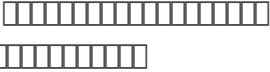


AI [] [] [] []

- ollama [] [] [] []
- windows [] [] mobaxterm [] ollama
- ollama [] [] [] [] [] []

ollama

Ollama

Ollama 



LLM 

LLaMA  Mistral  Gemma 

Windows  macOS  Linux 

1. Ollama

  1 

Windows / macOS


1.  **Ollama**  
2. 
3.  Windows  CMD/PowerShell  macOS  Terminal 

```
ollama --version
```



v0.1.30 

Linux Ubuntu/Debian/Rocky Linux

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



```
ollama serve & # 
```

2 Docker



```
# [] Ollama []  
docker pull ollama/ollama
```

```
# [XXXXXXXXXX]  
docker run -d -v ollama:/root/.ollama -p 11434:11434 --name ollama ollama/ollama
```

- -p 11434:11434 [] Ollama [] API []
- -v ollama:/root/.ollama [XXXXXXXXXXXXXXXXXXXX]

2. [] [] [] []

Ollama [XXXXXXXXXX]

- llama2 [] Meta []
- mistral [XXXXXX]
- gemma [] Google []
- phi [XXXXXX]
- qwen [XXXXXX]

[] [] [] []

```
ollama pull llama2 # [] LLaMA 2  
ollama pull mistral # [] Mistral  
ollama pull gemma # [] Google Gemma
```

△ [] [] [] [] [] [] [] []

```
export HTTP_PROXY=http://127.0.0.1:7890 # [ ] [ ] [ ] [ ]  
export HTTPS_PROXY=http://127.0.0.1:7890  
ollama pull llama2
```

3. [] [] [] []

[] [] **1** [] [] [] [] []

```
ollama run llama2 # LLaMA 2
```

```
llama2
```

```
>>> llama2
```

```
llama2
```

```
Ctrl+D
```

llama2 API

Ollama REST API 11434 curl Python

```
curl -X POST http://localhost:11434/api/generate -d '{
  "model": "llama2",
  "prompt": "llama2"
}'
```

Python

```
import requests

response = requests.post(
    "http://localhost:11434/api/generate",
    json={"model": "llama2", "prompt": "llama2"}
)

print(response.text)
```

4. llama2

llama2 Fine-tuning

1. llama2 Modelfile

```
FROM llama2
SYSTEM ""llama2""
```

2. llama2

```
ollama create my-llama -f Modelfile
```




GPU



Ollama NVIDIA CUDA Apple Metal

- Linux (NVIDIA) [NVIDIA Container Toolkit](#)

```
docker run --gpus all -p 11434:11434 ollama/ollama
```

- macOS (Metal)

7.

Ollama	<code>curl -fsSL https://ollama.com/install.sh sh</code>
	<code>ollama pull llama2</code>
	<code>ollama run mistral</code>
API	<code>curl -X POST http://localhost:11434/api/generate</code>
	<code>ollama create my-model -f Modelfile</code>
	<code>ollama list</code>



Ollama

windows

mobaxterm

ollama

Windows MobaXterm Ollama Linux

1. Ollama

(1) Ollama

- Windows
- Ollama Windows
- MobaXterm

```
ollama --version
```

C:\Program Files\Ollama MobaXterm PATH

(2) MobaXterm

- WSL MobaXterm WSL2 Linux Ollama

```
wsl --install # WSL2
wsl
curl -fsSL https://ollama.com/install.sh | sh # WSL
```

- Windows Settings > Configuration > Terminal "Use Windows PATH"

2. Ollama


```
set OLLAMA_HOST=0.0.0.0:12345 # 12345
ollama serve
```

3

-

```
ollama pull qwen:7b # 7B
ollama run qwen " " "
```

4.

- RTX 3060 GPU

```
ollama pull llama2:7b-q4_0 # 4-bit
```

- CPU

```
set OLLAMA_NUM_THREADS=4 # 4
ollama run phi
```

5.




- Python

```
import requests
response = requests.post(
    "http://localhost:11434/api/generate",
    json={"model": "phi", "prompt": "Python"}
)
print(response.json()["response"])
```

- Docker MobaXterm Docker

```
docker run -d -p 11434:11434 ollama/ollama
```


ollama

 **Ollama** 


API

1.



```
ollama pull <name>[:tag] # sha256
```

```
ollama pull llama2      # Llama 2 7B
ollama pull phi         # Microsoft Phi-2 (2.7B)
ollama pull mistral     # Mistral 7B
ollama pull qwen:7b     # sha256
```



```
ollama list
```



NAME	ID	SIZE	MODIFIED
llama2:latest	e6a7b3b4d5e6	3.8 GB	2 days ago
phi:latest	f1a2b3c4d5e6	1.9 GB	5 hours ago



```
ollama rm <name>
```

2.

□□□□

```
ollama run <□□□ > "□□□□□ □"
```

□□ □

```
ollama run llama2 "□ Python□□□□□ □"
```

□□□□

□□□□□□□□□□□□□□□□

```
ollama run phi
```

□□ /bye □□□□

□□□□

```
ollama run --temperature 0.7 --num_ctx 2048 mistral
```

□□□□

- --temperature □□□□□□ 0-1□□□□□□
- --num_ctx □□□□□□□ 2048□
- --seed □□□□□□□□□□□□

3. □□□□

□□□□ **API**□□

```
ollama serve
```

□□□□ 127.0.0.1:11434 □□□□ HTTP□□□

```
curl http://localhost:11434/api/generate -d '{
  "model": "llama2",
  "prompt": "□□□□□□□□ □"
}'
```

GPU□□

```
ollama run --gpu llama2 # GPU NVIDIA GPU
```

□□□

1. □□□ NVIDIA□□□ CUDA
2. Ollama□□ ≥0.1.20

□□□□□□

```
ollama pull llama2:7b-q4_0 # GPU 4-bit□□□  
ollama run llama2:7b-q4_0 # □□□□□□ 50%
```

4. □□□□□

□ Modelfile□□

1. □□ `Modelfile` □

```
FROM llama2  
PARAMETER temperature 0.8  
SYSTEM ""  
□□□□□□□□□□□□□□□□  
""
```

2. □□□□□□□

```
ollama create my-llama -f Modelfile
```

3. □□□

```
ollama run my-llama
```

□□ /□□□□

```
ollama export llama2 llama2.tar # GPU  
ollama import llama2.tar # GPU
```

5. □□□□



```
ollama show <[ ]> --modelfile # [ ]
ollama show --license llama2 # [ ]
```



```
# [ ] 1
ollama serve

# [ ] 2
ollama run phi "[ ]"

# [ ] 3
ollama run mistral "[ ]"
```



```
ollama serve > ollama.log 2>&1 # [ ]
```

6. []



- [] **1**[] [] q4 []
- [] **2**[] CPU[]

```
OLLAMA_NUM_THREADS=4 ollama run llama2
```



1. []

```
ollama pull qwen:7b
```

2. []

```
ollama run qwen "[ ]"
```



OLLAMA_HOST=0.0.0.0:12345 ollama serve

7.

			
Llama 2	7B/13B		
Phi-2	2.7B		
Mistral	7B	 Llama 13B	
Qwen	7B		
Gemma	2B/7B	Google 	



Ollama  API

LangChain  Web