

# Personalized Healthcare QA System (Chatbot) - Chatbot UI

Ed:215279167

## 1. Import libraries and model for Chatbot UI

```

1 #!pip install datasets requests bitsandbytes accelerate peft trl sentencepiece wandb transformers evaluate rouge_score bert_score gradio
2 !pip install requests bitsandbytes wandb transformers gradio langchain-huggingface

Collecting tomli<0.14.0,>=0.12.0 (from gradio)
  Downloading tomli-0.13.2-py3-none-any.whl.metadata (2.7 kB)
Requirement already satisfied: typer<1.0,>=0.12 in /usr/local/lib/python3.11/dist-packages (from gradio) (0.15.3)
Collecting uvicorn>=0.14.0 (from gradio)
  Downloading uvicorn-0.34.2-py3-none-any.whl.metadata (6.5 kB)
Requirement already satisfied: fsspec in /usr/local/lib/python3.11/dist-packages (from gradio-client==1.10.0->gradio) (2025.3.2)
Requirement already satisfied: websockets<16.0,>=10.0 in /usr/local/lib/python3.11/dist-packages (from gradio-client==1.10.0->gradio) (13.1)
Collecting langchain-core<1.0.0,>=0.3.59 (from langchain-huggingface)
  Downloading langchain_core-0.3.59-py3-none-any.whl.metadata (5.9 kB)
Requirement already satisfied: sentence-transformers>=2.6.0 in /usr/local/lib/python3.11/dist-packages (from langchain-huggingface) (3.0.1)
Requirement already satisfied: sniffio>=1.1 in /usr/local/lib/python3.11/dist-packages (from anyio<5.0,>=3.0->gradio) (1.3.1)
Requirement already satisfied: six>=1.4.0 in /usr/local/lib/python3.11/dist-packages (from docker-pycrds>=0.4.0->wandb) (1.17.0)
Requirement already satisfied: gitdb<5,>=4.0.1 in /usr/local/lib/python3.11/dist-packages (from gitpython!=3.1.29,>=1.0.0->wandb) (4.0.11)
Requirement already satisfied: httpcore==1.* in /usr/local/lib/python3.11/dist-packages (from httpx>=0.24.1->gradio) (1.0.9)
Requirement already satisfied: h11>=0.16 in /usr/local/lib/python3.11/dist-packages (from httpcore==1.*->httpx>=0.24.1->gradio) (0.14.0)
Requirement already satisfied: langsmith<0.4,>=0.1.125 in /usr/local/lib/python3.11/dist-packages (from langchain-core<1.0.0,>=0.3.59) (0.1.125)
Requirement already satisfied: tenacity!=8.4.0,<10.0.0,>=8.1.0 in /usr/local/lib/python3.11/dist-packages (from langchain-core<1.0.0,>=0.3.59) (9.0.0)
Requirement already satisfied: jsonpatch<2.0,>=1.33 in /usr/local/lib/python3.11/dist-packages (from langchain-core<1.0.0,>=0.3.59) (1.33)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas<3.0,>=1.0->gradio) (2.9.0)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas<3.0,>=1.0->gradio) (2025.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas<3.0,>=1.0->gradio) (2025.2)
Requirement already satisfied: annotated-types>=0.6.0 in /usr/local/lib/python3.11/dist-packages (from pydantic<3->wandb) (0.7.0)
Requirement already satisfied: pydantic-core==2.33.2 in /usr/local/lib/python3.11/dist-packages (from pydantic<3->wandb) (2.33.2)
Requirement already satisfied: typing-inspection>=0.4.0 in /usr/local/lib/python3.11/dist-packages (from pydantic<3->wandb) (0.4.0)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=2.6.0->langchain-huggingface) (1.5.2)
Requirement already satisfied: scipy in /usr/local/lib/python3.11/dist-packages (from sentence-transformers>=2.6.0->langchain-huggingface) (1.13.1)
Requirement already satisfied: networkx in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (3.4.2)
Collecting nvidia-cuda-nvrtc-cu12==12.4.127 (from torch<3,>=2.0->bitsandbytes)
  Downloading nvidia_cuda_nvrtc_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cuda-runtime-cu12==12.4.127 (from torch<3,>=2.0->bitsandbytes)
  Downloading nvidia_cuda_runtime_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cuda-cupti-cu12==12.4.127 (from torch<3,>=2.0->bitsandbytes)
  Downloading nvidia_cuda_cupti_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cudnn-cu12==9.1.0.70 (from torch<3,>=2.0->bitsandbytes)
  Downloading nvidia_cudnn_cu12-9.1.0.70-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cublas-cu12==12.4.5.8 (from torch<3,>=2.0->bitsandbytes)
  Downloading nvidia_cublas_cu12-12.4.5.8-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cufft-cu12==11.2.1.3 (from torch<3,>=2.0->bitsandbytes)
  Downloading nvidia_cufft_cu12-11.2.1.3-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-curand-cu12==10.3.5.147 (from torch<3,>=2.0->bitsandbytes)
  Downloading nvidia_curand_cu12-10.3.5.147-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cusolver-cu12==11.6.1.9 (from torch<3,>=2.0->bitsandbytes)
  Downloading nvidia_cusolver_cu12-11.6.1.9-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Collecting nvidia-cusparselt-cu12==0.6.2 (from torch<3,>=2.0->bitsandbytes)
  Downloading nvidia_cusparselt_cu12-0.6.2-py3-none-manylinux2014_x86_64.whl.metadata (1.6 kB)
Requirement already satisfied: nvidia-nvtx-cu12==12.4.127 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (12.4.127)
Collecting nvidia-nvjitlink-cu12==12.4.127 (from torch<3,>=2.0->bitsandbytes)
  Downloading nvidia_nvjitlink_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Requirement already satisfied: triton==3.2.0 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (3.2.0)
Requirement already satisfied: sympy==1.13.1 in /usr/local/lib/python3.11/dist-packages (from torch<3,>=2.0->bitsandbytes) (1.13.1)
Requirement already satisfied: mpmath<1.4,>=1.1.0 in /usr/local/lib/python3.11/dist-packages (from sympy==1.13.1->torch<3,>=2.0->bitsandbytes) (1.3.0)
Requirement already satisfied: shellingham>=1.3.0 in /usr/local/lib/python3.11/dist-packages (from typer<1.0,>=0.12->gradio) (1.5.4)
Requirement already satisfied: rich>=10.11.0 in /usr/local/lib/python3.11/dist-packages (from typer<1.0,>=0.12->gradio) (13.9.4)
Requirement already satisfied: smmap<6,>=3.0.1 in /usr/local/lib/python3.11/dist-packages (from gitdb<5,>=4.0.1->gitpython!=3.1.29,>=1.0.0->wandb) (5.0.0)
Requirement already satisfied: jsonpointer>=1.9 in /usr/local/lib/python3.11/dist-packages (from jsonpatch<2.0,>=1.33->langchain-core<1.0.0,>=0.3.59) (3.0.0)
Requirement already satisfied: requests-toolbelt<2.0.0,>=1.0.0 in /usr/local/lib/python3.11/dist-packages (from langsmith<0.4,>=0.1.125) (1.0.0)

```

IMPORTANT: Restart Colab runtime after PIP install!!!!

## 2. Implement base model and imports

First step is to import all necessary libraries and implement base model, to validate it can generate a response from prompt.

```

1 from google.colab import userdata
2 from huggingface_hub import login

```

```

3 from transformers import (
4     AutoTokenizer,
5     AutoModelForCausalLM,
6     BitsAndBytesConfig,
7     TrainingArguments,
8     logging,
9     pipeline
10 )
11 from peft import LoraConfig, PeftModel, prepare_model_for_kbit_training, get_peft_model, TaskType
12 import torch
13 import wandb
14 import pandas as pd
15 import os
16 import random
17 import re
18
19 # for chatbot UI
20 import gradio as gr


1 base_model_name = "meta-llama/Llama-2-7b-chat-hf"
2 #model_name = "digitalblue/model_e_merge_v2"
3 project_name = "ed_medical"
4 hf_username = "digitalblue"
5
6 # define peft adapters saved to huggingface hub
7 MODEL_A = "digitalblue/ed_medical-2025-04-03_09.11.29" # trained on 200 rows
8 MODEL_B = "digitalblue/ed_medical-2025-04-05_10.59.54" # trained on 800 rows
9 MODEL_C = "digitalblue/ed_medical-2025-04-05_11.36.32" # trained on 1600 rows
10 MODEL_D = "digitalblue/ed_medical-2025-04-06_10.40.29" # trained on 1600 rows with NEFTune
11 MODEL_E = "digitalblue/ed_medical-2025-04-15_12.43.48" # trained on 8000 rows of pubmedqa only w/ NEFTune
12 MODEL_F = "digitalblue/ed_medical-2025-05-05_23.51.25" # trained on 9000 row, lower learning rate, 3 epochs
13 MODEL_G = "digitalblue/ed_medical-2025-05-08_01.23.37" # trained on 9000 rows, NEFTune false
14 MODEL_H = "digitalblue/ed_medical-2025-05-11_04.05.52" # trained on 900 row, with context, 3 epochs
15 MODEL_I = "digitalblue/ed_medical-2025-05-12_02.31.47" # trained on 900 rows of MedDialog

```

```

1 # log into hugging face and wandb
2 hf_token = userdata.get('HF_TOKEN')
3 login(hf_token)
4
5 wandb_api_key = userdata.get('WANDB_API_KEY')
6 os.environ["WANDB_API_KEY"] = wandb_api_key
7 wandb.login()
8
9 # Configure Weights & Biases to record against our project
10 os.environ["WANDB_PROJECT"] = "ed_medical"
11 os.environ["WANDB_LOG_MODEL"] = "checkpoint" if True else "end"
12 os.environ["WANDB_WATCH"] = "gradients"

```

 **wandb**: Currently logged in as: **digitalblue (digitalblue-ai)** to <https://api.wandb.ai>. Use `wandb login --relogin` to force relogin

```


1 # quantisation config to use less memory when loading model
2 quant_config = BitsAndBytesConfig(
3     load_in_4bit=True,
4     bnb_4bit_use_double_quant=True,
5     bnb_4bit_compute_dtype=torch.bfloat16,
6     bnb_4bit_quant_type="nf4"
7 )

```

```

1 model = AutoModelForCausalLM.from_pretrained(
2     base_model_name,
3     device_map="auto",
4     quantization_config=quant_config)

```

 **config.json**: 100% 614/614 [00:00<00:00, 30.5kB/s]

**model.safetensors.index.json**: 100% 26.8k/26.8k [00:00<00:00, 2.56MB/s]

**Fetching 2 files**: 100% 2/2 [04:45<00:00, 285.89s/it]

**model-00002-of-00002.safetensors**: 100% 3.50G/3.50G [03:09<00:00, 71.6MB/s]

**model-00001-of-00002.safetensors**: 100% 9.98G/9.98G [04:45<00:00, 240MB/s]

**Loading checkpoint shards**: 100% 2/2 [01:23<00:00, 38.30s/it]

**generation\_config.json**: 100% 188/188 [00:00<00:00, 20.1kB/s]

```

1 # initialise tokenizer and pipeline
2 tokenizer = AutoTokenizer.from_pretrained(base_model_name, token=hf_token)

```

```
3 tokenizer.pad_token = tokenizer.eos_token
4 tokenizer.padding_side = "right"
```

tokenizer\_config.json: 100% 1.62k/1.62k [00:00<00:00, 152kB/s]

tokenizer.model: 100% 500k/500k [00:00<00:00, 16.2MB/s]

tokenizer.json: 100% 1.84M/1.84M [00:00<00:00, 2.12MB/s]

special\_tokens\_map.json: 100% 414/414 [00:00<00:00, 30.3kB/s]

```
1 #model_g = PeftModel.from_pretrained(model, MODEL_G)
2 model_i = PeftModel.from_pretrained(model, MODEL_I)
```

adapter\_config.json: 100% 813/813 [00:00<00:00, 55.2kB/s]

adapter\_model.safetensors: 100% 67.1M/67.1M [00:00<00:00, 250MB/s]

```
1 system_prompt = "You are a helpful personalised medical assistant. In your response do not include any personal names and ensure it
2
3 def get_model_response(model_x, prompt, chat_history):
4
5     messages = [
6         {"role": "system", "content": system_prompt},
7     ]
8     for turn in chat_history:
9         print(f"turn={turn}")
10        messages.append(turn)
11
12    messages.append({"role": "user", "content": prompt})
13
14    print(f"messages={messages}")
15
16    inputs = tokenizer.apply_chat_template(messages, return_tensors="pt").to("cuda")
17    outputs = model_x.generate(inputs, max_new_tokens=512, temperature=0.1)
18    response = tokenizer.decode(outputs[0], skip_special_tokens=False)
19
20    print(f"response={response}")
21    answer = response.split("[/INST]")[1].strip()
22    answer = answer.replace('</s>', '') # remove trailing special token if present
23    print(f"answer={answer}")
24    return answer
```

```
1 get_model_response(model_i, "Hi, does utility of preliminary bronchoalveolar lavage result in suspected ventilator-associated pneumo
```

messages=[{'role': 'system', 'content': 'You are a helpful personalised medical assistant. In your response do not include any pers  
response=<s> [INST] <<SYS>>  
You are a helpful personalised medical assistant. In your response do not include any personal names and ensure it is detailed. Do r  
<</SYS>>

Hi, does utility of preliminary bronchoalveolar lavage result in suspected ventilator-associated pneumonia? [/INST] Hello. I have re  
answer=Hello. I have read your query. For further information consult a pulmonologist online --> <https://www.icliniq.com/ask-a-doctor-online/pulmonologist>

```
1
2 def respond(message, chat_history):
3     try:
4         # Format prompt for Llama 2
5         print(f"chat_history={chat_history}")
6         print(f"message={message}")
7
8         # Get model output
9         model_output = get_model_response(model_i, message, chat_history)
10        print(f"model_output = {model_output}")
11        print("-----")
12
13        return model_output
14
15
16    except Exception as e:
17        print(f"Error in respond(): {str(e)}")
18        return "Error occurred", chat_history
```

```
1 force_dark_mode = ""
2 function refresh() {
3     const url = new URL(window.location);
4     if (url.searchParams.get('__theme') !== 'dark') {
5         url.searchParams.set('__theme', 'dark');
6         window.location.href = url.href;
```

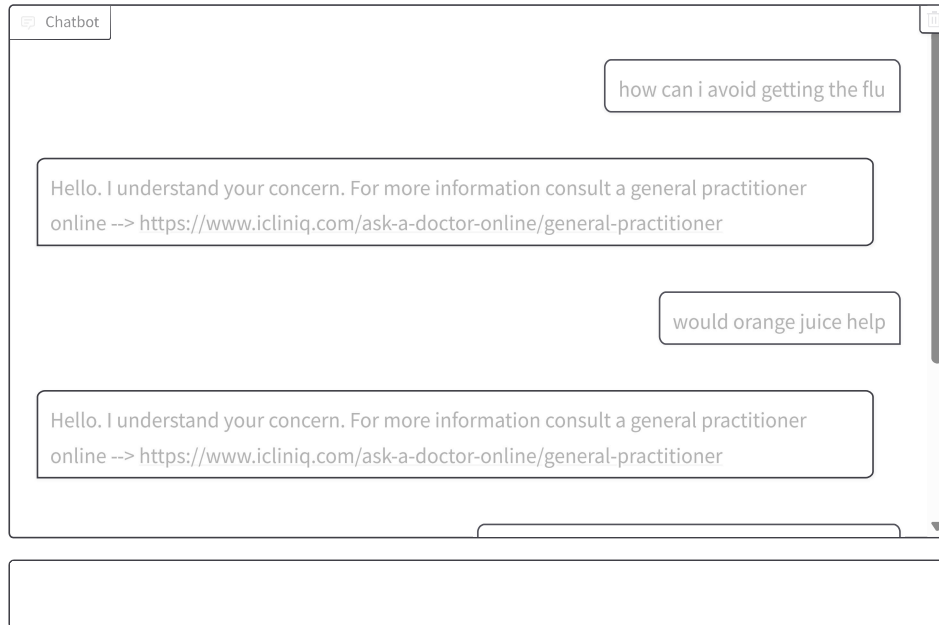
```

7     }
8 }
9 """
10
11 chatbot_ui = gr.ChatInterface(
12     fn=respond,
13     type="messages",
14     flagging_mode="never",
15     js=force_dark_mode)
16
17
18 chatbot_ui.launch(debug=True)

```

It looks like you are running Gradio on a hosted Jupyter notebook. For the Gradio app to work, sharing must be enabled. Automatic 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://563757c4598215d38e.gradio.live>

This share link expires in 1 week. For free permanent hosting and GPU upgrades, run `gradio deploy` from the terminal in the working



```

chat_history=[]
message=how can i avoid getting the flu
messages=[{'role': 'system', 'content': 'You are a helpful personalised medical assistant. In your response do not include any pers
response=<s> [INST] <<SYS>>
You are a helpful personalised medical assistant. In your response do not include any personal names and ensure it is detailed. Do r
<</SYS>>

```

```

how can i avoid getting the flu [/INST] Hello. I understand your concern. For more information consult a general practitioner online
answer=Hello. I understand your concern. For more information consult a general practitioner online --> https://www.icliniq.com/ask
model_output = Hello. I understand your concern. For more information consult a general practitioner online --> https://www.icliniq
-----

```

```

chat_history=[{'role': 'user', 'metadata': None, 'content': 'how can i avoid getting the flu', 'options': None}, {'role': 'assistant
message=would orange juice help
turn={'role': 'user', 'metadata': None, 'content': 'how can i avoid getting the flu', 'options': None}
turn={'role': 'assistant', 'metadata': None, 'content': 'Hello. I understand your concern. For more information consult a general pr
messages=[{'role': 'system', 'content': 'You are a helpful personalised medical assistant. In your response do not include any pers
response=<s> [INST] <<SYS>>
You are a helpful personalised medical assistant. In your response do not include any personal names and ensure it is detailed. Do r
<</SYS>>

```

```

how can i avoid getting the flu [/INST] Hello. I understand your concern. For more information consult a general practitioner online
answer=Hello. I understand your concern. For more information consult a general practitioner online --> https://www.icliniq.com/ask
model_output = Hello. I understand your concern. For more information consult a general practitioner online --> https://www.icliniq
-----

```

```

chat_history=[{'role': 'user', 'metadata': None, 'content': 'how can i avoid getting the flu', 'options': None}, {'role': 'assistant
message=ok what if I have a problem with a blood clot
turn={'role': 'user', 'metadata': None, 'content': 'how can i avoid getting the flu', 'options': None}
turn={'role': 'assistant', 'metadata': None, 'content': 'Hello. I understand your concern. For more information consult a general pr
turn={'role': 'user', 'metadata': None, 'content': 'would orange juice help', 'options': None}
turn={'role': 'assistant', 'metadata': None, 'content': 'Hello. I understand your concern. For more information consult a general pr
messages=[{'role': 'system', 'content': 'You are a helpful personalised medical assistant. In your response do not include any pers
response=<s> [INST] <<SYS>>
You are a helpful personalised medical assistant. In your response do not include any personal names and ensure it is detailed. Do r
<</SYS>>

```

```

how can i avoid getting the flu [/INST] Hello. I understand your concern. For more information consult a general practitioner online
answer=Hello. I understand your concern. For more information consult a general practitioner online --> https://www.icliniq.com/ask
model_output = Hello. I understand your concern. For more information consult a general practitioner online --> https://www.icliniq
-----

```

```
→ /usr/local/lib/python3.11/dist-packages/torch/cuda/memory.py:391: FutureWarning: torch.cuda.reset_max_memory_allocated now calls to
warnings.warn(
```

<https://colab.research.google.com/drive/1K3FpLxZXI81mon5KpIQIopWRQKuQ5EYk#scrollTo=Q0NvJyuYeGBJ&printMode=true>