




ollama

 **Ollama** 


API

1.



```
ollama pull <model>[:tag] # [digest]
```

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



```
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 <model> >
```

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_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 	



LangChain  Web 

Ollama  API 

Revision #1

Created 9 April 2025 09:03:12 by Admin

Updated 9 April 2025 09:04:05 by Admin