

1 !nvidia-smi

```
Thu May 30 14:57:54 2024
```

+-----+-----+-----+-----+-----+											
NVIDIA-SMI 535.104.05			Driver Version: 535.104.05			CUDA Version: 12.2					
+-----+-----+-----+-----+-----+											
GPU		Name		Persistence-M		Bus-Id		Disp.A		Volatile Uncorr. ECC	
Fan		Temp		Perf		Pwr:Usage/Cap		Memory-Usage		GPU-Util Compute M.	
										MIG M.	
+-----+-----+-----+-----+-----+											
0		NVIDIA L4		Off		00000000:00:03.0		Off		0	
N/A		48C P8		12W / 72W		1MiB / 23034MiB		0%		Default	
										N/A	
+-----+-----+-----+-----+-----+											

+-----+-----+-----+-----+-----+											
Processes:											
GPU		GI		CI		PID		Type		Process name	
		ID		ID						GPU Memory	
										Usage	
+-----+-----+-----+-----+-----+											
No running processes found											
+-----+-----+-----+-----+-----+											

```
1 !pip install openai
2 !pip install gradio
3
4 import IPython
5 import sys
6 import os
7
8 # Run the installation commands
9 if 'google.colab' in sys.modules:
10     print("Running in Google Colab")
11     !pip install bitsandbytes accelerate
12 else:
13     print("Not running in Google Colab")
14     !pip install transformers accelerate datasets bitsandbytes
15
16 import openai
17
18
19 def clean_notebook():
20     IPython.display.clear_output(wait=True)
21     print("Notebook cleaned.")
22
23 # Clean up the notebook
24 clean_notebook()
```

Notebook cleaned.


```
1
2 openai.api_key = "sk-proj-sQ00S1mGBrbhcg02bGpDT3B1bkFJp0G6rn90fMM1SLdq14y0"
3
4 os.environ['HF_TOKEN'] ="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"
5 hf_token = os.environ['HF_TOKEN']

1 # Add information from text file
2 if 'google.colab' in sys.modules:
3     from google.colab import drive
4     drive.mount('/content/drive')
5
6     with open("/content/drive/My Drive/Colab Notebooks/MasterAI KMITL Training/stx_history.txt", 'r') as f:
7         history = f.read()
8
9     drive.flush_and_unmount()
10
11 else:
12     print("Not running in Google Colab")
13     with open("./stx_history.txt", 'r', errors="ignore") as f:
14         history = f.read()
15
16 #IPython.display.Markdown(history)
17
18 # Define system expertise
19 system_prompt = """You are a technology blogger whos has more than 15 years \
20 experience working in IT company. Your answers should be \
21 short, concise, and clear. If you don't know the answer, say Sorry, I can't \
22 find a good answer for you.
23 """
```

Mounted at /content/drive

```
1 del model_typhoon, tokenizer_typhoon
2 torch.cuda.empty_cache()
```

```
1 from transformers import AutoModelForCausalLM,AutoTokenizer,BitsAndBytesConfig
2 import torch
3
4 #####
5 # bitsandbytes parameters
6 #####
7
8 # Activate 4-bit precision base model loading
9 use_4bit = True
10
11 # Compute dtype for 4-bit base models
12 bnb_4bit_compute_dtype = "float16"
13
14 # Quantization type (fp4 or nf4)
15 bnb_4bit_quant_type = "nf4"
16
17 # Activate nested quantization for 4-bit base models (double quantization)
18 use_nested_quant = False
19
20
21 # Load tokenizer and model with QLoRA configuration
22 compute_dtype = getattr(torch, bnb_4bit_compute_dtype)
23
24 bnb_config = BitsAndBytesConfig(
25     load_in_4bit=use_4bit,
26     bnb_4bit_quant_type=bnb_4bit_quant_type,
27     bnb_4bit_compute_dtype=compute_dtype,
28     bnb_4bit_use_double_quant=use_nested_quant,
29 )
30
31
32 model_typhoon = AutoModelForCausalLM.from_pretrained("scb10x/llama-3-typhoon-v1.5-8b-instruct", quantization_config=bnb_config,device_map="auto",token=hf_token)
33 tokenizer_typhoon = AutoTokenizer.from_pretrained("scb10x/llama-3-typhoon-v1.5-8b-instruct")
34
35 model_llama = AutoModelForCausalLM.from_pretrained("meta-llama/Meta-Llama-3-8B-Instruct", quantization_config=bnb_config,device_map="auto",token=hf_token)
36 tokenizer_llama = AutoTokenizer.from_pretrained("meta-llama/Meta-Llama-3-8B-Instruct")
```

 /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://huggingface.co/settings/tokens>), set it as secret in y
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 models or datasets.

config.json: 100%	616/616 [00:00<00:00, 54.0kB/s]
model.safetensors.index.json: 100%	23.9k/23.9k [00:00<00:00, 1.96MB/s]
Downloading shards: 100%	4/4 [00:47<00:00, 10.17s/it]
model-00001-of-00004.safetensors: 100%	4.98G/4.98G [00:16<00:00, 372MB/s]
model-00002-of-00004.safetensors: 100%	5.00G/5.00G [00:12<00:00, 429MB/s]
model-00003-of-00004.safetensors: 100%	4.92G/4.92G [00:14<00:00, 419MB/s]
model-00004-of-00004.safetensors: 100%	1.17G/1.17G [00:02<00:00, 411MB/s]
Loading checkpoint shards: 100%	4/4 [00:08<00:00, 1.85s/it]
generation_config.json: 100%	194/194 [00:00<00:00, 18.2kB/s]
tokenizer_config.json: 100%	51.0k/51.0k [00:00<00:00, 4.51MB/s]
tokenizer.json: 100%	9.09M/9.09M [00:00<00:00, 9.44MB/s]
special_tokens_map.json: 100%	449/449 [00:00<00:00, 37.4kB/s]
Special tokens have been added in the vocabulary, make sure the associated word embeddings are fine-tuned or trained.	
config.json: 100%	654/654 [00:00<00:00, 57.5kB/s]
model.safetensors.index.json: 100%	23.9k/23.9k [00:00<00:00, 2.06MB/s]
Downloading shards: 100%	4/4 [00:43<00:00, 9.51s/it]
model-00001-of-00004.safetensors: 100%	4.98G/4.98G [00:11<00:00, 442MB/s]
model-00002-of-00004.safetensors: 100%	5.00G/5.00G [00:13<00:00, 411MB/s]
model-00003-of-00004.safetensors: 100%	4.92G/4.92G [00:13<00:00, 439MB/s]
model-00004-of-00004.safetensors: 100%	1.17G/1.17G [00:02<00:00, 391MB/s]
Loading checkpoint shards: 100%	4/4 [00:08<00:00, 1.88s/it]
generation_config.json: 100%	187/187 [00:00<00:00, 17.0kB/s]
tokenizer_config.json: 100%	51.0k/51.0k [00:00<00:00, 4.54MB/s]
tokenizer.json: 100%	9.09M/9.09M [00:00<00:00, 12.9MB/s]
special_tokens_map.json: 100%	73.0/73.0 [00:00<00:00, 6.66kB/s]
Special tokens have been added in the vocabulary, make sure the associated word embeddings are fine-tuned or trained.	

```

1 #Typhoon and Llama
2
3 def generate_Typhoon_Llama(model, tokenizer, user_prompt, temperature, max_token=256, top_p=0.9):
4
5     messages = [
6         {"role": "system", "content": system_prompt},
7         {"role": "user", "content": user_prompt},
8     ]
9
10    input_ids = tokenizer.apply_chat_template(messages,add_generation_prompt=True,return_tensors="pt").to(model.device)
11
12    terminators = [
13        tokenizer.eos_token_id,
14        tokenizer.convert_tokens_to_ids("<|eot_id|>")
15    ]
16
17    outputs = model.generate(
18        input_ids,
19        max_new_tokens=max_token,
20        pad_token_id=model.config.eos_token_id,
21        eos_token_id=terminators,
22        do_sample=True,
23        temperature=temperature,
24        top_p=top_p,
25    )
26
27    response = outputs[0][input_ids.shape[-1]:]
28    return tokenizer.decode(response, skip_special_tokens=True)

```

```

1 # Open AI
2 def generate_OpenAI(user_prompt, temperature, max_token=256):
3     completion = openai.chat.completions.create(
4         model='gpt-3.5-turbo',
5         messages=[
6             {"role": "system", "content": system_prompt},
7             {"role": "user", "content": user_prompt},
8         ],
9         temperature=temperature,
10        max_tokens=max_token,
11    )
12    return completion.choices[0].message.content

```

```

1 def answer(question, temperature, max_token, model_list):
2     prompt = f""""Please answer the following question:
3
4     Question:
5
6     ```{question}```
7
8     Use the following context to find the answer:
9
10    ```{history}```
11    """
12    #print(prompt)
13
14    result1 = result2 = result3 = "Model is not selected."
15
16    if "OpenAI" in model_list:
17        result1 = generate_OpenAI(prompt, temperature, max_token)
18
19    if "Typhoon" in model_list:
20        result2 = generate_Typhoon_Llama(model_typhoon, tokenizer_typhoon, prompt, temperature, max_token)
21
22    if "Llama" in model_list:
23        result3 = generate_Typhoon_Llama(model_Llama, tokenizer_Llama, prompt, temperature, max_token)
24
25    result = [result1, result2, result3]
26    return result

```

```

1 import gradio as gr
2
3 # Example prompts
4 examples = [
5     ["What is Seagate?"],
6     ["What are the Seagate products?"],
7     ["Current Seagate stock price?"]
8 ]
9
10 # Create Gradio interface
11 interface = gr.Interface(
12     fn=answer,
13     inputs=[
14         gr.Textbox(lines=2, placeholder="Enter your prompt here..."),
15         gr.Slider(minimum=0.1, maximum=1.0, step=0.1, value=0.85, label="Temperature"),
16         gr.Slider(minimum=0.1, maximum=512, step=1, value=256, label="token"),
17         gr.CheckboxGroup(["OpenAI", "Typhoon", "Llama"], label="Model selections", value=["OpenAI", "Typhoon", "Llama"])
18     ],
19     outputs=[gr.Textbox(label="OpenAI result"),gr.Textbox(label="Typhoon result"),gr.Textbox(label="Llama result")],
20     title="Q&A about Seagate",
21     description="Ask anything that you want to know about Seagate",
22     examples=examples
23 )
24
25 # Launch the interface
26 interface.launch(share=True)
27

```

Colab notebook detected. To show errors in colab notebook, set debug=True in launch()
Running on public URL: <https://c7fc2994fb4c747855.gradio.live>

This share link expires in 72 hours. For free permanent hosting and GPU upgrades, run `gradio deploy` from Terminal to deploy to Spaces (https://huggingface.co/docs/gradio/guides/deploy_to_spaces)

Q&A about Seagate

Ask anything that you want to know about Seagate

question

What is Seagate?

Temperature

0.85

token

256

Model selections

☐ OpenAI

☐ Typhoon

☐ Llama

Clear

Submit

OpenAI result

Seagate is an American data storage company that specializes in producing hard disk drives (HDDs), hybrid drives, and solid-state drives (SSDs). It was founded in 1979 as Shugart Technology and has since become a major player in the computer storage industry.

Typhoon result

Seagate Technology Holdings plc is an American data storage company. It was initially incorporated in 1978 as Shugart Technology and commenced business in 1979. Seagate developed the first 5.25-inch hard disk drive in 1980. They were a major supplier in the microcomputer market during the 1980s, especially after the introduction of the IBM XT in 1983. Much of their growth has come from the microcomputer market.

Llama result

Seagate is an American data storage company that was incorporated in 1978. It is one of the largest and most well-known hard drive manufacturers in the world.