

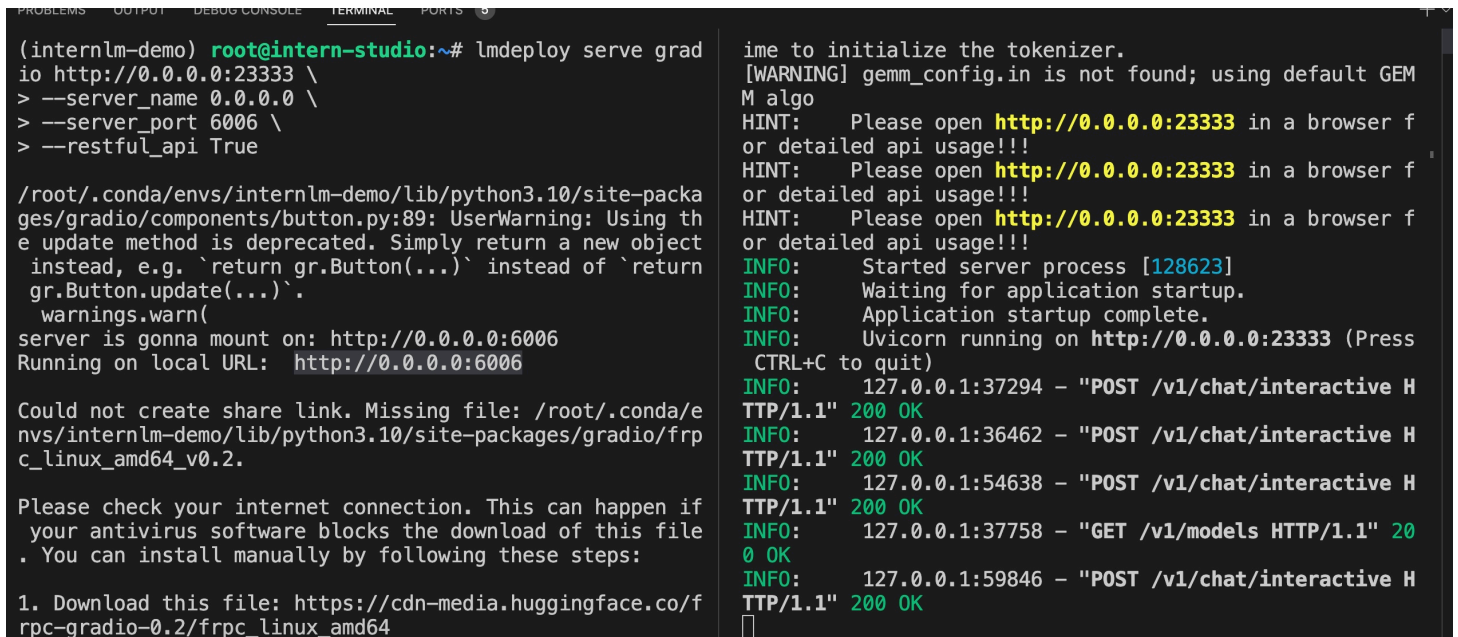
# 基础作业：

使用 LMDeploy 以本地对话、网页Gradio、API服务中的一种方式部署 InternLM-Chat-7B 模型，生成 300 字的小故事（需截图）

## 本地对话 (TurboMind推理+API服务)

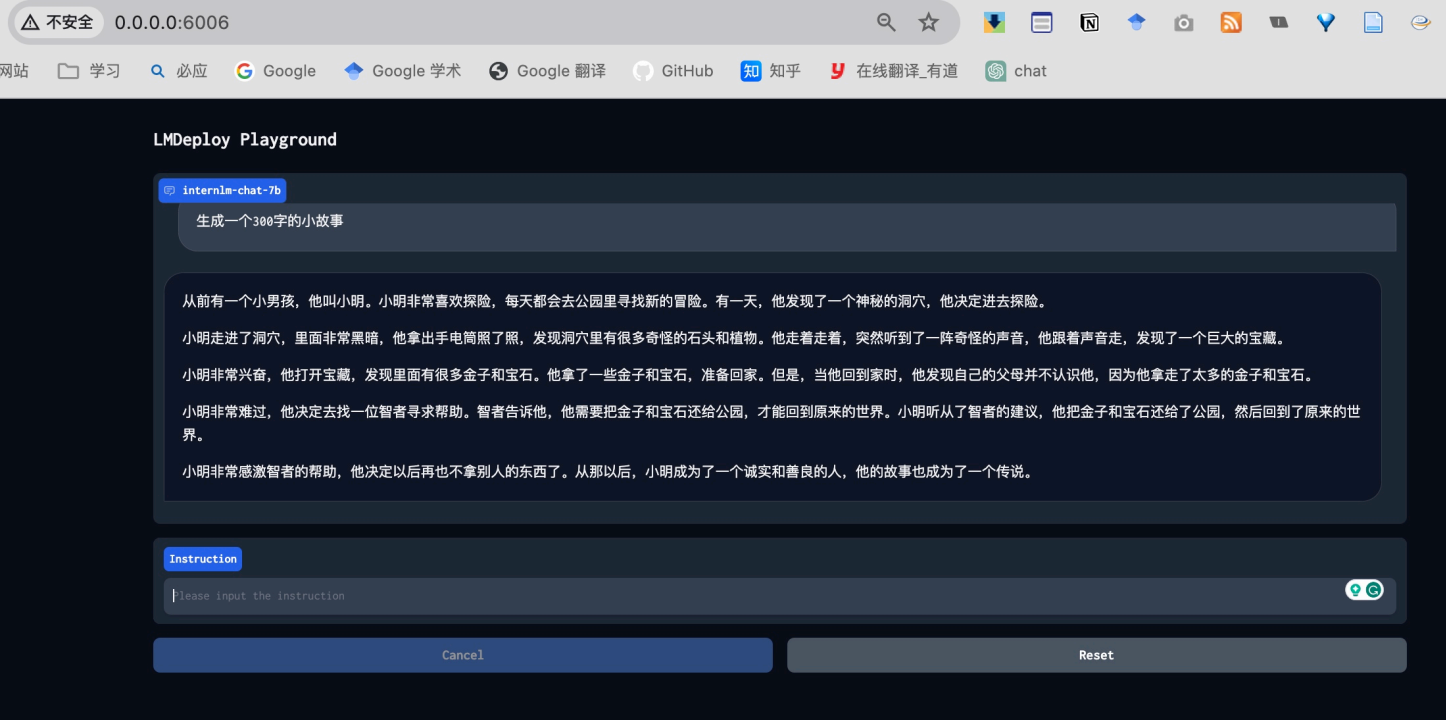


## TurboMind 服务作为后端 + gradio





## TurboMind 推理作为后端 + gradio



## 进阶作业（可选做）

无量化 显存14G

CPU 6.98%GPU-1: Nvidia A100(1/4) 0%内存 2.63 / 56 GB 4.7%显存 14662 / 20470 MiB 71.63%

root

infer\_compare.pyconfig.ini ~/workspace/triton\_models/...config.ini ~/workspace\_quant/...

workspace\_quant > triton\_models > weights > config.ini

```
10 attn_bias = 1
11 start_id = 1
12 end_id = 2
13 session_len = 512
14 weight_type = int4
15 rotary_embedding = 128
16 rope_theta = 10000.0
17 size_per_head = 128
18 group_size = 128
19 max_batch_size = 1
20 max_context_token_num = 1
21 step_length = 1
22 cache_max_entry_count = 0.5
23 cache_block_seq_len = 128
24 cache_chunk_size = 1
25 use_context_fmha = 1
26 quant_policy = 0
27 max_position_embeddings = 2048
28 rope_scaling_factor = 0.0
```

PROBLEMSOUTPUTDEBUG CONSOLETERMINALPORTS 5

double enter to end input >>  
> hi  
  
<|User|>:hi  
<|Bot|>: 书生·浦语：你好，有什么我可以帮助你的吗？  
(internlm-demo) root@intern-  
studio:~# lmdeploy serve api  
\_server ./workspace \  
> --server\_name 0.0.0.0 \  
> --server\_port 23333 \  
> --instance\_num 1 \  
> --tp 1  
  
model\_source: workspace

o -studio:~# lmdeploy serve a  
pi\_client http://localhost:  
23333  
  
double enter to end input >  
>> 书生·浦语：好的，我会尽  
力帮助您。请问您有什么需要  
我帮忙的吗？  
double enter to end input >  
>> nihao  
  
书生·浦语：你好！有什么我可  
以帮助你的吗？  
double enter to end input >  
>> □

py...  
py...  
bash

W4A16量化



root



infer\_compare.py

config.ini ~/workspace/triton\_models/...

config.ini ~/workspace\_quant/...



workspace\_quant > triton\_models > weights > config.ini

```
1  [llama]
2  model_name = internlm-chat-7b
3  tensor_para_size = 1
4  head_num = 32
5  kv_head_num = 32
6  vocab_size = 103168
7  num_layer = 32
8  inter_size = 11008
9  norm_eps = 1e-06
10 attn_bias = 1
11 start_id = 1
12 end_id = 2
13 session_len = 2056
14 weight_type = int4
15 rotary_embedding = 128
16 rope_theta = 10000.0
17 size_per_head = 128
18 group_size = 128
19 max_batch_size = 64
```



PROBLEMS

OUTPUT

DEBUG CONSOLE

TERMINAL

PORTS

4



python



developed by Shanghai AI Laboratory (上海人工智能实验室). It is designed to be helpful, honest, and harmless.

– InternLM (书生·浦语) can understand and communicate fluently in the language chosen by the user such as English and 中文.

<|User|>:

<|Bot|>: I'm sorry, I'm not sure what you're asking for. Could you please provide more context or information?

double enter to end input >>> hello

<|User|>:hello

<|Bot|>: Hello! How can I assist you today?

double enter to end input >>>