


```

4 if torch.cuda.is_available():
5     print("GPU:", torch.cuda.get_device_name(0))
6 print("Python:", platform.python_version())
7

```

```

PyTorch: 2.6.0+cu124
CUDA available: True
GPU: Tesla T4
Python: 3.11.13

```

```

1 !pip -q install -U "transformers>=4.44" "accelerate>=0.34"
2

```

```

WARNING: Ignoring invalid distribution ~vidia-cudnn-cu12 (/usr/local/lib/python3.11/dist-packages)
WARNING: Ignoring invalid distribution ~vidia-cudnn-cu12 (/usr/local/lib/python3.11/dist-packages)
WARNING: Ignoring invalid distribution ~vidia-cudnn-cu12 (/usr/local/lib/python3.11/dist-packages)
WARNING: Ignoring invalid distribution ~vidia-cudnn-cu12 (/usr/local/lib/python3.11/dist-packages)

```

1. Error Loading

```

1 from transformers import pipeline, AutoTokenizer
2 import torch
3
4 def main():
5     try:
6         model_id = "TinyLlama/TinyLlama-1.1B-Chat-v1.0"
7
8         print("Loading tokenizer...")
9         tok = AutoTokenizer.from_pretrained(model_id, use_fast=True)
10
11         prompt = "write a poem about I am Just a Girl."
12         messages = [{"role": "system", "content": prompt}]
13
14         # Build the input text using the model's chat template if present.
15         try:
16             input_text = tok.apply_chat_template(
17                 messages,
18                 tokenize=False,
19                 add_generation_prompt=True
20             )
21         except Exception as e:
22             print(f" Chat template not available, falling back to raw prompt. Details: {e}")
23             input_text = prompt # fallback
24
25         use_gpu = torch.cuda.is_available()
26         dtype = torch.float16 if use_gpu else torch.float32
27
28         print("⚙ Loading model & pipeline...")
29         try:
30             generator = pipeline(
31                 "text-generation",
32                 model=model_id,
33                 tokenizer=tok,
34                 device=0 if use_gpu else -1, # GPU if available
35                 torch_dtype=dtype,
36             )
37         except Exception as e:
38             print(f" Error loading pipeline: {e}")
39             return
40
41         print(" Generating text...")
42         try:
43             out = generator(
44                 input_text,
45                 max_new_tokens=300,      # keep modest for speed/RAM
46                 do_sample=True,
47                 temperature=0.8,
48                 top_p=0.9,
49                 repetition_penalty=1.2,
50                 pad_token_id=tok.eos_token_id if tok.eos_token_id is not None else None,
51                 eos_token_id=tok.eos_token_id if tok.eos_token_id is not None else None,
52                 return_full_text=False,
53                 truncation=True,
54             )
55             print("\n Output:\n")
56             print(out[0]["generated_text"])

```

```

57     except Exception as e:
58         print(f"❌ Error during generation: {e}")
59
60     except Exception as e:
61         print(f"❗ Unexpected error: {e}")
62
63 if __name__ == "__main__":
64     main()
65

```

⚙️ Loading tokenizer...

/usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: 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
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.

warnings.warn(
tokenizer_config.json: 1.29k/? [00:00<00:00, 36.9kB/s]
tokenizer.model: 100% 500k/500k [00:00<00:00, 607kB/s]
tokenizer.json: 1.84M/? [00:00<00:00, 19.5MB/s]
special_tokens_map.json: 100% 551/551 [00:00<00:00, 11.6kB/s]
⚙️ Loading model & pipeline...
config.json: 100% 608/608 [00:00<00:00, 17.7kB/s]
model.safetensors: 100% 2.20G/2.20G [00:32<00:00, 140MB/s]
generation_config.json: 100% 124/124 [00:00<00:00, 7.09kB/s]
Device set to use cuda:0
Generating text...

Output:

I Am Just a Girl, with all my heart and soul,
Without a doubt or hesitation,
Facing life's challenges head-on each day;
Growing in ways that make me proud to be known.

From the smallest of things, to the biggest of dreams,
I strive for excellence and always try to excel higher.
The world is full of wonders, but they are nothing compared to me,
For every step I take, every word spoken,
My spirit soars high above and beyond.

As an individual, there is no one like me,
Each part of me unique and lovely too.
There is something within, waiting to unfold,
And though it may take time, I know where I belong.

Through thick and thin, through good times and bad,
I hold on tightly to my inner strength.
It's not easy, but I never give up or lose sight,
Of what makes me who I truly am, unstoppable.

So here I stand, a girl from earth,
A little bit different yet just as brave,
Believing in myself and everything around,
I am Just a Girl, with all her might.

Oh, I Am Just a Girl, living life boldly too!

2. Text Generation

```

1 from transformers import pipeline, AutoTokenizer
2 import torch
3
4 model_id = "TinyLlama/TinyLlama-1.1B-Chat-v1.0"
5
6 # Load tokenizer
7 tok = AutoTokenizer.from_pretrained(model_id)
8
9 # Messages in chat format
10 messages = [
11     {"role": "system", "content": "You are a creative storyteller."},
12     {"role": "user", "content": "Write a short story about a girl who dreams of touching the stars."}
13 ]
14

```

```

15 # Convert messages to model input using the chat template
16 input_text = tok.apply_chat_template(messages, tokenize=False, add_generation_prompt=True)
17
18 # Load pipeline
19 generator = pipeline(
20     "text-generation",
21     model=model_id,
22     tokenizer=tok,
23     device=0 if torch.cuda.is_available() else -1,
24     torch_dtype=torch.float16 if torch.cuda.is_available() else torch.float32
25 )
26
27 # Generate output
28 out = generator(
29     input_text,
30     max_new_tokens=300,
31     do_sample=True,
32     temperature=0.8,
33     top_p=0.9,
34     repetition_penalty=1.1,
35     return_full_text=False
36 )
37
38 print("📄 Generated Story:\n")
39 print(out[0]["generated_text"])
40

```

🔄 Device set to use cuda:0
 📄 Generated Story:

Sarah had always been fascinated by the night sky, ever since she was a little girl. She would spend hours stargazing, imagining herself One evening, as she walked home from school, Sarah saw a shooting star streak across the sky. It was the most breathtaking sight she had She ran home as fast as she could, opening the window to let in the cool breeze. As she stepped outside, she felt a sense of exhilaratic Sarah closed her eyes and took a deep breath, then stepped into the street. A warm breeze blew against her face, and the sound of rustli As she drew nearer to the stars, she felt a sudden jolt of electricity, like a shockwave coursing through her body. Her heart raced as s But just as she reached out, something happened. A bright light shone

Streaming Response

```

1 from transformers import AutoTokenizer, AutoModelForCausalLM, TextIteratorStreamer
2 import torch
3 from threading import Thread
4
5 model_id = "TinyLlama/TinyLlama-1.1B-Chat-v1.0"
6
7 # Load model & tokenizer
8 tokenizer = AutoTokenizer.from_pretrained(model_id)
9 model = AutoModelForCausalLM.from_pretrained(
10     model_id,
11     torch_dtype=torch.float16 if torch.cuda.is_available() else torch.float32,
12     device_map="auto"
13 )
14
15 # Chat messages
16 messages = [
17     {"role": "system", "content": "You are a creative storyteller."},
18     {"role": "user", "content": "Write a short story about a girl who dreams of touching the stars."}
19 ]
20
21 # Build chat input
22 input_text = tokenizer.apply_chat_template(messages, tokenize=False, add_generation_prompt=True)
23 inputs = tokenizer(input_text, return_tensors="pt").to(model.device)
24
25 # Setup streaming
26 streamer = TextIteratorStreamer(tokenizer, skip_prompt=True, skip_special_tokens=True)
27
28 # Launch generation in a background thread
29 generation_kwargs = dict(
30     **inputs,
31     max_new_tokens=300,
32     do_sample=True,

```

```
33     temperature=0.8,  
34     top_p=0.9,  
35     repetition_penalty=1.1,  
36     streamer=streamer  
37 )  
38 thread = Thread(target=model.generate, kwargs=generation_kwargs)  
39 thread.start()  
40  
41 # Print tokens as they stream in  
42 print(" 📄 Streaming story:\n")  
43 for new_text in streamer:  
44     print(new_text, end="", flush=True)  
45  
46 thread.join()  
47
```

🔄 📄 Streaming story:

Lena had always been fascinated by the cosmos, and one summer evening, she found herself staring at the sky with a sense of wonder. She
As she walked, she felt a sudden surge of energy in her muscles, as if her body was awakening from a deep sleep. Lena opened her eyes an
She stood there for a while, lost in thought. Suddenly, she heard a faint sound coming from the direction of the horizon. Curious, she t
Lena ran after it, her heart pounding in her chest as she caught sight of it moving further away. She could feel a wave of excitement wa
Without thinking, she pulled out her phone and began scrolling through her social media feed. She found a few articles about space explc
