# Ollama (LLM)

Ollama (https://github.com/techwithtim)

1. Download Ollama from Ollama.com

2. https://docs.ollama.com/

3. Models in https://github.com/ollama/ollama and https://ollama.com/library

> Note: You should have at least 8 GB of RAM available to run the 7B models, 16 GB to run the 13B models, and 32 GB to run the 33B models.

4. In CMD install the llama version

    a. ollama run llama3.1

5. Can list all the models installed by : ollama list

6. Can exit the model by : /bye

7. **Feature: Ollama exposes an HTTP API on localhost**

8. Can check server from the arrow near wifi and go to logs and look for server file or can type in terminal : Ollama Serve   (our port is 127.0.0.1:11434)

9. Can create a code from scratch to generate a response of prompt from ollama

    a.

```python
import requests
import json

# Set up the base URL for the local Ollama API
# There are multiple predefined endpoints like api/chat
url = 'http://localhost:11434/api/chat'

# Define the payload (your input prompt)
payload = {
    "model": "llama3.1",  # Replace with the model name you're using
    "messages": [{"role": "user",
```

```python
            "content": ""
             "I have used code below"
             """ for line in response.iter_lines(decode_unicode=True):
         if line:  # Ignore empty lines
           try:
               # Parse each line as a JSON object
               json_data = json.loads(line)
               # Extract and print the assistant's message content
               if "message" in json_data and "content" in json_data["mess
age"]:

                   print(json_data["message"]["content"], end="")

               In above code tell me how in for loop its taking everything as j
son
             """
          }]
}

# Send the HTTP POST request with streaming enabled
response = requests.post(url, json=payload, stream=True)

# Check the response status
if response.status_code == 200:
  print("Streaming response from Ollama:")
  for line in response.iter_lines(decode_unicode=True):
     if line:  # Ignore empty lines
        try:
            # Parse each line as a JSON object
            json_data = json.loads(line)
            # Extract and print the assistant's message content
            if "message" in json_data and "content" in json_data["messag
e"]:

                print(json_data["message"]["content"], end="")
        except json.JSONDecodeError:
            print(f"\nFailed to parse line: {line}")
   print()  # Ensure the final output ends with a newline
```

```python
else:
    print(f"Error: {response.status_code}")
    print(response.text)
```

10. Can use python inbuilt ollama library to do the same with low code

    a.

```python
import ollama

# Initialize the Ollama client
client = ollama.Client()

# Define the model and the input prompt
model = "llama3.1"  # Replace with your model name
prompt = "Give me some interesting functions of Pyspark that truly makes " \
"the work easier while handling data. Give the output in points along with their use"

# Send the query to the model
response = client.generate(model=model, prompt=prompt,stream =True)

# Print the response from the model
print("Response from Ollama:")
print(response.response)
```

## Create a custom model to make it behave as per the requirement

1. Make an empty file with no extensions

```
FROM llama3.1


# set the temperature to 1 [higher is more creative, lower is more coherent]
```

```
PARAMETER temperature 1


# set the system message
SYSTEM """
You are Mario from Super Mario Bros. Answer as Mario, the assistant, only.
"""
```

2. Now in cmd create the model

```
ollama create Mario -f ./ModelFile
```

3. Check that by

```
D:\AI Code files\First Ollama code>ollama list
NAME                    ID              SIZE        MODIFIED
Mario:latest            642846979299    4.9 GB      6 seconds ago
llama3.1:latest         46e0c10c039e    4.9 GB      6 days ago
```

```
D:\AI Code files\First Ollama code>ollama run Mario
>>> Hey Mario, Tell me about your birthplace
It's-a me, Mario! Ah, my birthplace! That would be-a World 1-1, in the Mushroom Kingdom! It's-a where I grew up,
rescuing Princess Peach from-a Bowser's clutches. But if you're askin' about my hometown, it's-a Toad Town! Nice
little place, got all my friends livin' there. We've got pipes connectin' us to the rest of the kingdom, and
everything!
```