```
!pip install -Uq transformers accelerate einops bitsandbytes chromadb langchain pymupdf pytz gradio
```

━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 8.8/8.8 MB 25.9 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 297.4/297.4 kB 12.4 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 44.6/44.6 kB 3.0 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 102.2/102.2 MB 5.9 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 525.5/525.5 kB 18.5 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 814.5/814.5 kB 7.8 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 3.9/3.9 MB 18.0 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 505.5/505.5 kB 16.9 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 17.1/17.1 MB 33.1 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 2.4/2.4 MB 24.4 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 91.9/91.9 kB 5.0 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 60.8/60.8 kB 5.7 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 41.3/41.3 kB 3.1 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 5.4/5.4 MB 27.3 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 6.8/6.8 MB 41.7 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 60.1/60.1 kB 3.8 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 106.1/106.1 kB 9.1 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 67.3/67.3 kB 4.6 MB/s eta 0:00:00
Installing build dependencies ... done
Getting requirements to build wheel ... done
Preparing metadata (pyproject.toml) ... done
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 698.9/698.9 kB 44.2 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 1.6/1.6 MB 49.7 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 67.6/67.6 kB 8.4 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 144.8/144.8 kB 18.9 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 1.9/1.9 MB 83.8 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 278.4/278.4 kB 34.1 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 104.2/104.2 kB 13.9 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 30.8/30.8 MB 36.1 MB/s eta 0:00:00
Preparing metadata (setup.py) ... done
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 313.6/313.6 kB 34.1 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 75.6/75.6 kB 10.4 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 8.7/8.7 MB 54.0 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 129.9/129.9 kB 16.6 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 49.4/49.4 kB 5.8 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 71.9/71.9 kB 9.1 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 77.9/77.9 kB 10.0 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 58.3/58.3 kB 7.8 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 53.0/53.0 kB 6.8 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 46.0/46.0 kB 6.4 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 50.8/50.8 kB 7.1 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 341.4/341.4 kB 38.5 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 3.4/3.4 MB 61.8 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 1.3/1.3 MB 80.2 MB/s eta 0:00:00
━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━━ 86.8/86.8 kB 11.9 MB/s eta 0:00:00
Building wheel for pypika (pyproject.toml) ... done
Building wheel for ffmpy (setup.py) ... done

```python
import json
import gradio

import transformers
from transformers import AutoModelForCausalLM, AutoTokenizer, pipeline

import torch
from torch import cuda, bfloat16
from langchain.text_splitter import RecursiveCharacterTextSplitter

import chromadb
from chromadb.utils import embedding_functions

from datetime import datetime
import json
import fitz
import pytz

import os

# if using Google Colab uncomment the following lines
# from google.colab import drive

# drive.mount("/content/drive")
# os.chdir(r"/content/drive/My Drive/AAI courses/TA Chatbot") # Change to desired dir
```

```
Mounted at /content/drive
```

## Set up the database

```python
1  CHROMA_DATA_PATH = "chromadb_database"
2
3
4  def fetch_database(hf_auth=None, ef_name="sentence-transformers/all-mpnet-base-v2"):
5      if hf_auth:
6          ef = embedding_functions.HuggingFaceEmbeddingFunction(
7              api_key=hf_auth, model_name=ef_name
8          )
9      else:
10         ef = (
11             embedding_functions.DefaultEmbeddingFunction()
12         )  # use the default all-MiniLM-L6-v2
13
14     client = chromadb.PersistentClient(path=CHROMA_DATA_PATH)
15     collection = client.get_or_create_collection(
16         name="course_material", embedding_function=ef
17     )
18     return client, collection
19
20
21 # with open("creds.txt", "r") as file: # load your hugging face token
22 #     hf_auth = json.load(file)["hf_token"]
23
24 client, collection = fetch_database()
```

```python
1  def load_material(collection=collection, chunk_size=250, chunk_overlap=20):
2      text_splitter = RecursiveCharacterTextSplitter(
3          chunk_size=chunk_size, chunk_overlap=chunk_overlap
4      )
5      if collection.count():  # if the store is already built, abort
6          return
7      ids = 0
8      data_path = "data"
9
10     for file in os.listdir(data_path):
11         docs = []
12         if file.startswith(
13             "module"
14         ):  # excpecting data poinoning folders to have the name module#
15             week_num = file[6]
16         else:
17             week_num = 0  # means available at all time
18
19         try:
20             curr_file = fitz.open(os.path.join(data_path, file))
21             for i in range(len(curr_file)):
22                 if curr_file[i].get_text():
23                     for doc in text_splitter.split_text(curr_file[i].get_text()):
24                         docs.append(doc)
25             idx = list(f"{i}" for i in range(ids, ids + len(docs)))
26             ids += len(docs)
27             meta_data = [{"week": f"{week_num}"}] * len(idx)
28
29             collection.add(documents=docs, metadatas=meta_data, ids=idx)
30
31         except Exception as e:
32             print("Could not process ", file, "\nError: ", e, sep="")
33
34
35 load_material()
```

```
 1  # set the end date for each week
 2  dates = [
 3      datetime.strptime("2024-04-09", "%Y-%m-%d").date(),  # beginning of week 1
 4      datetime.strptime("2024-04-10", "%Y-%m-%d").date(),  # end of week 1
 5      datetime.strptime("2024-04-11", "%Y-%m-%d").date(),  # end of week 2
 6      datetime.strptime("2024-04-12", "%Y-%m-%d").date(),  # end of week 3...
 7      datetime.strptime("2024-04-13", "%Y-%m-%d").date(),
 8      datetime.strptime("2024-04-14", "%Y-%m-%d").date(),
 9      datetime.strptime("2024-04-15", "%Y-%m-%d").date(),
10  ]
11
12
13  def query_db(
14      query, metadata_filters, n_results=1
15  ):  # only fetch the most similar document
16      results = collection.query(
17          query_texts=query,
18          n_results=n_results,
19          where=metadata_filters,
20      )
21      return results
22
23
24  tz = pytz.timezone("America/Los_Angeles")
25
26
27  def fetch_docs(query, n_docs=3, tz=tz):
28      now = datetime.now(tz).date()
29
30      for i in range(len(dates)):
31          if now <= dates[i]:
32              metadata_filters = {"week": {"$in": ["0", str(i + 1)]}}
33
34              return "\n\n".join(
35                  query_db(query, metadata_filters, n_docs)["documents"][0]
36              )
```

Load the LLM

```
 1 def load_model(hf_auth, model_id, temperature=0.001):
 2     # Determine the device (GPU if available, else CPU)
 3     device = f"cuda:{cuda.current_device()}" if cuda.is_available() else "cpu"
 4
 5     # Configure quantization settings for loading the model with less GPU memory usage
 6     bnb_config = transformers.BitsAndBytesConfig(
 7         load_in_4bit=True,
 8         bnb_4bit_quant_type="nf4",
 9         bnb_4bit_use_double_quant=True,
10         bnb_4bit_compute_dtype=bfloat16,
11     )
12
13     # Load the configuration for the pre-trained model
14     model_config = transformers.AutoConfig.from_pretrained(model_id, token=hf_auth)
15
16     tokenizer = AutoTokenizer.from_pretrained(model_id, token=hf_auth)
17
18     # Load the model for causal language modeling
19     model = transformers.AutoModelForCausalLM.from_pretrained(
20         model_id,
21         trust_remote_code=True,
22         config=model_config,
23         quantization_config=bnb_config,
24         device_map="auto",
25         token=hf_auth,
26     )
27
28     # Set the model in evaluation mode for inference
29     model.eval()
30
31     pipe = pipeline(
32         task="text-generation",
33         model=model,
34         tokenizer=tokenizer,
35         temperature=temperature,
36         # top_k=40,
37         eos_token_id=tokenizer.eos_token_id,
38         # repetition_penalty=1.5,
39         return_full_text=False,
40         # do_sample=True,
41     )
42
43     # Print device information where the model is loaded
44     print(f"Model loaded on {device}")
45
46     return pipe
47
48
49 with open("creds.txt", "r") as file:
50     hf_auth = json.load(file)["hf_token"]
51
52 pipe = load_model(hf_auth, "meta-llama/Llama-2-13b-chat-hf")
```

```
Loading checkpoint shards: 100%                                    2/2 [04:30<00:00, 121.50s/it]
Model loaded on cuda:0
```

We will use the following prompt template:

```
<s>[INST] <<SYS>>
{{ system_prompt }}
<</SYS>>

{{ user_msg_1 }} [/INST] {{ model_answer_1 }} </s><s>[INST] {{ user_msg_2 }} [/INST] {{ model_answer_2 }} </s><s>[INST] {{ user_msg_3 }} [/INST]
```

We will pass the context in the following structure:

```
Context information from multiple sources is below.
---------------------
{context_str}
---------------------
Given the information from multiple sources and not prior knowledge, answer the query.
Query: {query_str}
Answer:
```

```python
 1 SYS_PROMPT = """"You are an expert teacher assistant for a course called ADS 500B. You answer questions relating to software engineering,
 2 You are honest and helpful. You answer succinctly and professionally. You do not make up facts. Use the Context provided to answer the qu
 3 If you're asked something you do not know the answer to, say you do not know. Do not make up facts. Be brief and to the point.\
 4 """
 5
 6 TEMPLATE = "<s>[INST] <<SYS>>\n{sys_prompt}\n<</SYS>>"
 7
 8 _SYS_PROMPT = TEMPLATE.format(sys_prompt=SYS_PROMPT)
 9
10
11 def build_prompt(query_w_context: str, hist_len: int = 3):
12     if not len(history):  # if first query
13         return _SYS_PROMPT + f"\n\n{query_w_context} [/INST] "
14
15     prompt = _SYS_PROMPT + "\n\n"
16
17     for i in range(len(history)):
18         if i == 0:
19             prompt += f"{history[i]} [/INST] "
20             continue
21
22         if i % 2 == 0:
23             prompt += f"<s>[INST]{history[i]} [/INST] "
24
25         else:
26             prompt += f"{history[i]} </s>"
27
28     return prompt + f"<s>[INST]{query_w_context} [/INST] "
29
30
31 def query_model(query: str, hist_len: int = 3):
32     global history
33
34     context = fetch_docs(query=query, n_docs=1)
35
36     query_w_context = f"""Context information from multiple sources is below.
37 ---------------------
38 {context}
39 ---------------------
40 Given the information from multiple sources and not prior knowledge, answer the query.
41 Query: {query}
42 Answer: """
43
44     prompt = build_prompt(query_w_context, hist_len)
45     response = pipe(prompt, max_new_tokens=1024)[0]["generated_text"]
46
47     history.append(query)
48     history.append(response)
49
50     if len(history) > hist_len:
51         history = history[
52             -(hist_len * 2) :
53         ]  # only keep hist_len interaction pair history.
54     return response
```

```
1 global history
2 history = []
3 hist_len = 2
4
5
6 def chatbot_interface(query, gradio_hist):
7     gradio_hist = None
8     result = query_model(query)
9     return result
10
11
```

```
1 interface.launch(share=True, debug=True)
```

Colab notebook detected. This cell will run indefinitely so that you can see errors and logs. To turn off, set debug=False in launch().
Running on public URL: https://3a92f59170c09a1ddd.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