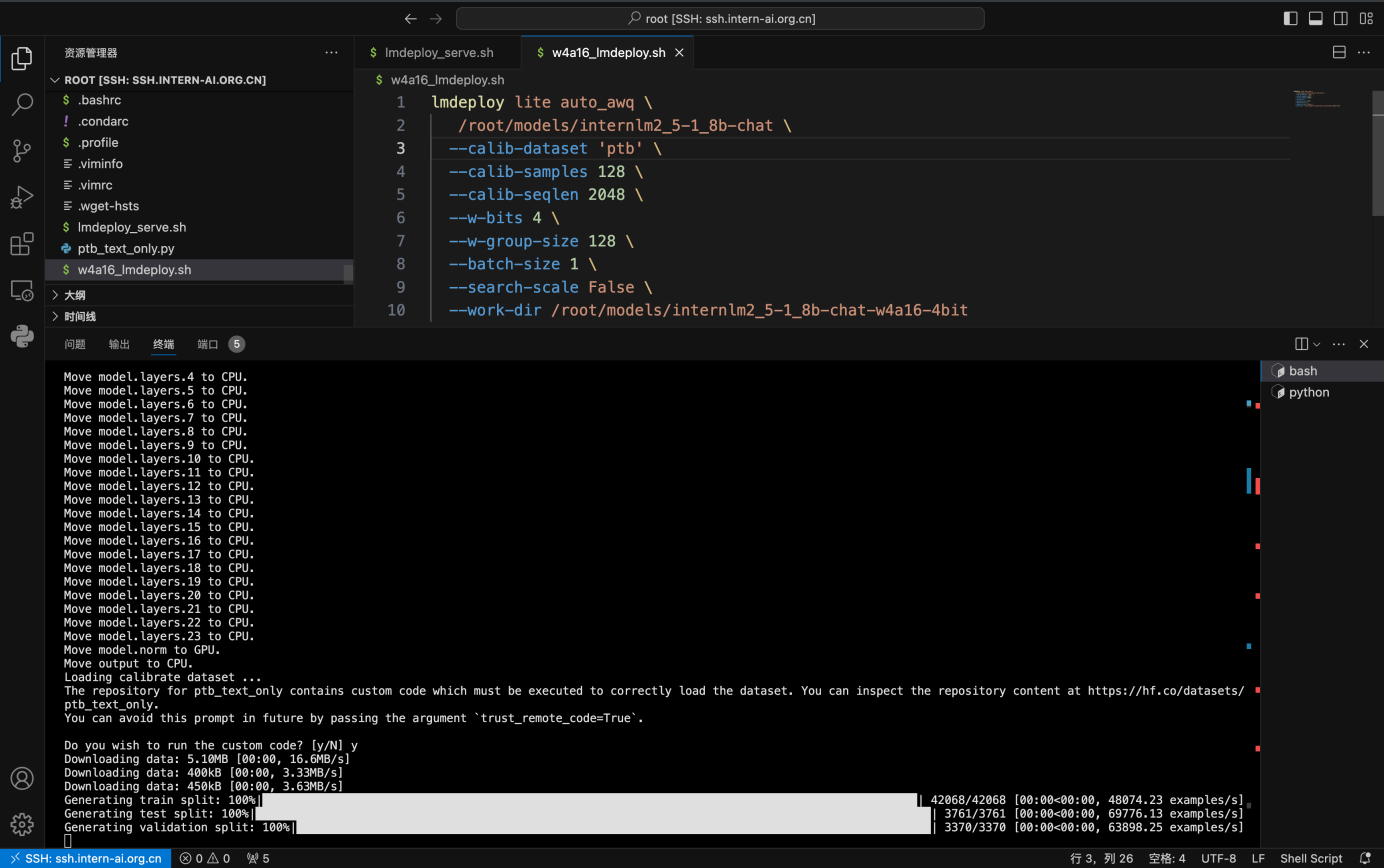


# 量化internlm2\_5-1\_8b-chat



```
root [SSH: ssh.intern-ai.org.cn]
lmdeploy lite auto_awq \

model.layers.15.feed_forward.w2 weight packed.
model.layers.16.attention.wqkv weight packed.
model.layers.16.attention.wo weight packed.
model.layers.16.feed_forward.w1 weight packed.
model.layers.16.feed_forward.w3 weight packed.
model.layers.16.feed_forward.w2 weight packed.
model.layers.17.attention.wqkv weight packed.
model.layers.17.attention.wo weight packed.
model.layers.17.feed_forward.w1 weight packed.
model.layers.17.feed_forward.w3 weight packed.
model.layers.17.feed_forward.w2 weight packed.
model.layers.18.attention.wqkv weight packed.
model.layers.18.attention.wo weight packed.
model.layers.18.feed_forward.w1 weight packed.
model.layers.18.feed_forward.w3 weight packed.
model.layers.18.feed_forward.w2 weight packed.
model.layers.19.attention.wqkv weight packed.
model.layers.19.attention.wo weight packed.
model.layers.19.feed_forward.w1 weight packed.
model.layers.19.feed_forward.w3 weight packed.
model.layers.19.feed_forward.w2 weight packed.
model.layers.20.attention.wqkv weight packed.
model.layers.20.attention.wo weight packed.
model.layers.20.feed_forward.w1 weight packed.
model.layers.20.feed_forward.w3 weight packed.
model.layers.20.feed_forward.w2 weight packed.
model.layers.21.attention.wqkv weight packed.
model.layers.21.attention.wo weight packed.
model.layers.21.feed_forward.w1 weight packed.
model.layers.21.feed_forward.w3 weight packed.
model.layers.21.feed_forward.w2 weight packed.
model.layers.22.attention.wqkv weight packed.
model.layers.22.attention.wo weight packed.
model.layers.22.feed_forward.w1 weight packed.
model.layers.22.feed_forward.w3 weight packed.
model.layers.22.feed_forward.w2 weight packed.
model.layers.23.attention.wqkv weight packed.
model.layers.23.attention.wo weight packed.
model.layers.23.feed_forward.w1 weight packed.
model.layers.23.feed_forward.w3 weight packed.
model.layers.23.feed_forward.w2 weight packed.
(agent) root@intern-studio-5021523:~#
(agent) root@intern-studio-5021523:~# cd models
(agent) root@intern-studio-5021523:~/models# ls
InternVL2-26B  internlm2_5-1_8b-chat  internlm2_5-1_8b-chat-w4a16-4bit  internlm2_5-7b-chat
(agent) root@intern-studio-5021523:~/models# du -sh *
0      InternVL2-26B
0      internlm2_5-1_8b-chat
1.5G   internlm2_5-1_8b-chat-w4a16-4bit
0      internlm2_5-7b-chat
(agent) root@intern-studio-5021523:~/models#
```

# 量化模型API对话

量化后的模型

资源管理器

ROOT [SSH: SSH.INTERN-AI.ORG.CN]

.bash\_history

.bashrc

! .condarc

.profile

.viminfo

.vimrc

大纲

时间线

lmdeploy\_serve.sh

w4a16\_lmdeploy.sh

1

2

3

4

5

6

7

8

lmdeploy lite auto\_awq \

/root/models/internlm2\_5-1\_8b-chat \

--calib-dataset 'ptb' \

--calib-samples 128 \

--calib-seqlen 2048 \

--w-bits 4 \

--w-group-size 128 \

batch\_size 1 \

问题

输出

终端

端口

python models

python

root@intern-studio-5021523:~/models# lmdeploy chat /root/models/internlm2\_5-1\_8b-chat-w4a16-4bit/ --model-format awq

chat\_template\_config:

ChatTemplateConfig(model\_name='internlm2', system=None, meta\_instruction=None, eosys=None, user=None, eoh=None, assistant=None, eoa=None, separator=None, capability='chat',

stop\_words=None)

engine\_cfg:

TurbomindEngineConfig(model\_name='/root/models/internlm2\_5-1\_8b-chat-w4a16-4bit/', model\_format='awq', tp=1, session\_len=32768, max\_batch\_size=1, cache\_max\_entry\_count=0.8,

cache\_block\_seq\_len=64, enable\_prefix\_caching=False, quant\_policy=0, rope\_scaling\_factor=0.0, use\_logn\_attn=False, download\_dir=None, revision=None, max\_prefill\_token\_num=81

92, num\_tokens\_per\_iter=0, max\_prefill\_iters=1)

[WARNING] gemm\_config.in is not found; using default GEMM algo

double enter to end input >>> hello

<|im\_start|>system

You are an AI assistant whose name is InternLM (书生·浦语).

- InternLM (书生·浦语) is a conversational language model that is developed by Shanghai AI Laboratory (上海人工智能实验室). It is designed to be helpful, honest, and harmles

s.

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

<|im\_end|>

<|im\_start|>user

hello<|im\_end|>

<|im\_start|>assistant

你好，很高兴见到你。有什么我能帮助你的吗？

double enter to end input >>> 你是谁

<|im\_start|>user

你是谁<|im\_end|>

<|im\_start|>assistant

我是由上海人工智能实验室开发的对话式语言模型，名为 InternLM。我致力于通过自然语言交互为人类提供帮助和解答问题。如果你有任何问题或需要帮助，请随时告诉我。

double enter to end input >>> 1+100=?

<|im\_start|>user

1+100=?<|im\_end|>

<|im\_start|>assistant

好的，我会帮你计算一下。1加100等于101。请问还有其他需要帮助的吗？

double enter to end input >>>

SSH: ssh.intern-ai.org.cn

行 3, 列 26 空格: 4 UTF-8 LF Shell Script

## 量化前的模型

资源管理器

ROOT [SSH: SSH.INTERN-AI.ORG.CN]

.bash\_history

.bashrc

! .condarc

.profile

.viminfo

.vimrc

大纲

时间线

lmdeploy\_serve.sh

w4a16\_lmdeploy.sh

1

2

3

4

5

6

7

8

lmdeploy lite auto\_awq \

/root/models/internlm2\_5-1\_8b-chat \

--calib-dataset 'ptb' \

--calib-samples 128 \

--calib-seqlen 2048 \

--w-bits 4 \

--w-group-size 128 \

batch\_size 1 \

问题

输出

终端

端口

python models

python

root@intern-studio-5021523:~/models# lmdeploy chat /root/models/internlm2\_5-1\_8b-chat/

chat\_template\_config:

ChatTemplateConfig(model\_name='internlm2', system=None, meta\_instruction=None, eosys=None, user=None, eoh=None, assistant=None, eoa=None, separator=None, capability='chat',

stop\_words=None)

engine\_cfg:

TurbomindEngineConfig(model\_name='/root/models/internlm2\_5-1\_8b-chat/', model\_format=None, tp=1, session\_len=32768, max\_batch\_size=1, cache\_max\_entry\_count=0.8, cache\_block

seq\_len=64, enable\_prefix\_caching=False, quant\_policy=0, rope\_scaling\_factor=0.0, use\_logn\_attn=False, download\_dir=None, revision=None, max\_prefill\_token\_num=8192, num\_toks

ns\_per\_iter=0, max\_prefill\_iters=1)

[WARNING] gemm\_config.in is not found; using default GEMM algo

double enter to end input >>> hello

<|im\_start|>system

You are an AI assistant whose name is InternLM (书生·浦语).

- InternLM (书生·浦语) is a conversational language model that is developed by Shanghai AI Laboratory (上海人工智能实验室). It is designed to be helpful, honest, and harmles

s.

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

<|im\_end|>

<|im\_start|>user

hello<|im\_end|>

<|im\_start|>assistant

Hello! How can I assist you today?

double enter to end input >>> 你是谁？

<|im\_start|>user

你是谁？<|im\_end|>

<|im\_start|>assistant

我是书生·浦语，一个由上海人工智能实验室开发的语言模型，旨在通过智能和人性化的交互，为您提供帮助和娱乐。我可以通过回答问题、提供定义和解释、将文本从一种语言翻译成另一种语言、

总结文本、生成文本、编写故事、分析情感、提供推荐、开发算法、编写代码以及其他任何基于语言的任务来帮助您。我很乐意帮助您解答问题，并提供有趣和有用的信息。

double enter to end input >>> 1+101=?

<|im\_start|>user

1+101=?<|im\_end|>

<|im\_start|>assistant

1+101=102

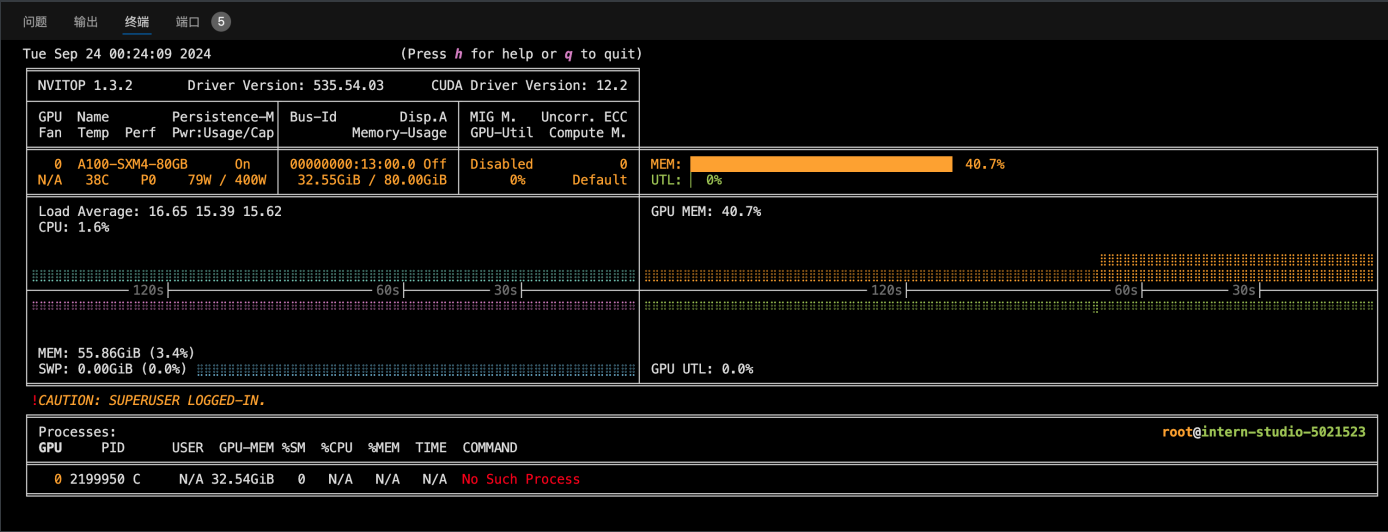
double enter to end input >>>

SSH: ssh.intern-ai.org.cn

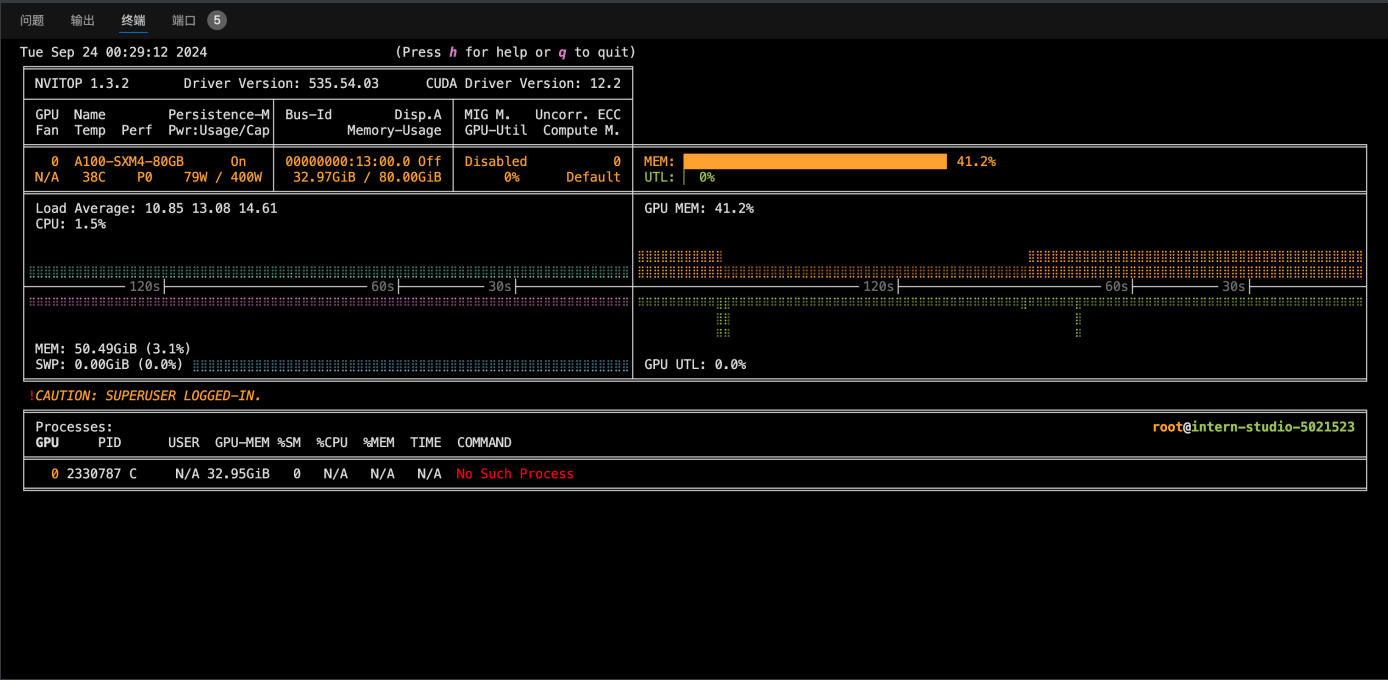
行 8, 列 19 空格: 4 UTF-8 LF Shell Script

# 50%A100运行现存占用

量化后的模型



量化前的模型



Function call

资源管理器

ROOT [SSH: SSH.INTERN-AI.ORG.CN]

models

internlm2\_5-1.8b-chat

internlm2\_5-1.8b-chat-w4a16-4bit

internlm2\_5-7b-chat

InternVL2-26B

nltk\_data

openccompass

share

workspace

.aide\_drive\_storage\_lock

.aide\_notify\_send\_lock

.aide\_usage\_lock

.bash\_history

.bashrc

.condarc

.profile

.viminfo

.vimrc

.wget-hsts

internlm2\_5\_func.py

internlm2\_5.py

lmdeploy\_serve.sh

ptb\_text\_only.py

大纲

时间线

internlm2\_5\_func.py

```
71
72 messages.append({
73     'role': 'assistant',
74     'content': response.choices[0].message.content
75 })
76 messages.append({
77     'role': 'environment',
78     'content': f'3+5={func1_out}',
79     'name': 'plugin'
80 })
81 response = client.chat.completions.create(
82     model=model_name,
83     messages=messages,
84     temperature=0.8,
85     top_p=0.8,
86     stream=False,
87     tools=tools)
88 print(response)
89 func2_name = response.choices[0].message.tool_calls[0].function.name
90 func2_args = response.choices[0].message.tool_calls[0].function.arguments
91 func2_out = eval(f'{func2_name}({func2_args})')
92 print(func2_out)
93
```

问题

输出

终端

端口

● (agent) root@intern-studio-5021523:~# python internlm2\_5\_func.py

ChatCompletion(id='4', choices=[Choice(finish\_reason='tool\_calls', index=0, logprobs=None, message=ChatCompletionMessage(content='I will call the API to calculate the result.', refusal=None, role='assistant', function\_call=None, tool\_calls=[ChatCompletionMessageToolCall(id='0', function=Function(arguments='{\"a\": 3, \"b\": 5}', name='add'), type='function'))]), created=1727110235, model='/root/models/internlm2\_5-7b-chat', object='chat.completion', service\_tier=None, system\_fingerprint=None, usage=CompletionUsage(completion\_tokens=35, prompt\_tokens=263, total\_tokens=298))

8

ChatCompletion(id='5', choices=[Choice(finish\_reason='tool\_calls', index=0, logprobs=None, message=ChatCompletionMessage(content='The result of the first calculation is 8. Now I will call the API again to calculate the final result.', refusal=None, role='assistant', function\_call=None, tool\_calls=[ChatCompletionMessageToolCall(id='1', function=Function(arguments='{\"a\": 8, \"b\": 2}', name='mul'), type='function'))]), created=1727110235, model='/root/models/internlm2\_5-7b-chat', object='chat.completion', service\_tier=None, system\_fingerprint=None, usage=CompletionUsage(completion\_tokens=48, prompt\_tokens=292, total\_tokens=340))

16

[{'role': 'user', 'content': 'Compute (3+5)\*2'}, {'role': 'assistant', 'content': 'I will call the API to calculate the result.'}, {'role': 'environment', 'content': '3+5=8', 'name': 'plugin'}]

● (agent) root@intern-studio-5021523:~#

bash

python

bash

SSH: ssh.intern-ai.org.cn 0 0 5 行 67, 列 1 空格: 4 UTF-8 LF Python