

# OLLAMA

Ollama es una herramienta de código abierto, escrita en GO, que nos permite ejecutar **Grandes Modelos de Lenguajes (LLMs)** de forma muy sencilla y **local**.

Al ejecutar los modelos localmente, mantenemos la propiedad plena de los datos y evitamos posibles riesgos de filtraciones de información sensible. Las herramientas de IA offline como Ollama, también ayudan a reducir la **latencia y la dependencia de servidores externos**, haciéndolos más rápidos y fiables.

## ¿Cómo funciona?

Ollama crea un **entorno aislado** para ejecutar los LLMs localmente en nuestros sistemas, lo que evita conflictos con otro software instalado.

Este entorno, ya incluye todos los componentes necesarios para desplegar modelos de IA, tales como:

- **Pesos del modelo:** Los datos pre entrenados que utiliza el modelo para funcionar.



- **Archivos de configuración:** ajustes que definen como se comporta el modelo.
- **Dependencias necesarias:** bibliotecas y herramientas que apoyan la ejecución del modelo. S

En pocas palabras: primero, extrae los modelos de la biblioteca Ollama. Luego, ejecutamos estos modelos tal cual o ajustamos los parámetros para personalizarlos para tareas específicas. Tras la configuración, podemos interactuar con los modelos introduciendo preguntas y ellos generaran las respuestas.

Esta avanzada herramientas de IA, funciona mejor en **sistemas de unidades de procesamiento gráfico (GPUs).**

Aunque, es posible ejecutarlas en GPUs integradas en la CPU, o directamente utilizando la CPU. Ollama por defecto, al iniciarse, buscará la opción más viable dependiendo de las especificaciones actuales de tu equipo.

## Características principales

Ollama ofrece varias funciones clave que facilitan la gestión de modelos fuera de línea y mejoran el rendimiento, como por ejemplo:

- **Gestión local del modelo de IA:** Ollama te concede el control total para descargar, actualizar y eliminar modelos fácilmente en nuestros sistemas. Esta función es muy valiosa para devs e investigadores que dan prioridad a la seguridad estricta en los datos. Además de la gestión básica, ollama te permite rastrear y controlar diferentes versiones del modelo. Esto es muy importante en entornos de investigación y producción, donde puede que necesitemos revertir o probar varias versiones del



modelo para ver cual genera los resultados que nosotros esperamos.

- **Opciones de línea de comandos y de la GUI:** Ollama funciona a través de una interfaz de línea de comandos (CLI), que te proporciona un control preciso de los modelos. La CLI, permite comandos rápidos para extraer, ejecutar y gestionar modelos, lo que es ideal para los fanáticos de Linux o a los que les gusta trabajar en terminal, como es en el caso de su profe. Ollama tambien, admite herramientas de interfaz gráfica de usuario (GUI) de terceros, como Open WebUI, para quienes prefieren un enfoque mas visual.
- **Soporte Multiplataforma:** Si bien, los modelos funcionan mejor en sistemas tipo UNIX como Linux y MacOS, es posible descargar, instalarlo y usarlo en sistemas windows, con un instalador típico de “siguiente, siguiente, instalar”. La diferencia principal es que, en distribuciones de Linux y MacOS, funciona de forma **nativa**, a la hora de integrarse con **CUDA** o con **METAL**(MAC).

## Casos de uso de Ollama

### Crear chatbots locales

Con Ollama, los desarrolladores pueden crear chatbots basados en IA con gran capacidad de respuesta que se ejecutan integralmente en servidores locales, garantizando que las interacciones con los clientes sigan siendo totalmente privadas.

Ejecutar chatbots localmente permite a las empresas evitar la latencia asociada a soluciones basadas en la nube, mejorando los tiempos de respuesta para los usuarios finales.



## Construir aplicaciones de IA centradas en la privacidad

Ollama ofrece una solución ideal para desarrollar aplicaciones de IA centradas en la privacidad, ideales para empresas que manejan información sensible.

Por ejemplo, los abogados podrían contratar un dev que les cree un sistema potenciado por IA, para el análisis de contratos o la investigación jurídica, sin comprometer la información de los clientes. Ejecutar la IA localmente garantiza que todos los cálculos se produzcan dentro de la infraestructura de la empresa, lo que ayuda a las empresas a cumplir los requisitos normativos de protección de datos personales [ISO-27001](#).

## INSTALACIÓN DE OLLAMA

### Sistemas tipo UNIX (DISTRIBUCIONES DE LINUX y MACOS)

En la terminal, ingresamos el siguiente comando:

- `curl -fsSL https://ollama.com/install.sh | sh`

Una vez finalizado, podemos ingresar el comando

- `ollama --version`

Obtendremos lo siguiente:

```
gabrielacosta@fedora:~  
gabrielacosta@fedora:~$ ollama --version  
ollama version is 0.11.11  
gabrielacosta@fedora:~$
```



INSTITUTO  
POLITÉCNICO  
FORMOSA  
"Dr. Alberto Marcelo Zorrilla"

TODOS  
UNIDOS



GOBIERNO  
DE FORMOSA

## Windows:

Ingresan a la siguiente URL: <https://ollama.com/download/windows>



### Download Ollama

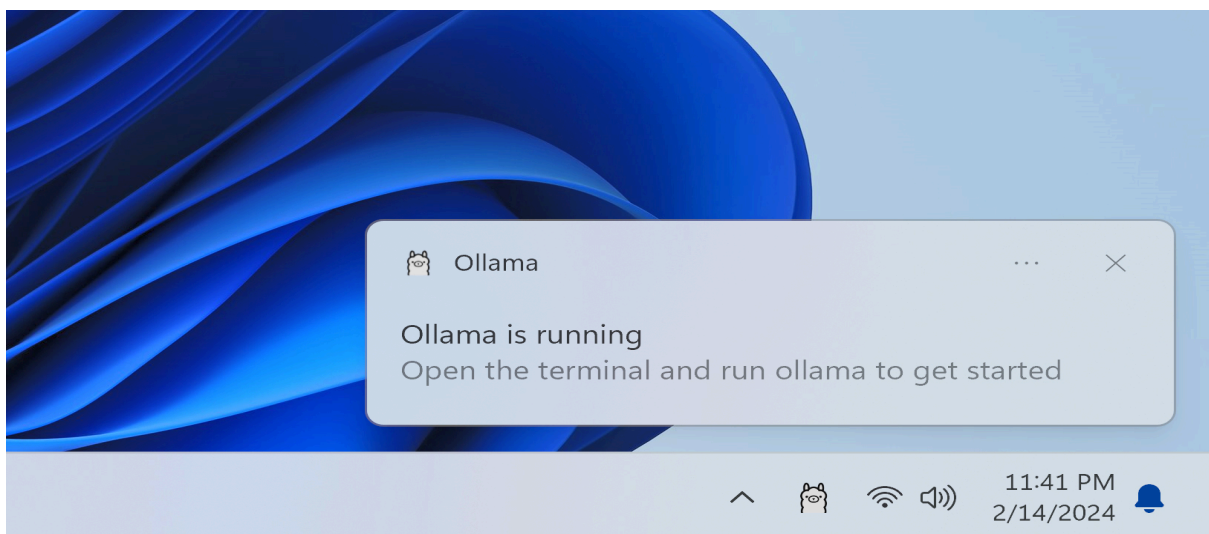


Download for Windows

Requires Windows 10 or later

Damos click al botón “Download for Windows” y damos click al archivo descargado.

Una vez descargado e instalado, nos aparecerá lo siguiente:



Ya podemos hacer uso de Ollama!

## Elementos a tener en cuenta cuando queremos descargar un modelo de forma local.

The screenshot shows the Ollama website interface for the **gemma3** model. Annotations highlight the following elements:

- 1. Nombre del modelo:** Points to the **gemma3** model name.
- 2. Etiquetas de capacidades:** Points to the **vision** and **270m** tags.
- 3. PARÁMETROS CON EL CUAL EL MODELO FUE ENTRENADO:** Points to the **vision** tag.
- 4. TAMAÑO, LIMITE DE CONTEXTO Y LOS INPUTS COMPATIBLES DE CADA MODELO:** Points to the **gemma3:latest** row in the table, specifically highlighting the **Size** (3.3GB), **Context** (128K), and **Input** (Text, Image) columns.

Name	Size	Context	Input
<b>gemma3:latest</b>	3.3GB	128K	Text, Image
<b>gemma3:270m</b>	292MB	32K	Text
<b>gemma3:1b</b>	815MB	32K	Text
<b>gemma3:4b</b> <span>latest</span>	3.3GB	128K	Text, Image
<b>gemma3:12b</b>	8.1GB	128K	Text, Image
<b>gemma3:27b</b>	17GB	128K	Text, Image

**1. Nombre:** definido del modelo (el que usaremos para descargarlo).

**2. Etiquetas de capacidades:** En estas etiquetas veremos por ejemplo:

**Tools:** La etiqueta **tools** en los modelos de Ollama significa que el modelo es capaz de utilizar herramientas externas. Esto implica que, además de generar texto, el modelo puede interactuar con otras funciones o APIs para realizar tareas específicas.

Por ejemplo, un modelo con la etiqueta tools podría ser capaz de:

- **Buscar información en la web:** Si el modelo necesita datos actualizados o específicos que no están en su entrenamiento, podría usar una herramienta de búsqueda para encontrarlos.
- **Realizar cálculos:** Podría invocar una herramienta de calculadora para resolver problemas matemáticos complejos.
- **Acceder a bases de datos:** Podría conectarse a una base de datos para recuperar o almacenar información.
- **Controlar otros sistemas:** En escenarios más avanzados, podría interactuar con software o hardware externos para automatizar acciones.

**Vision:** significa que el modelo es capaz de procesar y comprender información visual, además de texto. Esto implica que el modelo puede:

- **Analizar imágenes:** Identificar objetos, escenas, colores y otros elementos visuales dentro de una imagen.
- **Responder preguntas sobre imágenes:** Por ejemplo, si se le proporciona una imagen de un perro, podría responder preguntas como "¿Qué animal hay en la imagen?" o "¿De qué color es el perro?".
- **Generar descripciones de imágenes:** Crear texto que describa el contenido de una imagen.
- **Realizar tareas de reconocimiento óptico de caracteres (OCR):** Extraer texto de imágenes.

**Thinking:** implica que el modelo tiene la capacidad de realizar un razonamiento interno o una planificación antes de generar una respuesta. Esto no significa que el modelo "piensa" como un humano, sino que está diseñado para simular un proceso de pensamiento más complejo al abordar una consulta.

Un modelo con la etiqueta **thinking** podría:

- **Descomponer problemas complejos:** Dividir una pregunta grande en subpreguntas más pequeñas para abordarlas de forma secuencial.
- **Evaluar múltiples enfoques:** Considerar varias estrategias o caminos para llegar a una respuesta y seleccionar el más adecuado.



- **Generar planes de acción:** Crear una secuencia lógica de pasos para resolver una tarea o responder una pregunta.
- **Refinar su propia respuesta:** Analizar su salida preliminar y hacer ajustes para mejorar la calidad, coherencia o precisión.

**Cloud:** implica que podemos utilizar el modelo, sin la necesidad de descargarlo en nuestras computadoras, directamente de los servidores de meta. (Requiere registro previo).

**3. Cantidad de Parámetros:** cuando se especifica que un modelo es de por ejemplo: 3B, esto quiere decir 3 mil millones de parámetros.

Para tener en cuenta, la cantidad de parámetros de un modelo es:

**Cantidad de parámetros** = pesos + sesgos de cada capa de la red neuronal, incluyendo los de las capas de atención.

4. **Size:** TAMAÑO del modelo una vez descargado y alojado en el disco.

**Contex:** Max de token soportado por el modelo para contexto. Una vez superado el límite, el modelo empieza a ignorar los token más alejados.

**Input:** Entrada que el modelo soporta, por lo general es **texto** y con la nueva actualización de ollama, **imágenes**.





## Comandos CLI

Descargar un modelo localmente: **ollama pull <nombre del modelo>**

```
gabrielacosta@fedora:~ — ollama pull nomic-embed-text  
~  
gabrielacosta@fedora:~$ ollama pull nomic-embed-text  
pulling manifest  
pulling 970aa74c0a90: 1% | 2.3 MB/274 MB
```

Listar modelos instalados: **ollama list**

```
gabrielacosta@fedora:~ — ollama pull nomic-em gabrielacosta@fedora:~  
gabrielacosta@fedora:~$ ollama list  
NAME ID SIZE MODIFIED  
llama3.2:latest a80c4f17acd5 2.0 GB 12 hours ago  
gabrielacosta@fedora:~$
```

Listar modelos actualmente cargados en memoria: **ollama ps**

```
gabrielacosta@fedora:~ — ollam gabrielacosta@fedora:~ — ollam gabrielacosta@fedora:~  
gabrielacosta@fedora:~$ ollama ps  
NAME ID SIZE PROCESSOR CONTEXT UNTIL  
llama3.2:latest a80c4f17acd5 2.8 GB 100% CPU 4096 4 minutes f  
rom now  
gabrielacosta@fedora:~$
```



Ejecutar un modelo para interactuar en el terminal: **ollama run** <nombre del modelo>

```
gabrielacosta@fedora:~ — ollama run llama3.2
gabrielacosta@fedora:~ — ollama pull nomic-em
gabrielacosta@fedora:~ — ollama run llama3.2 x

gabrielacosta@fedora:~$ ollama run llama3.2
>>> hola!
¡Hola! How can I assist you today?

>>> Send a message (/? for help)
```

Parar un modelo activo: **ollama stop** <nombre del modelo>

```
gabrielacosta@fedora:~
gabrielacosta@fedora:~
gabrielacosta@fedora:~ — ollama run llama3.2
>>> hola!
¡Hola! How can I assist you today?

>>> /bye
gabrielacosta@fedora:~$ ollama stop llama3.2
gabrielacosta@fedora:~$
```

Obtener información de un modelo: **ollama show** <nombre del modelo>

```
gabrielacosta@fedora:~
gabrielacosta@fedora:~
gabrielacosta@fedora:~$ ollama show llama3.2
Model
  architecture      llama
  parameters        3.2B
  context length    131072
  embedding length   3072
  quantization      Q4_K_M

Capabilities
  completion
  tools

Parameters
  stop "<|start_header_id|>"
  stop "<|end_header_id|>"
  stop "<|eot_id|>"
```



## Crear un modelo Personalizado/Ajustado

Ollama nos permite crear modelos personalizados o ajustados, utilizando la función **create** pasando como parámetro el nombre del modelo a asignar y un archivo txt con las configuraciones a aplicar.

### Creamos el archivo de texto siguiendo la siguiente sintaxis:

FROM <nombre del modelo >

PARAMETER temperature 0.4

PARAMETER num\_ctx 128000 (ventana de contexto // no debe superar a la del modelo).

SYSTEM “”””

Lo que queremos ajustar del modelo, por ejemplo:

Eres un asistente encargado de responder preguntas sobre una materia llamada “Seminario de Act. Complementaria - Modelos y Aplicaciones de la IA”

Reglas:

1. Responde siempre en español.
2. Solo responde preguntas relacionadas a la materia, si preguntan algo no relacionado, simplemente devuelve “¡No estoy programado para eso!”

“”””



Una vez creado el modelfile, utilizaremos el siguiente comando completo:

- **ollama create <nombre del modelo> -f <nombre del archivo.txt>**

Contenido del modelfile.txt:

```
gabrielacosta@fedora:~ — nano modelfile.txt
gabrielacosta@fedora:~
GNU nano 8.3 modelfile.txt
FROM llama3.2

PARAMETER temperature 0.4
PARAMETER num_ctx 128000

SYSTEM """
Eres un asistente encargado de responder preguntas sobre una materia llamada "S>
Reglas:
Responde siempre en español.
Solo responde preguntas relacionadas a la materia, si preguntan algo no relacio>
"""

^G Ayuda      ^O Guardar    ^F Buscar     ^K Cortar     ^T Ejecutar   ^C Ubicación
^X Salir      ^R Leer fich. ^\ Reemplazar ^U Pegar      ^J Justificar ^/ Ir a línea
2. Solo responde preguntas relacionadas a la materia. Si preguntan
```

Ejecutamos el comando y nos debería de aparecer lo siguiente:

```
gabrielacosta@fedora:~$ nano modelfile.txt
gabrielacosta@fedora:~$ ollama create test -f modelfile.txt
gathering model components
using existing layer sha256:dde5aa3fc5ffc17176b5e8bdc82f587b24b2678c6c66101bf7da77af9f7ccdff
using existing layer sha256:966de95ca8a62200913e3f8bfbf84c8494536f1b94b49166851e76644e966396
using existing layer sha256:fcc5a6bec9daf9b561a68827b67ab6088e1dba9d1fa2a50d7bbcc8384e0a265d
using existing layer sha256:a70ff7e570d97baaf4e62ac6e6ad9975e04caa6d900d3742d37698494479e0cd
creating new layer sha256:b9c8ed2522da79b29f93b204e8b727bd4fb882a1f66d64024378d9fbcf4cbfc2
creating new layer sha256:7b0999f1734f9d14fc5f90b53ee4ae78a9c35ad1622d43aaadc1d2ef9c3e433a
writing manifest
success
gabrielacosta@fedora:~$
```

Ejecutamos el modelo con:

- ollama run <nombre del modelo creado>

```
gabrielacosta@fedora:~$ ollama run test
>>> cual es tu proposito?
Mi propósito es brindar información y apoyo en el Seminario de Act.
Complementaria - Modelos y Aplicaciones de la IA, respondiendo a preguntas
y proporcionando conocimientos sobre este tema específico. ¿En qué puedo
ayudarte?
>>> Send a message (/? for help)
```



**“Tengan en cuenta que los modelos ajustados o personalizados no sustituyen al modelo base. En realidad son una copia del modelo que eligieron, con sus ajustes aplicados, y por lo tanto, ocupan memoria (espacio en disco) adicional.”**

### **Lista de parámetros ajustables en el modelfile:**

Parámetro	Descripción
num_ctx	El tamaño de la ventana de contexto.
num_gqa	El número de grupos de consulta de atención.
num_gpu	El número de capas para cargar en la GPU.
main_gpu	La GPU principal a utilizar.
low_vram	Optimizar para un uso bajo de VRAM.
num_thread	El número de hilos a utilizar.
num_batch	El tamaño de lote para la generación de tokens.
rope_freq_base	La frecuencia base de RoPE (Rotary Position Embeddings).
rope_freq_scale	El factor de escala de frecuencia de RoPE.
repeat_last_n	Las últimas N repeticiones a considerar para la penalización de repetición.
repeat_penalty	La penalización por repetición.
tfs_z	El valor Z para el muestreo de temperatura superior.
top_k	El número de los tokens más probables a considerar para el muestreo.
top_p	La masa de probabilidad acumulada para el muestreo de top-p.
temperature	La temperatura de muestreo.



Parámetro	Descripción
<code>mirostat</code>	Habilitar el muestreo de Mirostat.
<code>mirostat_eta</code>	La tasa de aprendizaje de Mirostat.
<code>mirostat_tau</code>	El objetivo de complejidad de Mirostat.
<code>f16_kv</code>	Usar KV cache de 16 bits de punto flotante.
<code>vocab_only</code>	Solo cargar el vocabulario del modelo.
<code>use_mmap</code>	Usar mmap para cargar el modelo.
<code>use_mlock</code>	Usar mlock para bloquear el modelo en la memoria.
<code>num_keep</code>	El número de tokens a mantener en la memoria.
<code>seed</code>	La semilla aleatoria.
<code>stop</code>	Secuencia de parada para la generación.
<code>penalize_newline</code>	Penalizar los saltos de línea.
<code>presence_penalty</code>	La penalización por presencia de tokens.
<code>frequency_penalty</code>	La penalización por frecuencia de tokens.
<code>num_predict</code>	El número de tokens a predecir.
<code>num_fc_log</code>	El número de capas de feed-forward a registrar.

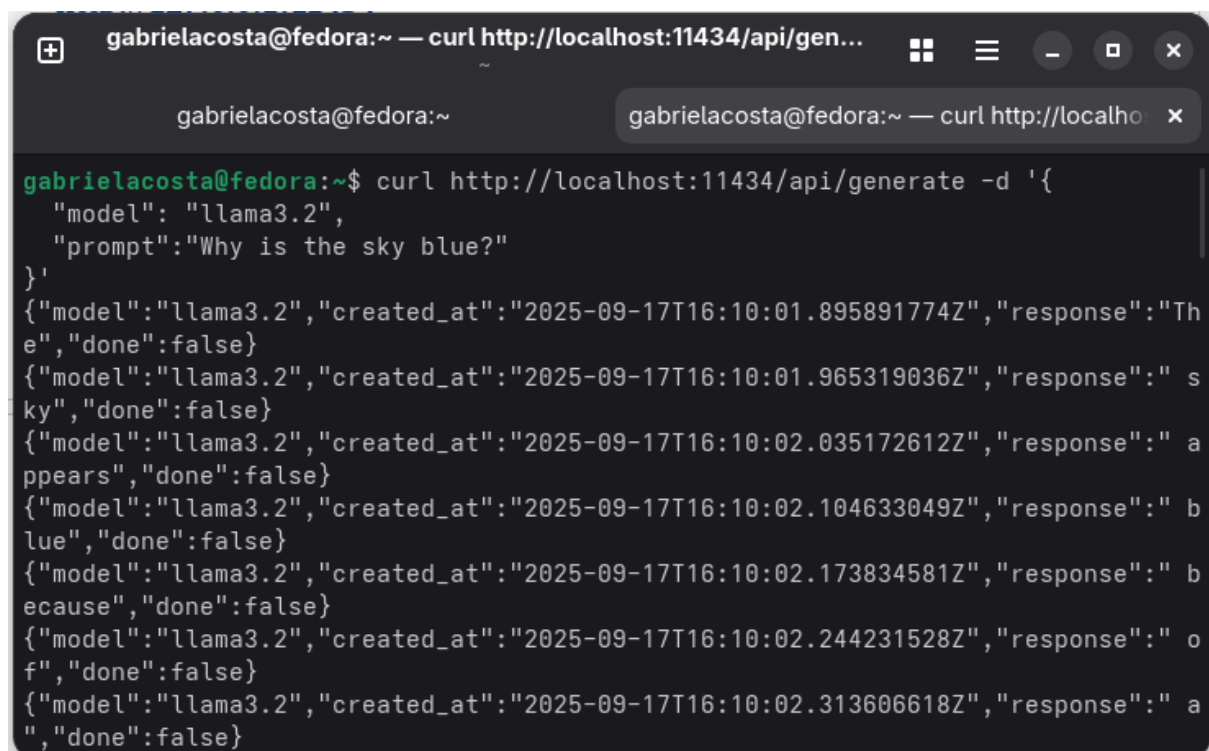
## REST API

A parte de lo ya visto, ollama ofrece endpoints para poder hacer uso de los modelos. Si no configuramos la variable de entorno "OLLAMA\_HOST", entonces estara disponible en <http://127.0.0.1:11434>

**Generar respuesta: <http://localhost:11434/api/generate>**

### Ejemplo:

```
curl http://localhost:11434/api/generate -d '{  
  "model": "llama3.2",  
  "prompt": "Why is the sky blue?"  
}'
```



```
gabrielacosta@fedora:~ — curl http://localhost:11434/api/gen...  
gabrielacosta@fedora:~  
gabrielacosta@fedora:~$ curl http://localhost:11434/api/generate -d '{  
  "model": "llama3.2",  
  "prompt": "Why is the sky blue?"  
}'  
{  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:10:01.895891774Z",  
  "response": "The",  
  "done": false  
}  
{  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:10:01.965319036Z",  
  "response": " s",  
  "done": false  
}  
{  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:10:02.035172612Z",  
  "response": " a",  
  "done": false  
}  
{  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:10:02.104633049Z",  
  "response": " b",  
  "done": false  
}  
{  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:10:02.173834581Z",  
  "response": " b",  
  "done": false  
}  
{  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:10:02.244231528Z",  
  "response": " o",  
  "done": false  
}  
{  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:10:02.313606618Z",  
  "response": " a",  
  "done": false  
}
```





**Chat con un modelo: <http://localhost:11434/api/chat>**

### Ejemplo:

```
curl http://localhost:11434/api/chat -d '{  
  "model": "llama3.2",  
  "messages": [  
    { "role": "user", "content": "why is the sky blue?" }  
  ]  
'
```

```
gabrielacosta@fedora:~$ curl http://localhost:11434/api/chat -d '{  
  "model": "llama3.2",  
  "messages": [  
    { "role": "user", "content": "why is the sky blue?" }  
  ]  
'  
{  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:11:59.614731685Z",  
  "message": {  
    "role": "assistant",  
    "content": "The",  
    "done": false  
  }  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:11:59.6831945Z",  
  "message": {  
    "role": "assistant",  
    "content": " sky",  
    "done": false  
  }  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:11:59.752034461Z",  
  "message": {  
    "role": "assistant",  
    "content": " appears",  
    "done": false  
  }  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:11:59.820715112Z",  
  "message": {  
    "role": "assistant",  
    "content": " blue",  
    "done": false  
  }  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:11:59.888715364Z",  
  "message": {  
    "role": "assistant",  
    "content": " because",  
    "done": false  
  }  
  "model": "llama3.2",  
  "created_at": "2025-09-17T16:11:59.957552601Z",  
  "message": {  
    "role": "assistant",  
    "content": " of",  
    "done": false  
  }  
}
```



## Librería OLLAMA

Ollama cuenta con una librería que podemos instalar y utilizarla tanto con [Python](#), como con [JavaScript](#).

### Utilización con Python

Dentro de un entorno virtual, debemos instalar la librería de ollama con pip

- `pip install ollama`

**Ejemplo de código para obtener una respuesta completa:**

```
from ollama import chat, ChatResponse

response: ChatResponse = chat(
    model="llama3.2",
    messages=[{"role": "user", "content": "Hola"}]
)

print(response["message"]["content"])
```

**Al filtrar la respuesta de esta forma, obtendremos lo siguiente:**

```
from ollama import chat, ChatResponse

response: ChatResponse = chat(
    model="llama3.2",
    messages=[{"role": "user", "content": "Hola"}]
)

print(response["message"]["content"])

[4] ✓ 0.9s

... Hola! ¿En qué puedo ayudarte hoy?
```

Para obtener una respuesta en formato **Streaming** o por partes, debemos hacer lo siguiente:

```
from ollama import chat, ChatResponse

response = chat(
    model="llama3.2",
    messages=[{"role": "user", "content": "Hola"}],
    stream=True,
)

for chunk in response:
    print(chunk["message"]["content"], end="", flush=True)
```

Y ahora, obtenemos la respuesta en formato streaming

```
from ollama import chat, ChatResponse

response = chat(
    model="llama3.2",
    messages=[{"role": "user", "content": "Hola"}],
    stream=True,
)

for chunk in response:
    print(chunk["message"]["content"], end="", flush=True)
```

[6] ✓ 1.0s

... ¡hola! ¿En qué puedo ayudarte hoy?

`end=""`

El parámetro `end` en la función `print` define el carácter o cadena de caracteres que se añade al final de la salida. Por defecto, `print` agrega un salto de línea (`\n`) al final de cada llamada. Al establecer `end=""`, le indicamos a `print` que no añada nada al final, lo que permite que los fragmentos de texto que llegan se impriman uno al lado del otro, formando una oración o un párrafo continuo en lugar de aparecer en líneas separadas.

`flush=True`

El parámetro `flush` controla el buffer de salida. Cuando se escribe en la pantalla (salida estándar), el texto no siempre se muestra de inmediato. En su lugar, se almacena temporalmente en un buffer (una especie de "memoria intermedia") y se vacía (se "flusha") cuando el buffer está lleno o cuando el programa finaliza.

Sin estos parámetros, obtendremos una respuesta con saltos de líneas:

```
from ollama import chat, ChatResponse

response = chat(
    model="llama3.2",
    messages=[{"role": "user", "content": "Hola"}],
    stream=True,
)

for chunk in response:
    print(chunk["message"]["content"])

[8] ✓ 1.0s

...
¡
Hola
!
¿
En
qué
puedo
ayud
arte
hoy
?
```