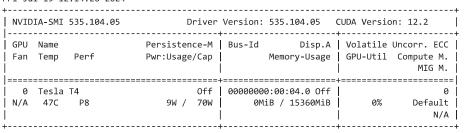
!nvidia-smi

→ Fri Jul 19 12:14:20 2024



!pip install -r requirement.txt

```
Requirement already satisfied: accelerate==0.29.3 in /usr/local/lib/python3.10/dist-packages (from -r requirement.txt (line 1)) (0.29.3)
Requirement already satisfied: bitsandbytes==0.43.1 in /usr/local/lib/python3.10/dist-packages (from -r requirement.txt (line 2)) (0.43.
Requirement already satisfied: transformers==4.40.0 in /usr/local/lib/python3.10/dist-packages (from -r requirement.txt (line 3)) (4.40.
Requirement already satisfied: numpy>=1.17 in /usr/local/lib/python3.10/dist-packages (from accelerate==0.29.3->-r requirement.txt (line
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.10/dist-packages (from accelerate==0.29.3->-r requirement.txt (
Requirement already satisfied: psutil in /usr/local/lib/python3.10/dist-packages (from accelerate==0.29.3->-r requirement.txt (line 1))
Requirement already satisfied: pyyaml in /usr/local/lib/python3.10/dist-packages (from accelerate==0.29.3->-r requirement.txt (line 1))
Requirement already satisfied: torch>=1.10.0 in /usr/local/lib/python3.10/dist-packages (from accelerate==0.29.3->-r requirement.txt (li
Requirement already satisfied: huggingface-hub in /usr/local/lib/python3.10/dist-packages (from accelerate==0.29.3->-r requirement.txt (
Requirement already satisfied: safetensors>=0.3.1 in /usr/local/lib/python3.10/dist-packages (from accelerate==0.29.3->-r requirement.tx
Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from transformers==4.40.0->-r requirement.txt (line
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.10/dist-packages (from transformers==4.40.0->-r requirement.t
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from transformers==4.40.0->-r requirement.txt (line
Requirement already satisfied: tokenizers<0.20,>=0.19 in /usr/local/lib/python3.10/dist-packages (from transformers==4.40.0->-r requirement already satisfied: tokenizers<0.20, token
Requirement already satisfied: tqdm>=4.27 in /usr/local/lib/python3.10/dist-packages (from transformers==4.40.0->-r requirement.txt (lin
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub->accelerate==0.29.3->-r
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub->accelerate==
Requirement already satisfied: sympy in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate==0.29.3->-r requirement.
Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate==0.29.3->-r requirement
Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate==0.29.3->-r requirement
Requirement already satisfied: nvidia-cuda-nvrtc-cu12==12.1.105 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelera
Requirement already satisfied: nvidia-cuda-runtime-cu12==12.1.105 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accele
Requirement already satisfied: nvidia-cuda-cupti-cu12==12.1.105 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelera
Requirement already satisfied: nvidia-cudnn-cu12==8.9.2.26 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate==0
Requirement already satisfied: nvidia-cublas-cul2==12.1.3.1 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate==
Requirement already satisfied: nvidia-cufft-cu12==11.0.2.54 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate==
Requirement already satisfied: nvidia-curand-cu12==10.3.2.106 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate
Requirement already satisfied: nvidia-cusolver-cu12==11.4.5.107 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelera
Requirement already satisfied: nvidia-cusparse-cu12==12.1.0.106 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelera
Requirement already satisfied: nvidia-nccl-cu12==2.20.5 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate==0.29
Requirement already satisfied: nvidia-nvtx-cu12==12.1.105 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate==0.
Requirement already satisfied: triton==2.3.1 in /usr/local/lib/python3.10/dist-packages (from torch>=1.10.0->accelerate==0.29.3->-r requirement already satisfied: triton==0.29.3->-r requirement already satisfied: triton=0.29.3->-r 
Requirement already satisfied: nvidia-nvjitlink-cu12 in /usr/local/lib/python3.10/dist-packages (from nvidia-cusolver-cu12==11.4.5.107->
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->transformers==4.40.0-
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->transformers==4.40.0->-r requirement
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->transformers==4.40.0->-r re
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->transformers==4.40.0->-r re
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2->torch>=1.10.0->accelerate==0.29.
Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.10/dist-packages (from sympy->torch>=1.10.0->accelerate==0.2
```

import json
import torch
from transformers import (AutoTokenizer, AutoModelForCausalLM,BitsAndBytesConfig,pipeline)

config_data = json.load(open('config.json'))

model_name="meta-llama/Meta-Llama-3-8B"

HF TOKEN = config data['HF TOKEN']

```
bnb_config = BitsAndBytesConfig(
   load_in_4bit=True,
    bnb_4bit_use_double_quant=True,
   bnb_4bit_quant_type="nf4",
   bnb_4bit_compute_dtype=torch.bfloat16
)
tokenizer = AutoTokenizer.from_pretrained(model_name, token=HF_TOKEN)
tokenizer.pad_token = tokenizer.eos_token
warnings.warn(
     /usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:89: UserWarning:
     The secret `HF\_TOKEN` does not exist in your Colab secrets.
     To authenticate with the Hugging Face Hub, create a token in your settings tab (https://
     You will be able to reuse this secret in all of your notebooks.
     Please note that authentication is recommended but still optional to access public model
       warnings.warn(
                                                                  73.0/73.0 [00:00<00:00, 1.29kB/s]
     special_tokens_map.json: 100%
     Special tokens have been added in the vocabulary, make sure the associated word embeddir
model=AutoModelForCausalLM.from_pretrained(
   model name.
   token=HF_TOKEN,
   quantization_config=bnb_config,
    device_map="auto"
)
     config.json: 100%
                                                            654/654 [00:00<00:00, 14.5kB/s]
     model.safetensors.index.json: 100%
                                                                 23.9k/23.9k [00:00<00:00, 392kB/s]
     Downloading shards: 100%
                                                                    4/4 [02:50<00:00, 38.27s/it]
     model-00001-of-
                                                               4.98G/4.98G [00:50<00:00, 20.8MB/s]
     00004.safetensors: 100%
     model-00002-of-
                                                               5.00G/5.00G [00:46<00:00, 268MB/s]
     00004.safetensors: 100%
     model-00003-of-
                                                               4.92G/4.92G [00:54<00:00, 51.8MB/s]
     00004.safetensors: 100%
     model-00004-of-
                                                               1.17G/1.17G [00:18<00:00, 67.8MB/s]
text_generator = pipeline(
   'text-generation',
   model=model,
  tokenizer=tokenizer,
   max_new_tokens=128,
)
def get_response(prompt):
    sequences = text_generator(prompt)
   gen_text = sequences[0]['generated_text']
   return gen_text
prompt = " what is Llama3 "
1lama3_response = get_response(prompt)
llama3_response
    ' what is Llama3 2.0? It is an extension to Llama2.0 which provides additional features. The Llama3.0 is a new and improved version of
```

' what is Llama3 2.0? It is an extension to Llama2.0 which provides additional features. The Llama3.0 is a new and improved version of Llama2.0, which has been developed with the help of the Llama2.0 developers. Llama3.0 is the first version of Llama2.0 which is a new and improved version of Llama2.0. The Llama3.0 is the first version of Llama2.0. The Llama3.0 is the first version of Llama2.1.