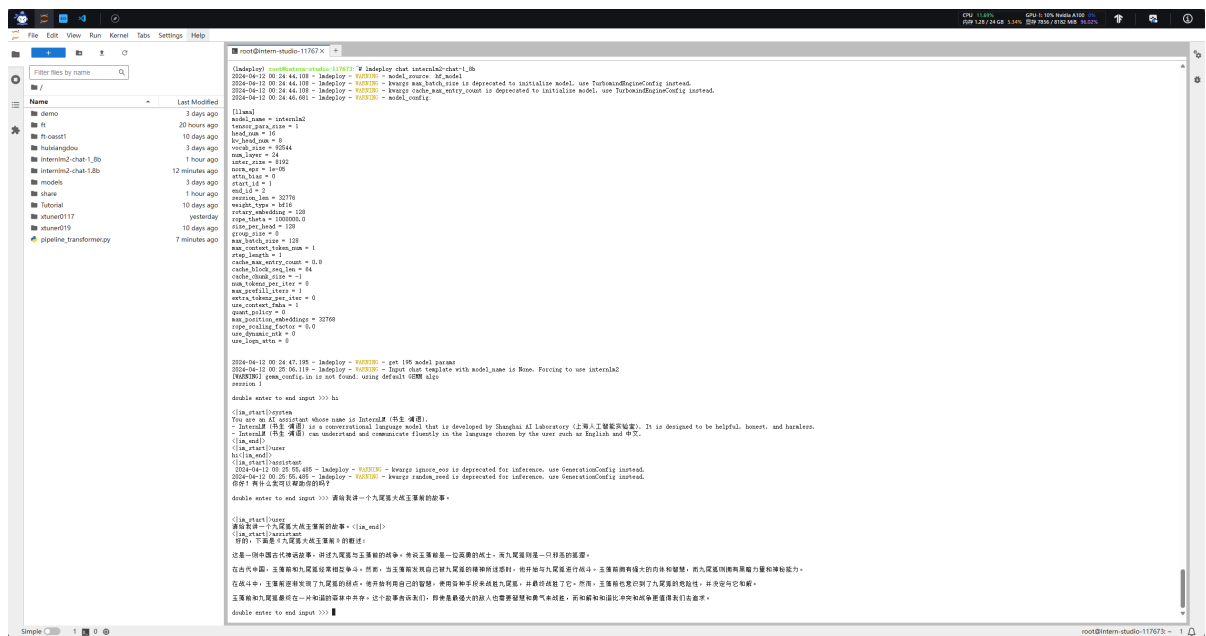


作业：LMDeploy量化部署

1、使用LMDeploy与模型对话



2、LMDeploy模型量化(lite)

```
(loadpolicy) root@ntware-studio-1187872: #
(loadpolicy) root@ntware-studio-1187872: # loadpolicy lite ntwc_sq_w \
  /root/internaln2-chst-1_0b-0bit
> --calib-dataset ptb \
> --calib-samples 128 \
> --calib-seqlen 1024 \
> --bits 4 \
> --group-size 128 \
> --work-dir /root/internaln2-chst-1_0b-0bit

loading checkpoint shards: 100%
Move model.vok embeddings to GPU.
Move model.layers.0 to CPU.
Move model.layers.1 to CPU.
Move model.layers.2 to CPU.
Move model.layers.3 to CPU.
Move model.layers.4 to CPU.
Move model.layers.5 to CPU.
Move model.layers.6 to CPU.
Move model.layers.7 to CPU.
Move model.layers.8 to CPU.
Move model.layers.9 to CPU.
Move model.layers.10 to CPU.
Move model.layers.11 to CPU.
Move model.layers.12 to CPU.
Move model.layers.13 to CPU.
Move model.layers.14 to CPU.
Move model.layers.15 to CPU.
Move model.layers.16 to CPU.
Move model.layers.17 to CPU.
Move model.layers.18 to CPU.
Move model.layers.19 to CPU.
Move model.layers.20 to CPU.
Move model.layers.21 to CPU.
Move model.layers.22 to CPU.
Move model.layers.23 to CPU.
Move model.norm to CPU.
Move output to GPU.

Loading calibrate dataset ...
root/.conda/envs/loadpolicy/lib/python3.10/site-packages/datasets/load.py:1461: FutureWarning: The repository for ptb_text_only contains custom code which must be executed to correctly load the dataset. You can inspect the repository content at https://hf.co/datasets/ptb_text_only
You can avoid this message in future by passing the argument 'trust_remote_code=True'.
Passing 'trust_remote_code=True' will be mandatory to load this dataset from the next major release of 'datasets'.
warning: warn()
Downloading builder script: 6.50MB [00:00, 26.3MB/s]
Downloading readme: 4.21kB [00:00, 19.3MB/s]
Downloading data: 5.10MB [00:00, 17.0MB/s]
Downloading data: 409B [00:00, 3.77MB/s]
Downloading data: 459B [00:00, 3.68MB/s]
Generating train split: 100%
Generating test split: 100%
Computing validation split: 100%
model.layers.0, sampler: 128, max gpu memory: 2.25 GB
model.layers.1, sampler: 128, max gpu memory: 2.75 GB
model.layers.2, sampler: 128, max gpu memory: 2.75 GB
model.layers.3, sampler: 128, max gpu memory: 2.75 GB
model.layers.4, sampler: 128, max gpu memory: 2.75 GB
model.layers.5, sampler: 128, max gpu memory: 2.75 GB
model.layers.6, sampler: 128, max gpu memory: 2.75 GB
model.layers.7, sampler: 128, max gpu memory: 2.75 GB
model.layers.8, sampler: 128, max gpu memory: 2.75 GB
model.layers.9, sampler: 128, max gpu memory: 2.75 GB
model.layers.10, sampler: 128, max gpu memory: 2.75 GB
model.layers.11, sampler: 128, max gpu memory: 2.75 GB
[2/2] [00:31:00.00, 15.63s/1]
```

3、LMDeploy服务(serve)

```
(lmdeploy) root@intern-studio-117673: # lmdeploy serve api_server \
> /root/internlm2-chat-1_8b \
> --model-format hf \
> --quant-policy 0 \
> --server-name 0.0.0.0 \
> --server-port 23333 \
> --tp 1

[WARNING] gemm_config.in is not found; using default GEMM algo
HINT: Please open http://0.0.0.0:23333 in a browser for detailed api usage!!!
HINT: Please open http://0.0.0.0:23333 in a browser for detailed api usage!!!
HINT: Please open http://0.0.0.0:23333 in a browser for detailed api usage!!!
INFO: Started server process [128982]
INFO: Waiting for application startup.
INFO: Application startup complete.
INFO: Uvicorn running on http://0.0.0.0:23333 (Press CTRL+C to quit)
INFO: 127.0.0.1:59960 - "GET / HTTP/1.1" 200 OK
INFO: 127.0.0.1:59960 - "GET /openapi.json HTTP/1.1" 200 OK
```

4、使用LMDeploy运行视觉多模态大模型

```
(lddeploy) root@intern-studio-117673:~# python pipeline_llava.py
[WARNING] gpt4_config.in is not found, using default GPT4 also
You are using a model of type llava to instantiate a model of type llava_llama. This is not supported for all configurations of models and can yield errors.
You are using a model of type llava to instantiate a model of type llava_llama. This is not supported for all configurations of models and can yield errors.
preprocessor_config.json: 100% | 316/316 [00:00:00.00, 2.01MB/s]
config.json: 4.76kB [00:00:18.00MB/s]
processor_model.bin: 100% | 1.11G/1.11G [00:33:00.00, 61.0MB/s]
Loading checkpoint shards: 100% | 3/3 [00:00:00.00, 4.56it/s]
Response(text="This image features a tiger lying on the grass. The tiger has distinctive stripes, with darker bands running horizontally across a lighter background. The animal's eyes are open and directed slightly to the side, giving it a watchful expression. The tiger's head is turned slightly to the left, and its mouth is closed. The background is a blurred green, suggesting a natural outdoor setting with grass. The lighting appears to be natural, possibly from the sun, casting soft shadows on the ground. The overall impression is one of a calm and alert predator in a natural environment.") generate_token_len=130, input_token_len=1023, session_id=0, finish_reason="stop")
(lddeploy) root@intern-studio-117673:~#
```



请用中文描述

这张图片显示了一只小猫在它的家中。猫咪身上有毛，脚有白色的毛，而耳朵则是浅棕色的。它正在弯腰地站在一个蓝色的扫帚里，扫帚充满了某种物品，猫咪看起来很有趣，似乎想挖掘这些物品。扫帚下方的地面是地板的，看起来是材质较为丰富的地板。在猫咪的周围，可以看到一些家具，如桌子和柜子，这些家具都是暴露在光线下，显示出家居的温馨氛围。

Flag

5、定量比较LMDeploy与Transformer库的推理速度差异

```
(lddeploy) root@intern-studio-117673:~# python benchmark_transformer.py
Loading checkpoint shards: 0% | 0/2 [00:00:00.00, ?it/s]
Loading checkpoint shards: 100% | 2/2 [00:19:00.00, 9.17it/s]
Warm up... [1/5]
Warm up... [2/5]
Warm up... [3/5]
Warm up... [4/5]
Warm up... [5/5]
Speed: 53.687 words/s
(lddeploy) root@intern-studio-117673:~#
(lddeploy) root@intern-studio-117673:~# touch benchmark_lddeploy.py
(lddeploy) root@intern-studio-117673:~# python benchmark_lddeploy.py
[WARNING] gpt4_config.in is not found, using default GPT4 also
Warm up... [1/5]
Warm up... [2/5]
Warm up... [3/5]
Warm up... [4/5]
Warm up... [5/5]
Speed: 481.161 words/s
(lddeploy) root@intern-studio-117673:~#
```

可以看到，Transformer库的推理速度约为53.687words/s，LMDeploy的推理速度约为481.161 words/s，是Transformer库的9倍。