

ollama

Ollama

Ollama   **LLM**  LLaMA  Mistral  Gemma 
 **Windows**  **macOS**  **Linux** 

1. Ollama

  **1** 



Windows / macOS

-  **Ollama**  
- 
-  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
```

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

- -p 11434:11434 ollama API
- -v ollama:/root/.ollama

2. ollama

Ollama

- llama2 Meta
- mistral
- gemma Google Gemma
- phi
- qwen



```
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
```

[] [] [] [] [] []

```
>>> [ ] [ ] [ ] [ ] [ ] [ ] [ ] [ ] [ ] [ ]
```

[] [] [] [] [] [] [] [] Ctrl+D [] []

[] [] 2 [] **API** [] []

Ollama [] REST API [] [] [] [] [] 11434 [] [] [] curl [] Python [] []

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

Python [] [] []

```
import requests

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

print(response.text)
```

4. [] [] [] []

[] [] [] [] [] [] [] [] [] [] Fine-tuning [] []

1. `Modelfile`

```
FROM llama2
SYSTEM ""
```

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

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

3. `ollama run my-llama`

```
ollama run my-llama
```

5. `ollama`

`ollama`

- **Windows** `PATH`
- **Linux/macOS**

```
export PATH=$PATH:/usr/local/bin
```

`ollama` `ollama` **2**

- `1`
- `2` `~/.ollama/models` Linux/macOS `C:\Users\<`
`>\.ollama\models` Windows



`ollama` `ollama` **3**

`11434`

```
OLLAMA_HOST=0.0.0.0:11435 ollama serve # 11435
```

6.



```
ollama list #   
ollama rm llama2 # 
```

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  

Revision #2

Created 3 April 2025 14:22:23 by Admin

Updated 9 April 2025 03:10:28 by Admin