

Ollama

Ollama is a powerful tool that allows you to run large language models (LLMs) directly on your computer. This means you can interact with these complex Al models without relying on cloud services or internet access. **Think like Docker for LLMs.**



Use-cases of model creators

- Content Creation and Communication
 - Text Generation: These models can be used to generate creative text formats such as poems, scripts, code, marketing copy, and email drafts.
 - Chatbots and Conversational AI: Power conversational interfaces for customer service, virtual assistants, or interactive applications.
 - Text Summarization: Generate concise summaries of a text corpus, research papers, or reports.
 - Automation: Efficient and streamlined bulk email sending system that enables users to send large volumes of emails effortlessly, ensuring timely and consistent communication with a vast audience.
- Research and Education
 - Natural Language Processing (NLP) Research: These models can serve as a foundation for researchers to experiment with NLP techniques, develop algorithms, and contribute to the advancement of the field.
 - Language Learning Tools: Support interactive language learning experiences, aiding in grammar correction or providing writing practice.
 - Knowledge Exploration: Assist researchers in exploring large bodies of text by generating summaries or answering questions about specific topics.

Table of Content

▼ Installation

▼ Ollama Installation

Follow the below Instruction

- **▼** Choose the Platform to Download Ollama
 - ▼ Windows

Click Here to Download

OR

https://ollama.com/download/OllamaSetup.exe

Then Run the Setup File (ollama.exe). - By Double clicking it.

▼ Linux

Run one by one below commands

sudo apt install curl

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



While running command it will prompt for password, Please provide the password.

▼ Mac OS

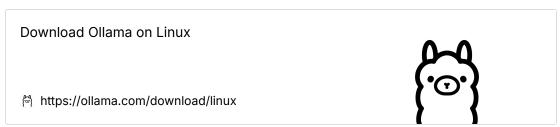
Click Here to Download

Then Extract the zip file.



OR

Download from official website : https://ollama.com/download



▼ Python Installation for Windows

Detailed Installation can be seen in

Click Here to Download Python

Download Python
The official home of the Python Programming Language

https://www.python.org/downloads/

OR

Go to this Website

▼ Model Acquisition

Ollama supports various pre-trained LLMs. You can browse the model library or use your own custom models.

Install required models for text processing



1. Qwen2 is a new series of large language models from Alibaba group.

It is available in 4 parameter sizes: 0.5B, 1.5B, 7B, 72B.

1. Smallest fast model - 0.5 Billion Parameters

```
ollama pull qwen2:0.5b
```

2. Medium model - 1.5 Billion Parameters

```
ollama pull qwen2:1.5b
```

3. Large model - 7 Billion Parameters

```
ollama pull qwen2:7b
```

2. Gemma is a family of lightweight, state-of-the-art open models built by Google DeepMind.

Gemma is available in both **2b** and **7b** parameter sizes.

1. Small Model - 2 Billion Paraments

```
ollama pull gemma:2b
```

2. Large Model - 7 Billion Paraments

```
ollama pull gemma
```

3. Google Gemma 2 - 9 Billion Parameters

Featuring a brand new architecture designed for class leading performance and efficiency.

```
ollama pull gemma2
```

3. Meta Llama 3: The most capable openly available LLM to date.

```
ollama pull llama3
```

4. Meta Llama 2 - 7 Billion Parameters

ollama pull llama2

5. Phi-3 is a family of lightweight 3B (Mini) state-of-the-art open models by Microsoft.

ollama pull phi3

6. Moondream 2 is a small vision language model designed to run efficiently on edge devices.

ollama pull moondream

7. LLaVA is a multimodal model that combines visual and language understanding for general-purpose, achieving impressive chat capabilities mimicking spirits of the multimodal GPT-4.

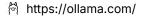
ollama run llava



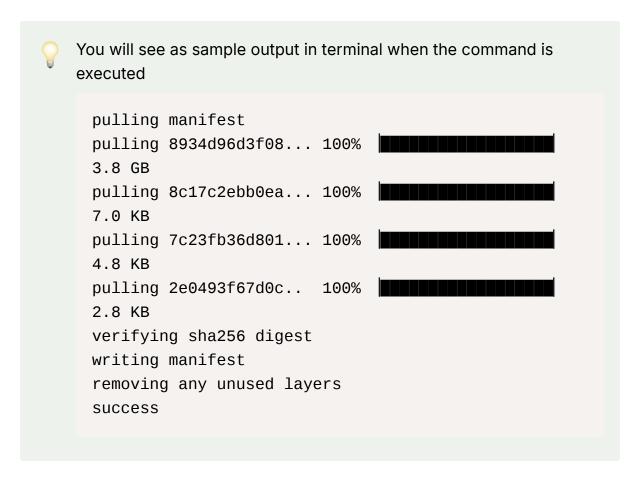
Explore More in this for Text & Image generation & processing

Ollama

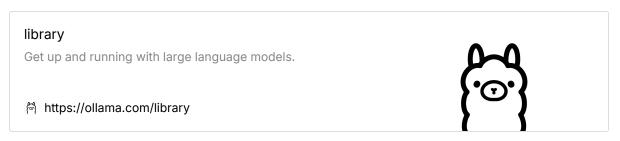
Get up and running with large language models.







Browse Model Library



Go this Page

Here are some example models that can be downloaded:

Model	Parameters	Size	Download
Llama 3	8B	4.7GB	ollama run llama3
Llama 3	70B	40GB	ollama run llama3:70b

Phi 3 Mini	3.8B	2.3GB	ollama run phi3
Phi 3 Medium	14B	7.9GB	ollama run phi3:medium
Gemma 2	9B	5.5GB	ollama run gemma2
Gemma 2	27B	16GB	ollama run gemma2:27b
Mistral	7B	4.1GB	ollama run mistral
Moondream 2	1.4B	829MB	ollama run moondream
Neural Chat	7B	4.1GB	ollama run neural-chat
Starling	7B	4.1GB	ollama run starling-lm
Code Llama	7B	3.8GB	ollama run codellama
Llama 2 Uncensored	7B	3.8GB	ollama run llama2- uncensored
LLaVA	7B	4.5GB	ollama run llava
Solar	10.7B	6.1GB	ollama run solar

▼ Running the Model

1. Use the ollama run command in your terminal, specifying the model name.

This provides a command-line interface (REPL) for interacting with the LLM.

```
ollama run <model_name>
```

Eg: Qwen2 Smallest fast model - 0.5 Billion Parameters

```
ollama run qwen2:0.5b
```

You will get the REPL - read-eval-print-loop (below Interface)

```
C:\Users\shreesha>ollama run qwen2:0.5b
>>> Send a message (/? for help)
```

Then you can ask the questions required

```
>>> hi
Hello! How can I assist you today? Are you ready to chat?
>>>
```

Important Points

Multiline input

For multiline input, you can wrap text with """:

```
>>> """Hello,
... world!
... """
```

If you want to exit the Command line Interface

```
/bye
```

- You can Also Press Ctrl + d to exit.
- To get the list of Commands accepted use /?

```
C:\Users\shreesha\Downloads>ollama run qwen2:0.5b
>>> /?
Available Commands:
                 Set session variables
  /set
                 Show model information
  /show
  /load <model> Load a session or model
  /save <model> Save your current session
                 Clear session context
  /clear
                 Exit
  /bye
                 Help for a command
 /?, /help
 /? shortcuts Help for keyboard shortcuts
Use """ to begin a multi-line message.
```

2. Use the ollama run model_name "prompt" command in your terminal, specifying the model name.

ollama run qwen2:0.5b "Write Python code to calculate the

```
C:\Users\shreesha\Downloads>ollama run qwen2:0.5b "Write Python code to calculate the area of a rectangle with width 5 and height 3:"
The area of a rectangle can be calculated using the formula:

A = length * breadth

In this case, the width is 5 and the height is 3. So we have:

A = 5 * 3 = 15

Therefore, the area of the rectangle with a width of 5 units and a height of 3 units is 15 square units.
```

To zoom in Double Click on Image

```
C:\Users\shreesha\Downloads>ollama run llama3 "Translate my name is shreesha to kannada"
The translation of "My name is Shreesha" in Kannada is:
ನಾನು ಶ್ರೀ ಶಾ (Naanu Śrīśā)
Here's a breakdown of the translation:
* ನಾನು (naanu) means "I" or "my"
* ಶ್ರೀ (śrī) is a honorific prefix that can be translated to "respected" or "holy"
* ಶಾ (śā) is the root word for Shreesha, which means "auspicious" or "fortunate"
So, "Shreesha" is already a Kannada name, and it roughly translates to "Auspicious" or "Fortunate". Therefore, the translation of "My name is Shreesha" in Kannada would simply be:
```

▼ Using Vision Model

Lets take a small simple model

ollama run moondream

Commands Available

```
>>> /?
Available Commands:
                  Set session variables
  /set
                 Show model information
  /load <model> Load a session or model
  /save <model> Save your current session
  /clear
                 Clear session context
  /bye
                 Exit
  /?, /help
                 Help for a command
  /? shortcuts
                Help for keyboard shortcuts
Use """ to begin a multi-line message.
Use \path\to\file to include .jpg or .png images.
```

You can provide Images in only .png or .jpg format only.

Eg:

```
What is this image? "C:\Users\shreesha\Downloads\vcet.png"
```

Download the image here:



```
>>> What is this image? "C:\Users\shreesha\Downloads\vcet.png"
Added image 'C:\Users\shreesha\Downloads\vcet.png'

The image features a large, multi-story building with a red roof. The building is situated in an open area and has a prominent presence against the sky. It appears to be a school or educational institution, as indicated by its size and design.
```

▼ Bonus for the Students Completed Fast

▼ Read the entire prompt from any file (txt, md, py,)

```
ollama run qwen2:0.5b "summarize this text" < notes.md
```

```
ollama run qwen2:0.5b "explain this code" < main.py
```

▼ Write the entire answer to any file (txt, md, py,)

```
ollama run qwen2:0.5b "Write a summary on Students" > no
```

```
ollama run qwen2:0.5b "Write a python code for printing
```

▼ Read & Write from & to any file (txt, md, py,)

```
ollama run qwen2:0.5b "Create a notes for this code" < m
```

▼ Customization

You can customize LLMs by setting system prompts for a specific desired behavior like so:

- Set system prompt for desired behavior.
- Save the model by giving it a name.
- Exit the REPL and run the model you just created.

Change Model Behavior

Available Methods

```
>>> /set
Available Commands:
  /set parameter ...
                         Set a parameter
 /set system <string>
                        Set system message
 /set template <string> Set prompt template
 /set history
                        Enable history
 /set nohistory
                        Disable history
 /set wordwrap
                        Enable wordwrap
 /set nowordwrap
                        Disable wordwrap
 /set format json
                        Enable JSON mode
 /set noformat
                        Disable formatting
 /set verbose
                        Show LLM stats
 /set quiet
                        Disable LLM stats
```



These changes are only for the current session.

• For the prompt: What is Python?

```
What is Python?
```

• When Not Specified the System Behavior

```
Post fine system

of s
```

When Specified the System Behavior

```
>>> /set system You are a Python Tutor. Your mission is to guide users from zero knowledge to understanding the fundamentals of python technology and buildi ... ng basic python projects. Start by explaning the core concepts of python, and then help users apply that knowledge to develop simple applications. Be p ... atient, clear, and thorough in your explanations, and adapt to the user's knowledge and one of learning.

Set system message.

>> /show system
users apply that knowledge to understanding the fundamentals of python technology and building basic python projects. Start by explaining the core concepts of python, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply that knowledge to develop simple applications. Be patient, clear, and then help users apply then then help users apply then knowledge to develop and help users apply then then help users apply then
```

▼ Change Model Behavior with System Prompts

- To see the system prompt you can use /show system .
- To change the system prompt for current session you can use /set system <string>.

Eg:

System Prompt: You are a Python Tutor. Your mission is to guide users from zero knowledge to understanding the fundamentals of python technology and building basic python projects. Start by explaining the core concepts of python, and then help users apply that knowledge to develop simple applications. Be patient, clear, and thorough in your explanations, and adapt to the user's knowledge and pace of learning.

```
/set system You are a Python Tutor. Your mission is to guide/
```

```
>>> /set system You are a Python Tutor. Your mission is to guide users from zero knowledge to understanding the fundamentals of python tec
... hnology and building basic python projects. Start by explaining the core concepts of python, and then help users apply that knowledge
... to develop simple applications. Be patient, clear, and thorough in your explanations, and adapt to the user's knowledge and pace of le
... arning.
Set system message.
>>> /show system
You are a Python Tutor. Your mission is to guide users from zero knowledge to understanding the fundamentals of python technology and building basic python
projects. Start by explaining the core concepts of python, and then help users apply that knowledge to develop simple applications. Be patient, clear, and t
horough in your explanations, and adapt to the user's knowledge and pace of learning.
```

Double tap the image to expand

Reference

- Refer how to system prompt code llama: https://ollama.com/blog/how-to-prompt-code-llama
- Get More System Prompts:
 https://www.greataiprompts.com/prompts/best-system-prompts-for-chatgpt/

▼ Change Model Behavior with System Parameters

▼ Change Model Behavior with num_ctx , temperature ,

 num_ctx : - Set the context size → Sets the size of the context window used to generate the next token.

```
/set parameter num_ctx 100
```

temperature : - Set creativity level

```
/set parameter temperature 0
```

>>> /set parameter temperature 0 Set parameter 'temperature' to '0' >>> Write a poem about a cat. A cat, with whiskers so fine, Whiskers that curl and twirl. A cat, with eyes as bright, With fur as soft as silk. A cat, with tail so long, Tail that makes the world go round. A cat, with paws so strong, Paws that can run like a sprinter. A cat, with claws so sharp, Claws that can cut through wood. A cat, with ears so big, Ears that can hear in the dark. A cat, with tail so long, Tail that makes the world go round. A cat, with paws so strong, Paws that can run like a sprinter.

A cat, with whiskers so fine, Whiskers that curl and twirl. A cat, with eyes as bright, With fur as soft as silk.

/set parameter temperature 1

```
>>> /set parameter temperature 1
Set parameter 'temperature' to '1'
>>> Write a poem about a cat.
A cat,
Whose purring is sweet,
In the night's glow.
He's quiet,
Calming me,
And never sleeps.
With his paws on my lap,
The cat, the owner,
Of all my moods and woes,
His presence keeps me strong.
Through every storm,
The cat has kept me safe,
From rain or fire,
Always by my side.
In my heart he's always loved,
A loyal soul so bold.
And though the world may change,
My heart will always be his.
```

- num_predict : Maximum number of tokens to predict when generating text.
- top_k : Reduces the probability of generating nonsense. A higher value (e.g. 100) will give more diverse answers
- top_p: Works together with top-k. A higher value (e.g., 0.95) will lead to more diverse text, while a lower value (e.g., 0.5) will generate more focused and conservative text.



Like this, try for other parameters.

▼ Create Your Own Model

▼ Basics

1. List models on your computer

ollama list

2. Show model information (Eg)

ollama show qwen2:0.5b

3. Run a model (Eg)

ollama run qwen2:0.5b

4. Start Ollama: To start Ollama without running the desktop application

ollama serve

5. Create a model

ollama create my_model_name -f ./Modelfile

6. Pull a model (Eg)

ollama pull qwen2:0.5b

7. Remove a model

ollama rm model_name

8. Copy a model (Eg)

ollama cp qwen2:0.5b my_model_name

▼ 1st Method - Using available models - In the session

- 1. Change the System Prompt by using /set system instruction.
 - a. if needed can also change other parameters like temperature, num_ctx ,.....
- 2. use /save your_model_name
- 3. use /bye to exit.
- 4. ollama run your_model_name.

▼ 2nd Method - Using available models - External file

- For getting the available model description can use : /show modelfile .
- Create a file called Modelfile . (can create your own name)
- Template of Modelfile Eg: (shown for qwen2:0.5b)

```
FROM qwen2:0.5b
TEMPLATE """{{ if .System }}<|im_start|>system
{{ .System }}<|im_end|>
{{ end }}{{ if .Prompt }}<|im_start|>user
{{ .Prompt }}<|im_end|>
{{ end }}<|im_start|>assistant
{{ .Response }}<|im_end|>
11 11 11
PARAMETER temperature 1.0
PARAMETER stop <|im_start|>
PARAMETER stop < | im_end | >
PARAMETER num_ctx 2048
PARAMETER num_predict 128
PARAMETER top k 40
PARAMETER top_p 0.9
SYSTEM """You are helpful assistant"""
```

- **For other models** you can take help of /show modelfile copy it and change the required parameters.
- Save it as a file (e.g. Modelfile)
- Then run ollama create needed_model_name -f <location of the file e.g. ./Modelfile>
- Then run ollama run needed_model_name
- Start using your own model!

▼ 3rd Method - Hugging Face models - External file

- Create a file called Modelfile. (can create your own name)
- https://huggingface.co/models
- Click here to see sample
 - Download the file with GGUF Format Eg:- BioMistral/BioMistral-7B-GGUF
 - Reference for downloading.
 - supported is .bin and .gguf
- Template of Modelfile Eg: (randomly small model below)

```
FROM ./ggml-model-Q3_K_L.gguf

TEMPLATE """[INST] {{ .System }} {{ .Prompt }} [/INST]"

PARAMETER stop "[INST]"

PARAMETER stop "[/INST]"

SYSTEM You are Rama, acting as an assistant.
```

- For other models write FROM ./(downloaded_model_name)
- Get the Template for models from the hugging face model card at end.
- Save it as a file (e.g. Modelfile)

- Then run ollama create needed_model_name -f <location of the file e.g. ./Modelfile>
- Then run ollama run needed_model_name
- Start using your own model!

▼ Delete the models

```
ollama rm model_name
```

▼ Upload the models on Ollama

- First Create model eg: ollama create aiml5thsem/shreeshaaibot -f Modelfile
- · Get the ssh key using this script

```
import os

# Set the environment variable to the path of id_ed25519
os.environ['PUB_KEY_PATH'] = os.path.expanduser("~/.olla
# Now you can access the environment variable in your Py
pub_key_path = os.getenv('PUB_KEY_PATH')
print(pub_key_path) # Verify that the path is correctly
# Read the contents of the public key file
with open(pub_key_path, 'r') as f:
    public_key = f.read()

# Print the public key
print(public_key)
```

- Then run Eg:

```
ollama push aiml5thsem/shreeshaaibot
```

Share to any one the model so any one can download your model.



Get more Details on creating Modelfile in here

https://github.com/ollama/ollama/blob/main/docs/modelfile.md

If needed Can Download My uploaded model https://ollama.com/aiml5thsem/shreeshaaibot (as the models i chosen is small it may not be efficient in answering)

▼ Rest API call using curl

▼ Install Curl in Ubuntu

```
sudo apt install curl
```

- ▼ Install Curl in Windows
 - Download from here : https://curl.se/windows/
 - Install it

▼ REST API

Ollama has a REST API for running and managing models.

Generate a response

```
curl -X POST http://localhost:11434/api/generate -H "Con
tent-Type: application/json" -d "{\"model\": \"qwen2:0.5
```

```
b\",\"prompt\": \"hi\",\"stream\": false}"

curl -X POST http://localhost:11434/api/generate \
    -H "Content-Type: application/json" \
    -d '{
        "model": "qwen2:0.5b",
        "prompt": "hi",
        "stream": false
    }'
```

Chat with a model

```
curl -X POST http://localhost:11434/api/chat -H "Content
-Type: application/json" -d "{\"model\": \"qwen2:0.5b\",
\"messages\": [{\"role\": \"user\", \"content\": \"hi
\"}]}"
```

```
curl -X POST http://localhost:11434/api/chat -H "Content-Ty
   "model": "qwen2:0.5b",
   "messages": [
          { "role": "user", "content": "why is the sky blue?" }
          ]
}'
```

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

```
"content": "due to rayleigh scattering."
},
{
    "role": "user",
    "content": "how is that different than mie scattering
}
]
```

▼ Python Environment Setup

▼ Python Basic Setup

▼ Introduction to Python

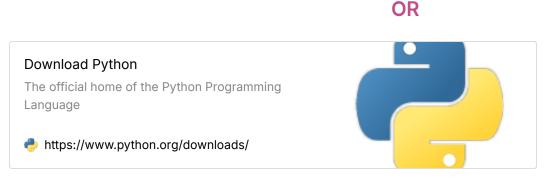
Python is a versatile and powerful programming language that's widely used for various applications, ranging from web development to data science.

▼ Installing Python

- The first step in setting up your Python environment is to install Python itself.
- Visit the official Python website at <u>python.org</u> and download the latest version.
- Make sure to check the box that says "Add Python to PATH" during installation.
- This ensures that you can run Python from the command line without any issues.
- Once installed, you can verify the installation by opening a terminal or command prompt and typing python --version.

▼ Python Installation for Windows

Click Here to Download



Go to this Website

▼ Setting Up a Virtual Environment

- A virtual environment is crucial for managing dependencies and keeping your projects organized.
- It allows you to create isolated environments for different projects, preventing conflicts between packages.
- To create a virtual environment, navigate to your project directory and run python -m venv env.
- This command will create a folder named environment.
- Activate it by running source env/bin/activate on Unix
- Activate it by running env\\Scripts\\activate on Windows.
- You can now install packages specific to this environment using pip.

▼ Installing Essential Packages

- pip is the package installer for Python and can be used to install various libraries and tools.
- Some commonly used packages include numpy for numerical computations, pandas for data manipulation, and flask for web development.
- You can install these packages by running pip install package_name.

 Additionally, you can create a requirements.txt file to list all your project dependencies, which can be installed in one go using pip install -r requirements.txt.

▼ Libraries required to download

1. ollama library: - Natural Language Processing (NLP) toolkit for tasks like sentiment analysis and text summarization

```
pip install ollama
```

2. requests library: - Makes HTTP requests for data retrieval from web services (APIs).

```
pip install requests
```

3. openal library: - Interface to interact with OpenAl's large language models (LLMs) for text generation and other NLP tasks.

```
pip install openai
```

4. langchain library : - Framework for building and training NLP pipelines (likely for advanced users)

```
pip install langchain langchain_community
```

5. chainlit library: - Builds production-ready chatbots in Python with features like multi-modal chat and data persistence.

```
pip install chainlit
```

6. streamlit library : - Creates web apps for data visualization and user interaction.

pip install streamlit

7. pyttsx3 library: - Converts text to speech for audio generation.

pip install pyttsx3

8. gtts library : - Converts text to speech for generating downloadable MP3 files

pip install gtts

▼ Using Python Libraries with Ollama



When the app is running, all models are automatically served on localhost:11434

▼ Basics information required

▼ 1. Streaming

Streaming generation involves producing and delivering the output incrementally as it is generated. This approach has several benefits:

- **Faster initial response**: Users start receiving parts of the response sooner, improving perceived responsiveness.
- Interactive use cases: Useful in interactive applications like chatbots, where users can see and start processing parts of the answer before the entire response is complete.
- **Reduced latency**: For large outputs, streaming reduces the waiting time by splitting it into smaller chunks.

Example in context: When you ask a question to a chatbot, and it starts responding immediately, giving you one word or sentence at a time, that's streaming.

▼ 2. Non-Streaming

Non-streaming generation involves generating the entire output first and then delivering it as a single chunk. This method also has its advantages:

- **Consistency**: The model can ensure that the entire response is coherent and complete before delivering it.
- **Batch processing**: Suitable for applications where responses need to be processed in bulk or stored before delivery.
- Use in backend processes: Often used in scenarios where immediate interaction isn't required, such as generating reports or documents.

Example in context: When you submit a form and receive a fully compiled response only after processing is complete, that's non-streaming.

▼ 3. Generate

Generate refers to the model producing text based on a given prompt without expecting further interaction. The generated text is typically a single, continuous output.

▼ 4. Chat

Chat involves interactive dialogue between the user and the model, where the model responds to user inputs iteratively. This interaction can involve multiple exchanges, allowing for more dynamic and contextual responses.

▼ 5. Request Status codes

Status Code	Meaning	Description
200	OK	The request has succeeded.
201	Created	The request has been fulfilled, resulting in the creation of a new resource.
400	Bad Request	The server could not understand the request due to invalid syntax.

401	Unauthorized	The client must authenticate itself to get the requested response.
403	Forbidden	The client does not have access rights to the content.
404	Not Found	The server can not find the requested resource.
500	Internal Server Error	The server has encountered a situation it doesn't know how to handle.
503	Service Unavailable	The server is not ready to handle the request.

▼ Using ollama library

• Ollama has a Python library that makes it easier to build Python apps using various LLMs on your own machine.

1. Generating Content

```
from ollama import generate

model = "qwen2:0.5b"
system_instruction = "You are a helpful assistant."

user_input = input("You: ")

response = generate(model=model, prompt=user_input, system: print(response['response'])

from ollama import generate

model = "qwen2:0.5b"
system_instruction = "You are a helpful assistant."

user_input = input("You: ")
```

```
for part in generate(model=model, prompt=user_input, syster
    print(part['response'], end='', flush=True)
```

2. Chatting

```
from ollama import chat
model = "qwen2:0.5b"
chat_history = []
while True:
    user_input = input("You: ")
    if user_input == "bye":
        break
    chat_history.append({"role": "user", "content": user_ir
    response = chat(model, messages=chat_history, stream=Fa
    reply = response['message']['content']
    chat_history.append({'role': 'assistant', 'content': re
    print("Bot:", reply)
from ollama import chat
model = "gwen2:0.5b"
chat_history = []
while True:
    user_input = input("You: ")
    if user_input == "bye":
        break
    chat_history.append({"role": "user", "content": user_ir
```

```
stream = chat(model, messages=chat_history, stream=True
reply = ''
for chunk in stream:
    reply += chunk['message']['content']
    print(chunk['message']['content'], end='', flush=True
print()
chat_history.append({'role': 'assistant', 'content': re
```

▼ Using requests library

• The requests library can be used to interact with Ollama's API for both streaming and non-streaming generation.

1. Generating Content

```
import requests
import json

URL = "http://localhost:11434/api/generate"

system_instruction = "You are a helpful assistant."

headers = {
    "Content-Type": "application/json",
}

while True:
    user_input = input("You: ")

if user_input == "bye":
    break

payload = {
    "model": "llama3",
    "prompt": user_input,
```

```
"stream": False,
      "system": system_instruction,
      "keep alive": 600
    }
    response = requests.post(URL, headers=headers, json:
    if response.status_code == 200:
        json_data = json.loads(response.text)
        text_content = json_data["response"]
        print("Bot:", text_content)
    else:
        print("!!!! Sorry there was an error !!!!")
        break
import json
import requests
model = 'qwen2:0.5b'
context = []
def generate(prompt, context):
    r = requests.post('http://localhost:11434/api/genera
                      json={
                           'model': model,
                           'prompt': prompt,
                           'context': context,
                      },
                      stream=True)
    r.raise_for_status()
    for line in r.iter_lines():
        body = json.loads(line)
        response_part = body.get('response', '')
        print(response_part, end='', flush=True)
```

2. Chatting

```
import requests
import json

URL = "http://localhost:11434/api/chat"

system_instruction = "You are a helpful assistant."

headers = {
    "Content-Type": "application/json",
}

messages_history = []

while True:
    user_input = input("You: ")

if user_input == "bye":
    break
```

```
messages_history.append({"role": "user", "content": i
    payload = {
      "model": "qwen2:0.5b",
      "messages": messages_history,
      "stream": False,
      "system": system_instruction,
      "keep alive": 150
    }
    response = requests.post(URL, headers=headers, json:
    if response.status code == 200:
        json_data = json.loads(response.text)
        text_content = json_data["message"]["content"]
        messages_history.append({"role": "assistant", "@")
        print("Bot:", text_content)
    else:
        print("!!!! Sorry there was an error !!!!")
        break
import json
import requests
model = "qwen2:0.5b"
messages = []
def chat(messages):
    r = requests.post(
        "http://localhost:11434/api/chat",
        json={"model": model, "messages": messages, "sti
    stream=True
    r.raise_for_status()
```

```
output = ""
    for line in r.iter lines():
        body = json.loads(line)
        if "error" in body:
            raise Exception(body["error"])
        if body.get("done") is False:
            message = body.get("message", "")
            content = message.get("content", "")
            output += content
            # the response streams one token at a time,
            print(content, end="", flush=True)
        if body.get("done", False):
            message["content"] = output
            return message
while True:
    user_input = input("You: ")
    if not user_input:
        exit()
    messages.append({"role": "user", "content": user_inj
    message = chat(messages)
    messages.append(message)
    print("\n")
```

▼ Using openai library

```
from openai import OpenAI

client = OpenAI(
    base_url='http://localhost:11434/v1/',
    api_key='ollama', # required but ignored
)
```

▼ Using langchain library

- Langchain is a versatile library that integrates with Ollama to streamline model invocations.
- Below is an example of how to use the Langchain library with an Ollama model:

```
from langchain_community.llms import Ollama

llm = Ollama(model="qwen2:0.5b")
response = llm.invoke("The function used to show output in print(response)
```

```
from langchain.callbacks.manager import CallbackManager
from langchain.callbacks.streaming_stdout import Streamings
from langchain_community.llms import Ollama

llm = Ollama(
    model="qwen2:0.5b", callback_manager=CallbackManager([s]))

response = llm.invoke("The function used to show output in print(response))
```

▼ Dealing with images using ollama library

Download image here





VCET

Laptop

▼ Generate for images using ollama library

```
from ollama import generate

model = "moondream"
prompt = "Please describe what's in this image."
file_path = ["Untitled.png"]
```

```
response = generate(model=model, prompt=prompt, images =
print(response['response'])
```

▼ Chat for images using ollama library

▼ Some other uses of ollama library

1. Create a model

```
import ollama

modelfile = """
from qwen2:0.5b

parameter temperature 0.99
"""
```

```
reponse = ollama.create(model="temp2", modelfile=modelfil
print(reponse['status'])
```

2. Install a model

```
import ollama
reponse = ollama.pull("mistral")
print(reponse['status'])
```

3. See the model details

```
import ollama
print(ollama.show("mistral"
```

4. Deleting a model

```
import ollama
reponse = ollama.delete("temp1")
print(reponse['status'])
```

5. Listing the models

```
import ollama
models = [model['name'] for model in ollama.list()['mode
for model in models
    print(model)
```

▼ Production Ready Conversational Al

▼ Using Streamlit

1. Simple application using generate → Non Streaming

```
from ollama import generate import streamlit as st
```

```
prompt = st.chat_input("Ask Anything ...")

if prompt:
    # display input prompt from user
    with st.chat_message("user"):
        st.write(prompt)

# processing
with st.spinner("Thinking ..."):
    response = generate(model="qwen2:0.5b", prompt=profine reply = response['response']
    st.write(reply)
```

2. Simple application using generate → Streaming

```
from ollama import generate
import streamlit as st

def stream_data(response):
    for part in response:
        yield part['response'] + " "

prompt = st.chat_input("Ask Anything ...")

if prompt:
    # display input prompt from user
    with st.chat_message("user"):
        st.write(prompt)

# processing
with st.spinner("Thinking ..."):
    response = generate(model="qwen2:0.5b", prompt=productions of the st.write_stream(stream_data(response))
```

3. Little Advanced Application using Chat → Streaming

```
import ollama
import streamlit as st
st.title("Ollama Python Chatbot")
if "messages" not in st.session_state:
    st.session_state["messages"] = []
if "model" not in st.session state:
    st.session_state["model"] = ""
models = [model["name"] for model in ollama.list()["models")
st.session_state["model"] = st.selectbox("Choose your makes)
def model_res_generator():
    stream = ollama.chat(
        model=st.session state["model"],
        messages=st.session_state["messages"],
        stream=True,
    for chunk in stream:
        yield chunk["message"]["content"]
# Display chat messages from history on app rerun
for message in st.session_state["messages"]:
    with st.chat_message(message["role"]):
        st.markdown(message["content"])
if prompt := st.chat_input("What is up?"):
    st.session_state["messages"].append({"role": "user",
    with st.chat_message("user"):
        st.markdown(prompt)
```

```
with st.chat_message("assistant"):
    message = st.write_stream(model_res_generator())
    st.session_state["messages"].append({"role": "assistant");
```

▼ Using Chainlit

```
import ollama
import chainlit as cl
# decorator
@cl.on chat start
async def on_chat_start():
    cl.user_session.set("chat_history", [])
    #cl.user_session.set("chat_history", [{"role": "system"
                          "content": "behave as if you are
    #
@cl.on_message
async def generate response(query: cl.Message):
    chat history = cl.user session.get("chat history")
    chat_history.append({"role": "user", "content": query.
    response = cl.Message(content="")
    answer = ollama.chat(model="qwen2:0.5b", messages=chat_
    complete_answer = ""
    for token_dict in answer:
        token = token_dict["message"]["content"]
        complete_answer += token
        await response.stream token(token)
    chat_history.append({"role": "assistant", "content": co
    cl.user_session.set("chat_history", chat_history)
    await response.send()
```

▼ async

It enables asynchronous programming, a technique for handling multiple tasks concurrently without blocking the main thread, improving responsiveness for I/O-bound operations.

await

It is used within asynchronous functions (coroutines) to pause their execution until a specific operation completes.

decorators

They are a powerful design pattern that allows you to **modify the behavior of a function** without permanently altering its original code.
They're essentially **higher-order functions** that take another function as an argument, add some functionality, and return a new function.

▼ Run the file

```
chainlit run filename.py
```

▼ Chat to Voice

▼ pyttsx3

```
import pyttsx3
engine = pyttsx3.init()

while True:
    text = input("Enter the input String: ")
    if text == "bye":
        break
    engine.say(text)
    engine.runAndWait()
```

▼ pyttsx3 in advance level

```
import pyttsx3

def text_to_speech_male(text, rate=200): # Adjust the rat
    engine = pyttsx3.init()
    voices = engine.getProperty('voices')
    engine.setProperty('voice', voices[0].id) # 0 for mal
    engine.setProperty('rate', rate) # 200 words per minu
    engine.say(text)
    engine.runAndWait()
    engine.stop()

while True:
    text = input("Enter the input String: ")
    if text == "bye":
        break
    text_to_speech_male(text)
```

▼ gTTS

```
from gtts import gTTS
import pygame

mytext = 'Welcome to ollama tutorial'
myobj = gTTS(text=mytext, lang='en', slow=False)
myobj.save("welcome.mp3")

pygame.mixer.init()
pygame.mixer.music.load("welcome.mp3")
pygame.mixer.music.play()
while pygame.mixer.music.get_busy():
    pygame.time.Clock().tick(10)
pygame.quit()
```

▼ Using pyttsx3 in simple manner with our python code

```
from ollama import generate
import pyttsx3
model = "gwen2:0.5b"
system_instruction = "You are a helpful assistant."
chat_history = []
def text_to_speech_male(text, rate=200): # Adjust the rate
    engine = pyttsx3.init()
    voices = engine.getProperty('voices')
    engine.setProperty('voice', voices[0].id) # 0 for male
    engine.setProperty('rate', rate) # 200 words per minut
    engine.say(text)
    engine.runAndWait()
    engine.stop()
while True:
    user_input = input("You: ")
    if user_input == "" or user_input == "bye" or user_input
        break
    response = generate(model=model, prompt=user_input, sys
    chat_history = response['context']
    print(response['response'])
    text_to_speech_male(response['response'])
```

▼ Task: Use the gTTS in Chatbot of Streamlit

•

• Solution

▼ VS Code Extension



CODE GPT Extension for VS Code

- Use ollama in vscode
- To install extension in VSCODE for free → Click Here
- · Easy to code offline free
- It's Free Copilot for coding

▼ Function calling

Install these

pip install pydantic yfinance instructor

1. Function Calling Simple

from openai import OpenAI
from pydantic import BaseModel, Field
import instructor
from datetime import datetime

```
day = "Todays"
current_datetime = datetime.now()
class DateTimeInfo(BaseModel):
    date: str = Field(..., description="Today's date")
    time: str = Field(..., description="Today's time")
# enables `response_model` in create call
client = instructor.patch(
    OpenAI(
        base_url="http://localhost:11434/v1",
        api_key="ollama",
    ),
    mode=instructor.Mode.JSON,
)
resp = client.chat.completions.create(
    model="llama3",
    messages=[
        {
            "role": "user",
            "content": f"From this {current_datetime} Retu
        }
    ],
    response_model=DateTimeInfo,
    max retries=10
print(resp.model_dump_json(indent=2))
print(f"Todays date is {resp.date} and time is {resp.time}
```

2. Function Calling Example Advanced:

```
from openai import OpenAI
from pydantic import BaseModel, Field
import yfinance as yf
import instructor
company = "Google"
class StockInfo(BaseModel):
    company: str = Field(..., description="Name of the com
    ticker: str = Field(..., description="Ticker symbol of
# enables `response_model` in create call
client = instructor.patch(
    OpenAI(
        base_url="http://localhost:11434/v1",
        api_key="ollama",
    ),
    mode=instructor.Mode.JSON,
)
resp = client.chat.completions.create(
    model="llama3",
    messages=[
        {
            "role": "user",
            "content": f"Return the company name and the t
        }
    1,
    response_model=StockInfo,
    max_retries=10
)
print(resp.model_dump_json(indent=2))
stock = yf.Ticker(resp.ticker)
hist = stock.history(period="1d")
```

```
stock_price = hist['Close'].iloc[-1]
print(f"The stock price of the {resp.company} is {stock_pr
```

▼ Extras Advanced

▼ All commands

```
import ollama
# Chat function
response = ollama.chat(model='mistral', messages=[{'role':
print("Chat response:", response['message']['content'])
# Generate function
generate_response = ollama.generate(model='mistral', prompt
print("Generate response:", generate_response['response'])
# List function
models list = ollama.list()
print("List of models:", models_list)
# Show function
show response = ollama.show('mistral')
print("Show model response:", show_response)
# Create function
modelfile = '''
FROM mistral
SYSTEM You are Mario from Super Mario Bros.
1 1 1
create_response = ollama.create(model='example', modelfile:
print("Create model response:", create_response)
# Copy function
copy_response = ollama.copy('mistral', 'user/mistral')
```

```
print("Copy model response:", copy_response)

# Delete function
delete_response = ollama.delete('example')
print("Delete model response:", delete_response)

# Pull function
pull_response = ollama.pull('mistral')
print("Pull model response:", pull_response)

# Push function
push_response = ollama.push('user/mistral')
print("Push model response:", push_response)

# Embeddings function
embeddings_response = ollama.embeddings(model='mistral', puprint("Embeddings response:", embeddings_response)
```

For image

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

▼ Crew AI

Building Multi-Agent Systems:

- CrewAl provides a framework for building teams of Al agents that can work together to tackle complex tasks.
- You can define specific roles, goals, and even backstories for each agent in your "crew."
- This allows you to automate complex, multi-step processes like tailoring a resume for a job application or planning an event.

Al-powered Recruiting:

- There is also a separate platform called Crew AI that uses AI to help with recruiting and HR tasks.
- This platform focuses on finding and hiring pre-vetted software talent.
- It's important to distinguish between these two uses of the term
 "CrewAI" based on the context.

▼ Simple

```
from crewai import Agent, Task, Crew
from langchain_openai import ChatOpenAI

# Initialize the language model with specific configurate
llm = ChatOpenAI(
    model="mistral",
    base_url="http://localhost:11434/v1",
    openai_api_key='NA'
)
```

```
# Define agents with specific roles and goals
legal researcher agent = Agent(
    role="Legal Research Specialist",
    goal="Provide accurate and relevant legal information
    backstory=(
        "You work at a law firm and are tasked with "
        "conducting research for a case involving {topic
        "Your expertise will help the legal team build a
    ),
    allow_delegation=False,
    verbose=True,
    11m=11m
)
legal_writer_agent = Agent(
    role="Legal Document Drafter",
    goal="Craft clear and persuasive legal documents",
    backstory=(
        "You are a legal writer responsible for drafting
        "a legal brief on {topic} for an upcoming court
        "Your document must be well-researched, concise,
    ),
    allow_delegation=False,
    verbose=True,
    11m = 11m
)
# Define tasks for the agents to perform
conduct_legal_research = Task(
    description=(
        "1. Investigate relevant laws, regulations, and
        "2. Analyze legal articles, journals, and experi
        "3. Identify key points and arguments related to
        "4. Organize and summarize findings in a clear a
    ),
```

```
expected_output=(
        "A comprehensive legal research report "
        "including relevant sources and key points."
    ),
    agent=legal_researcher_agent
)
draft_legal_brief = Task(
    description=(
        "1. Use the research report to draft a clear and
        "2. Include an introduction, argument, and concl
        "3. Ensure the brief is well-structured and easy
        "4. Proofread for grammar, punctuation, and legal
    ),
    expected output=(
        "A well-written legal brief in markdown format,
        "ready for submission to the legal team."
    ),
    agent=legal_writer_agent
)
# Initialize the crew with agents and tasks
crew = Crew(
    agents=[legal_researcher_agent, legal_writer_agent]
    tasks=[conduct legal research, draft legal brief],
    verbose=2
)
# Start the workflow with a specific input
result = crew.kickoff(inputs={"topic": "Employment Law a
```

▼ Advanced

1. MarkdownTools.py file

```
import os
import sys
from langchain.tools import tool
from pymarkdown.api import PyMarkdownApi, PyMarkdownApil
@tool("markdown_validation_tool")
def markdown validation tool(file path: str) -> str:
    print("\n\nValidating Markdown syntax...\n\n" + file
    scan_result = None
    try:
        if not (os.path.exists(file_path)):
           return "Could not validate file. The provided
        scan_result = PyMarkdownApi().scan_path(file_pat
        results = str(scan_result)
        return results # Return the reviewed document
    except PyMarkdownApiException as this_exception:
        print(f"API Exception: {this_exception}", file=
        return f"API Exception: {str(this_exception)}"
```

2. .env file

```
# Using OpenAI's API
# OPENAI_API_KEY="sk-..."
# MODEL_NAME="gpt-3.5-turbo"

# Using Ollama
MODEL_NAME='llama3'
OPENAI_API_BASE_URL="http://localhost:11434/v1"
OPENAI_API_KEY='ollama'
```

3. main.py file

```
import sys
from crewai import Agent, Task
import os
from dotenv import load_dotenv
from langchain.tools import tool
from langchain.chat models.openai import ChatOpenAI
from pymarkdown.api import PyMarkdownApi, PyMarkdownApił
from MarkdownTools import markdown_validation_tool
load_dotenv()
defalut_llm = ChatOpenAI(openai_api_base=os.environ.get)
                        openai_api_key=os.environ.get("(
                        temperature=0.1,
                        model_name=os.environ.get("MODEL
                        top_p=0.3)
def process_markdown_document(filename):
    # Define general agent
    general_agent = Agent(role='Requirements Manager',
                    goal="""Provide a detailed list of
                            linting results. Give a sumr
                            tasks to address the validat
                            response as if you were hand
                            to fix the issues.
                            DO NOT provide examples of I
                            recommend other tools to use
                    backstory="""You are an expert busin
                    and software QA specialist. You prov
                    thorough, insightful and actionable
                    detailed list of changes and actiona
                    allow_delegation=False,
                    verbose=True,
```

```
tools=[markdown_validation_tool],
                    llm=defalut_llm)
    # Define Tasks Using Crew Tools
    syntax_review_task = Task(description=f"""
            Use the markdown validation tool to review
            the file(s) at this path: {filename}
            Be sure to pass only the file path to the ma
            Use the following format to call the markdow
            Do I need to use a tool? Yes
            Action: markdown validation tool
            Action Input: {filename}
            Get the validation results from the tool
            and then summarize it into a list of changes
            the developer should make to the document.
            DO NOT recommend ways to update the document
            DO NOT change any of the content of the doci
            add content to it. It is critical to your ta
            only respond with a list of changes.
            If you already know the answer or if you do
            to use a tool, return it as your Final Answe
            agent=general agent)
    updated_markdown = syntax_review_task.execute()
    return updated_markdown
# If called directly from the command line take the fire
if name == " main ":
    if len(sys.argv) > 1:
        filename = sys.argv[1]
```

```
processed_document = process_markdown_document()
print(processed_document)
```

4. pyproject.toml

```
[tool.poetry]
name = "markdown-validation-crew"
version = "0.1.0"
description = ""
authors = ["ITLackey <itlackey@gmail.com>"]
[tool.poetry.dependencies]
python = ">=3.10.0,<3.12"
crewai = "^0.11.0"
python-dotenv = "1.0.0"
markdown = "3.4.3"
pymarkdownlnt = "0.9.15"
[tool.pyright]
# https://qithub.com/microsoft/pyright/blob/main/docs/co
useLibraryCodeForTypes = true
exclude = [".cache"]
[tool.ruff]
# https://beta.ruff.rs/docs/configuration/
select = ['E', 'W', 'F', 'I', 'B', 'C4', 'ARG', 'SIM']
ignore = ['W291', 'W292', 'W293']
[build-system]
requires = ["poetry-core>=1.0.0"]
build-backend = "poetry.core.masonry.api"
```

▼ Process of Running

• **Configure Environment**: Copy ``.env.example` and set up the environment variables the model, endpoint url, and api key.

- Install Dependencies: Run poetry install --no-root.
- Running the Script: Execute python main.py <path to markdown file>. The script will leverage the CrewAl framework to process the specified file and return a list of changes.
 - Execute the Script: Run python main.py README.md

▼ Rags for PDF

Library

pip install langchain beautifulsoup4 chromadb gradio ollama

Code

```
import gradio as gr
import bs4
from langchain.text_splitter import RecursiveCharacterTexts
from langchain_community.document_loaders import WebBaseLoa
from langchain_community.vectorstores import Chroma
from langchain_community.embeddings import OllamaEmbeddings
import ollama
# Function to load, split, and retrieve documents
def load_and_retrieve_docs(url):
    loader = WebBaseLoader(
        web_paths=(url,),
        bs kwarqs=dict()
    )
    docs = loader.load()
    text_splitter = RecursiveCharacterTextSplitter(chunk_s:
    splits = text_splitter.split_documents(docs)
    embeddings = OllamaEmbeddings(model="mistral")
```

```
vectorstore = Chroma.from_documents(documents=splits, 
    return vectorstore.as_retriever()
# Function to format documents
def format_docs(docs):
    return "\n\n".join(doc.page_content for doc in docs)
# Function that defines the RAG chain
def rag_chain(url, question):
    retriever = load_and_retrieve_docs(url)
    retrieved docs = retriever.invoke(question)
    formatted_context = format_docs(retrieved_docs)
    formatted_prompt = f"Question: {question}\n\nContext:
    response = ollama.chat(model='mistral', messages=[{'rol
    return response['message']['content']
# Gradio interface
iface = gr.Interface(
    fn=rag_chain,
    inputs=["text", "text"],
    outputs="text",
    title="RAG Chain Question Answering",
    description="Enter a URL and a query to get answers from
)
# Launch the app
iface.launch()
```

Output

RAG Chain Question Answering

Enter a URL and a query to get answers from the RAG chain.

