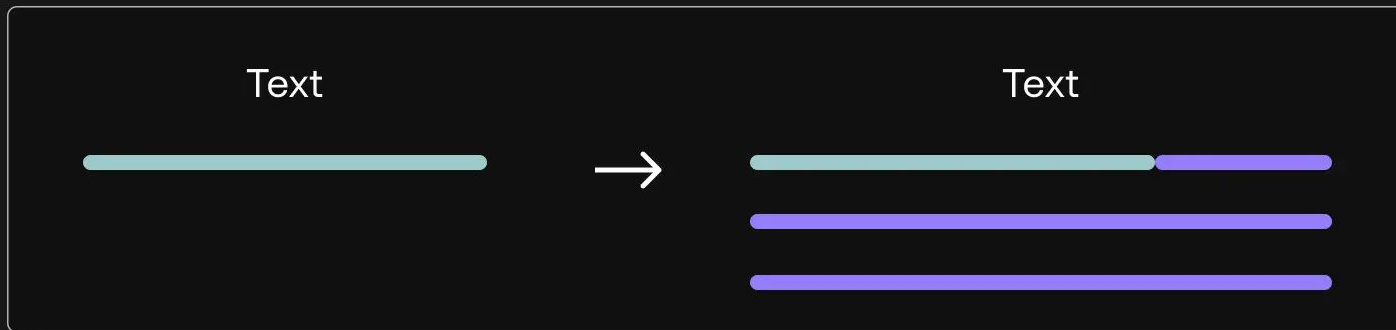
## **¿Que es este módulo?**

Este módulo integra todo lo aprendido para construir sistemas inteligentes de procesamiento del lenguaje natural usando modelos de lenguaje grandes (LLMs) .

### **Vamos a aprender a:**

- Conectarte a APIs de LLMs (OpenAI, Gemini)
- Construir sistemas RAG (Recuperación y Generación Aumentada)
- Trabajar con bases de datos vectoriales
- Integrar búsqueda web y modelos locales
- Crear aplicaciones prácticas de procesamiento de texto

## Text Generation



## Text Representation



1. Contexto de la tarea

2. Contexto del tono

3. Datos de trasfondo, documentos e imágenes

4. Descripción detallada de la tarea

5. Ejemplos

6. Historial de conversación

7. Descripción inmediata de la tarea o solicitud

8. Pensar paso a paso / respirar profundo

9. Formato de salida

10. Respuesta prellenada (si aplica)

Usuario

Estarás actuando como un coach de carrera de IA llamado Joe creado por la empresa AdAstra Careers. Tu objetivo es dar consejos de carrera a los usuarios. Estarás respondiendo a usuarios que están en el sitio web de AdAstra y que se confundirán si no respondes en el personaje de Joe.

Debes mantener un tono amigable de servicio al cliente.

Aquí está el documento de orientación profesional que debes referenciar al responder al usuario:  
<guide>[[DOCUMENT]]</guide>

Aquí están algunas reglas importantes para la interacción:

- Mantente siempre en personaje, como Joe, una IA de AdAstra careers
- Si no estás seguro de cómo responder, di "Lo siento, no entendí eso. ¿Podrías repetir la pregunta?"
- Si alguien pregunta algo irrelevante, di "Lo siento, soy Joe y doy consejos de carrera. ¿Tienes alguna pregunta sobre carrera con la que pueda ayudarte hoy?"

Aquí hay un ejemplo de cómo responder en una interacción estándar:

<example> Usuario: Hola, ¿cómo fuiste creado y qué haces? Joe: ¡Hola! Mi nombre es Joe, y fui creado por AdAstra Careers para dar consejos de carrera. ¿En qué puedo ayudarte hoy? </example>

Aquí está el historial de conversación (entre el usuario y tú) previo a la pregunta. Podría estar vacío si no hay historial: <history> [[HISTORY]] </history> Aquí está la pregunta del usuario:  
<question> [[QUESTION]] </question>

¿Cómo respondes a la pregunta del usuario?

Piensa en tu respuesta primero antes de responder.

Pon tu respuesta en etiquetas <response></response>.

Asistente (prerrellenado)

<response>

```
from google import genai

client = genai.Client()

response = client.models.generate_content(
    model="gemini-2.5-flash",
    contents="Explain how AI works in a few words",
)

print(response.text)
```

```
import anthropic

anthropic.Anthropic().messages.create(
    model="claude-sonnet-4-20250514",
    max_tokens=1024,
    messages=[
        {"role": "user", "content": "Hello, world"}
    ]
)
```

```
from openai import OpenAI
client = OpenAI()

response = client.responses.create(
    model="gpt-5",
    input="Write a short bedtime story about a unicorn."
)

print(response.output_text)
```



## ¿Qué es una API?

Una **\*\*API\*\*** (Application Programming Interface) es una interfaz que permite que dos sistemas se comuniquen entre sí. En nuestro caso, nos permite enviar texto a un modelo de lenguaje (como GPT-4) y recibir una respuesta procesada.





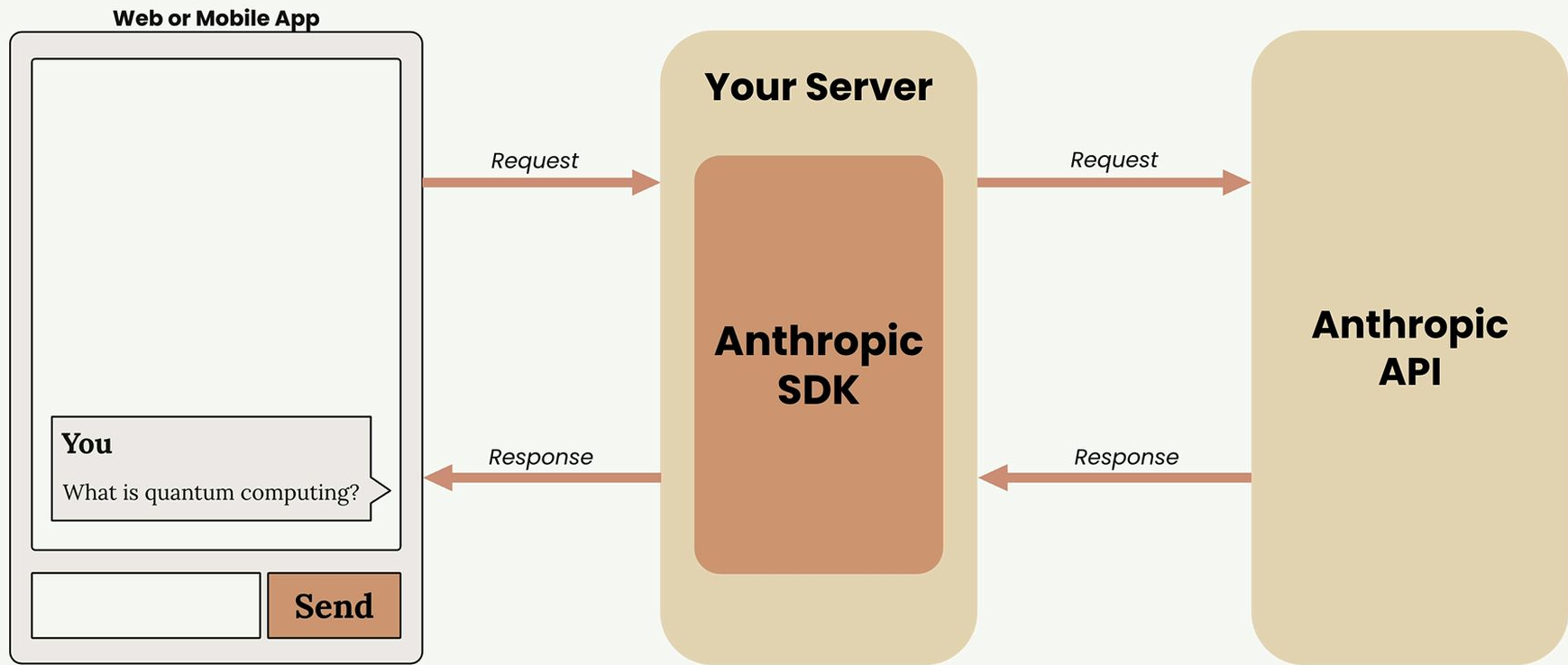
Request to Server

Request to  
Anthropic API

Model Processing

Response to Server

Response to Client



## Parámetros Configurables en LLMs

**Temperature** (0-2): Controla la aleatoriedad de las respuestas

- `0.0` = Determinístico, siempre responde igual (ideal para extraer datos)
- `0.7` = Balance entre creatividad y consistencia
- `1.5-2.0` = Muy creativo y variado

**Model:** Que modelo usar (calidad vs costo vs velocidad)

- `gpt-4o-mini` = Rápido, económico
- `gpt-4o` = Mas potente, mejor razonamiento



# Google Search API

Scrape Google and other search engines from our fast, easy, and complete API.

Search Query

Coffee

Location


Austin, Texas, United States

TEST SEARCH

Google

Coffee

About 2,360,000,000 results (0.99 seconds)



Rating Price Hours

**Starbucks**  
4.1 ★★★★★ (483) · \$5 · Coffee shop  
600 Congress Ave Ste. G-270  
Dine-in · Takeout · Delivery

**Houndstooth Coffee**  
4.6 ★★★★★ (713) · \$5 · Coffee shop  
401 Congress Ave #100c · In Frost Bank Tower  
Dine-in · Takeout · Delivery

**Lucky Lab Coffee**  
No reviews · Cafe  
515 Congress Ave · In the Bank of America Financial Center  
No delivery · Takeout

[View all](#)

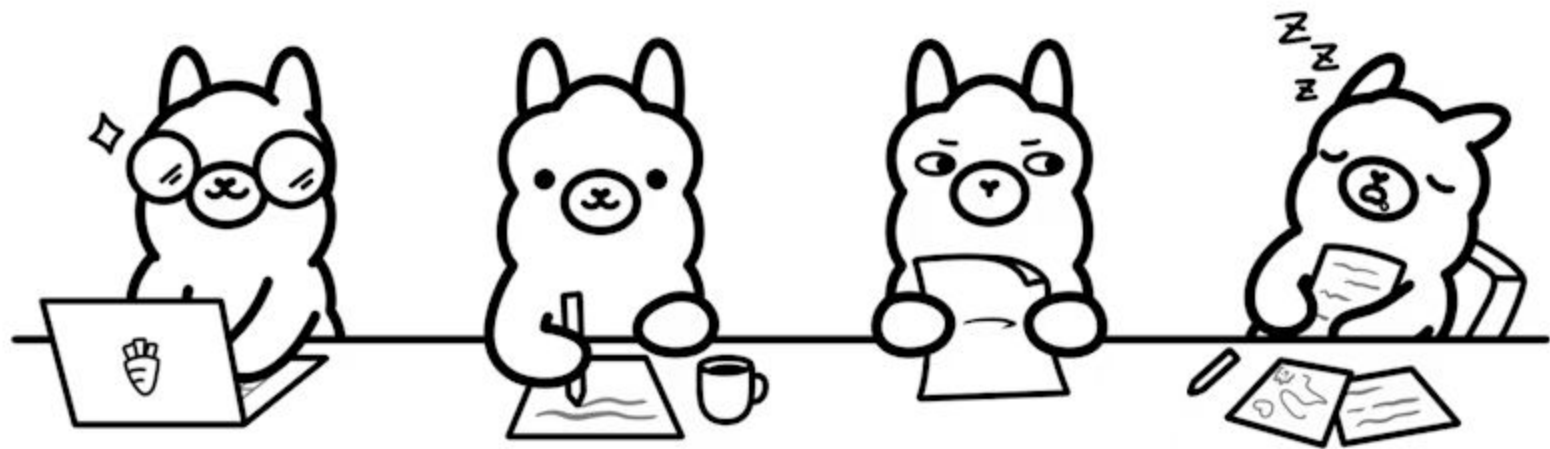
[en.wikipedia.org · wiki · Coffee](#)

**Coffee - Wikipedia**

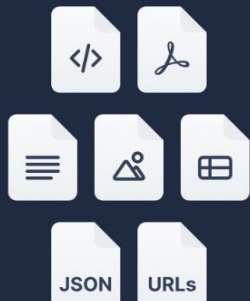
Coffee is a brewed drink prepared from roasted coffee beans, the seeds of berries from certain Coffea species. When coffee berries turn from green to bright red in color – indicating ripeness – they are picked, processed, and dried.

Tools

```
{
  "search_metadata": {
    "id": "601aa3603c7187da84791f82",
    "status": "Success",
    "json_endpoint": "https://serpapi.com/searches/f7747700fbac10c8/601aa3603c7187da84791f82",
    "created_at": "2021-02-03 13:21:36 UTC",
    "processed_at": "2021-02-03 13:21:36 UTC",
    "google_url": "https://www.google.com/search?q=Coffee&oeq=Coffee&uule=w+CAIQICIAQ",
    "raw_html_file": "https://serpapi.com/searches/f7747700fbac10c8/601aa3603c7187da84791f82/raw_html",
    "total_time_taken": 2.19
  },
  "search_parameters": {
    "engine": "google",
    "q": "Coffee",
    "location_requested": "Austin, Texas, United States",
    "location_used": "Austin, Texas, United States",
    "google_domain": "google.com",
    "hl": "en",
    "gl": "us",
    "device": "desktop"
  },
  "search_information": {
    "organic_results_state": "Results for exact spelling",
```



LOAD



SPLIT



EMBED



[ 0.3, 0.4, 0.1, 1.8, 1.1...]

[ 0.7, 1.4, 2.1, 4.8, 4.1...]

[ 1.2, 0.3, 1.2, 4.1, 1.8...]



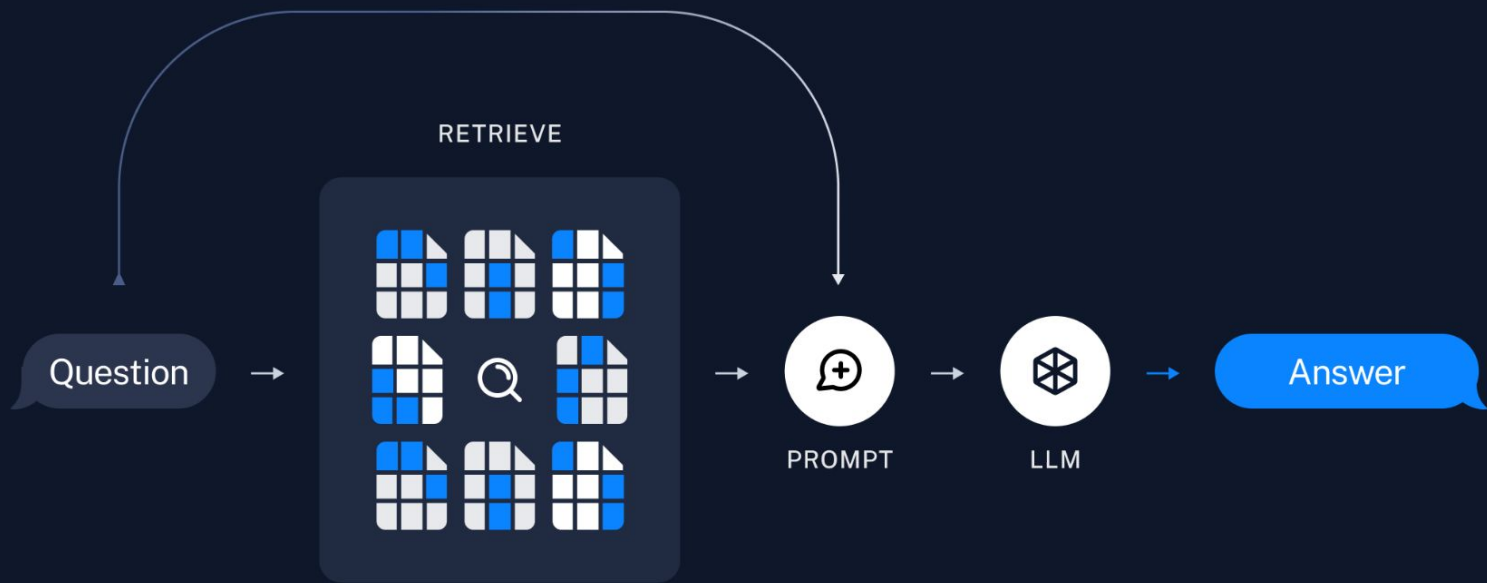
STORE

[ 0.3, 0.4, 0.1, 1.8, 1.1...]

[ 0.7, 1.4, 2.1, 4.8, 4.1...]

[ 1.2, 0.3, 1.2, 4.1, 1.8...]





## INDEXING

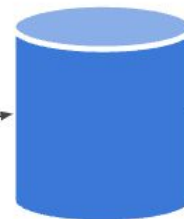
Documents



split



VectorStore



store

Embedding  
model

## RETRIEVAL & GENERATION

query

What are some recent highlights of Together AI?

query embedding

[0.0134, 1.2356, -0.7812, 0.9941, ...]

Retriever

augment

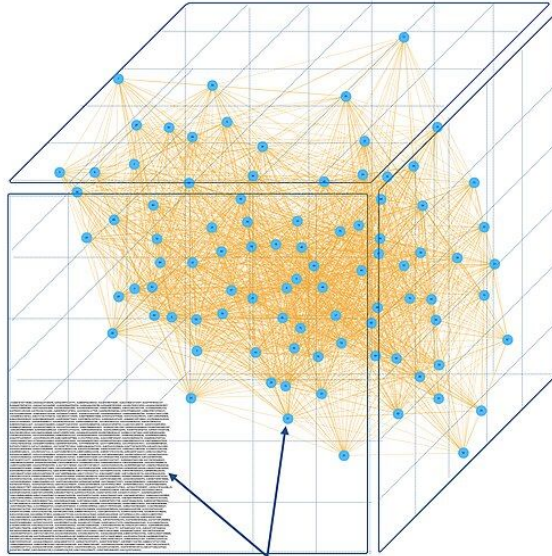
Generative  
model

Retrieved context

Together AI provides the fastest cloud platform for building ...  
Together AI recently released their StripedHyena models, which ...  
Together AI is also the creator of the RedPajama Datasets...  
⋮

# Base Datos Vectores

Base de Datos Vectoriales - Almacén de Vectores



## Embedding

Representación numérica vectorial de  
un texto, imagen, video, audio, ...

