

✓ Personalized Healthcare QA System (Chatbot) - V3 - PubMedQA

Ed:215279167

Objective:

Develop a healthcare chatbot that answers patient inquiries on medications, symptoms, and treatments while reducing hallucinations.

Update and changes:

This updated notebook reflects changes learned from weeks 1 -7 experiments with fine-tuning. Main changes include:

- Use single dataset of PubMedQA. Initial attempts to use a combination of MedDialog and PubMedQA didn't provide significant improvements in evaluation, and believe this is due to differing data quality in MedDialog dataset.
- Significant increase of data used for fine-tuning

Datasets:

- PubMedQA – Biomedical QA dataset - <https://huggingface.co/datasets/qiaojin/PubMedQA>

Task Breakdown:

1. Train models on medical QA datasets.
2. Fine-tune for patient-friendly responses (simplified, clear language).
3. Ensure context-aware, regulatory-compliant answers:
 - o Implement FDA/TGA guideline alignment.
 - o Develop a retrieval-based validation for generated answers.
4. Deploy as a chatbot interface (React-based UI + API integration).
5. Implement real-time fact-checking:
 - o Confidence score visualization (green = high confidence, yellow = medium, red = low confidence).
 - o Integration with PubMed and trusted medical sources.

Models to Use:

- Llama-2 7B

Evaluation Metrics:

- Rouge, BLEU, Meteor, BertScore, Manual inspection

✓ 1. Import libraries and model to prepare for fine tuning

```
1 !pip install datasets requests bitsandbytes accelerate peft trl sentencepiece wandb transformers evaluate rouge_score bert-score
```



```

  _____ 363.4/363.4 MB 3.5 MB/s eta 0:00:00
  Downloading nvidia_cuda_cupti_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl (13.8 MB)
  _____ 13.8/13.8 MB 53.2 MB/s eta 0:00:00
  Downloading nvidia_cuda_nvrtc_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl (24.6 MB)
  _____ 24.6/24.6 MB 42.7 MB/s eta 0:00:00
  Downloading nvidia_cuda_runtime_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl (883 kB)
  _____ 883.7/883.7 kB 44.2 MB/s eta 0:00:00
  Downloading nvidia_cudnn_cu12-9.1.0.70-py3-none-manylinux2014_x86_64.whl (664.8 MB)
  _____ 664.8/664.8 MB 2.2 MB/s eta 0:00:00
  Downloading nvidia_cufft_cu12-11.2.1.3-py3-none-manylinux2014_x86_64.whl (211.5 MB)
  _____ 211.5/211.5 MB 5.4 MB/s eta 0:00:00
  Downloading nvidia_curand_cu12-10.3.5.147-py3-none-manylinux2014_x86_64.whl (56.3 MB)
  _____ 56.3/56.3 MB 13.3 MB/s eta 0:00:00
  Downloading nvidia_cusolver_cu12-11.6.1.9-py3-none-manylinux2014_x86_64.whl (127.9 MB)
  _____ 127.9/127.9 MB 7.5 MB/s eta 0:00:00
  Downloading nvidia_cusparse_cu12-12.3.1.170-py3-none-manylinux2014_x86_64.whl (207.5 MB)
  _____ 207.5/207.5 MB 5.4 MB/s eta 0:00:00
  Downloading nvidia_nvjitlink_cu12-12.4.127-py3-none-manylinux2014_x86_64.whl (21.1 MB)
  _____ 21.1/21.1 MB 81.2 MB/s eta 0:00:00
  Downloading xxhash-3.5.0-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (194 kB)
  _____ 194.8/194.8 kB 16.9 MB/s eta 0:00:00

Building wheels for collected packages: rouge_score
  Building wheel for rouge_score (setup.py) ... done
    Created wheel for rouge_score: filename=rouge_score-0.1.2-py3-none-any.whl size=24934 sha256=61018e05c29de1ac53e312cf59c6999d77e
    Stored in directory: /root/.cache/pip/wheels/1e/19/43/8a442dc83660ca25e163e1bd1f89919284ab0d0c1475475148
Successfully built rouge_score

```

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 from trl import SFTTrainer, SFTConfig, DataCollatorForCompletionOnlyLM
13 from datasets import load_dataset, Dataset, DatasetDict
14 from datetime import datetime
15 import torch
16 import wandb
17 import pandas as pd
18 import os
19 import random
20 import re
21 import matplotlib.pyplot as plt
22 from tqdm.notebook import tqdm
23 import evaluate
24 import numpy as np
25 from bert_score import BERTScorer

1 model_name = "meta-llama/Llama-2-7b-chat-hf"
2 project_name = "ed_medical"
3 hf_username = "digitalblue"

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     model_name,
3     device_map="auto",
4     quantization_config=quant_config)

```

config.json: 100% 614/614 [00:00<00:00, 54.9kB/s]

model.safetensors.index.json: 100% 26.8k/26.8k [00:00<00:00, 1.98MB/s]

Fetching 2 files: 100% 2/2 [01:16<00:00, 76.76s/it]

model-00001-of-00002.safetensors: 100% 9.98G/9.98G [01:16<00:00, 304MB/s]

model-00002-of-00002.safetensors: 100% 3.50G/3.50G [00:38<00:00, 110MB/s]

Loading checkpoint shards: 100% 2/2 [01:00<00:00, 27.73s/it]

generation_config.json: 100% 188/188 [00:00<00:00, 11.1kB/s]

```

1 memory = model.get_memory_footprint() / 1e6
2 print(f"Memory footprint: {memory:.1f} MB")

```

Memory footprint: 3,762.8 MB

```

1 # model architecture
2 model

3 LlamaForCausalLM(
4     (model): LlamaModel(
5         (embed_tokens): Embedding(32000, 4096)
6         (layers): ModuleList(
7             (0-31): 32 x LlamaDecoderLayer(
8                 (self_attn): LlamaAttention(
9                     (q_proj): Linear4bit(in_features=4096, out_features=4096, bias=False)
10                    (k_proj): Linear4bit(in_features=4096, out_features=4096, bias=False)
11                    (v_proj): Linear4bit(in_features=4096, out_features=4096, bias=False)
12                    (o_proj): Linear4bit(in_features=4096, out_features=4096, bias=False)
13                )
14                (mlp): LlamaMLP(
15                    (gate_proj): Linear4bit(in_features=4096, out_features=11008, bias=False)
16                    (up_proj): Linear4bit(in_features=4096, out_features=11008, bias=False)
17                    (down_proj): Linear4bit(in_features=11008, out_features=4096, bias=False)
18                    (act_fn): SiLU()
19                )
20                (input_layernorm): LlamaRMSNorm((4096,), eps=1e-05)
21                (post_attention_layernorm): LlamaRMSNorm((4096,), eps=1e-05)
22            )
23            (norm): LlamaRMSNorm((4096,), eps=1e-05)
24            (rotary_emb): LlamaRotaryEmbedding()
25        )
26        (lm_head): Linear(in_features=4096, out_features=32000, bias=False)
27    )
28)

```

```

1 # initialise tokenizer and pipeline
2 tokenizer = AutoTokenizer.from_pretrained(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, 89.5kB/s]

tokenizer.model: 100% 500k/500k [00:00<00:00, 7.57MB/s]

tokenizer.json: 100% 1.84M/1.84M [00:00<00:00, 5.63MB/s]

special_tokens_map.json: 100% 414/414 [00:00<00:00, 28.9kB/s]

```

1 messages = [
2     {"role": "system", "content": "You are a helpful personalised medical assistant"},
3     {"role": "user", "content": "How can I get rid of the flu?"}
4 ]

```

```

1 # verify base model generates response
2 inputs = tokenizer.apply_chat_template(messages, return_tensors="pt").to("cuda")

```

```
3 outputs = model.generate(inputs, max_new_tokens=250)
4 print(tokenizer.decode(outputs[0]))
```

→ <s> [INST] <>SYS>>
You are a helpful personalised medical assistant
<>/SYS>>

How can I get rid of the flu? [/INST] Hello there! I'm glad you asked! 😊 As a helpful personalized medical assistant, I'm here to

First and foremost, it's important to understand that the flu is a viral infection, and there's no magic pill or potion that can make it go away.

- Stay hydrated: Drink plenty of fluids, such as water, tea, or soup, to help loosen congestion and replace lost fluids. Aim for at least 8 glasses a day.
- Rest: Give your body the rest it needs to fight off the infection. Try to get at least 8 hours of sleep each night and take naps when you need them.
- Use over-the-counter medications: Over-the-counter medications like acetaminophen (Tylenol) or ibuprofen (Advil, Motrin) can help relieve fever and pain.

3. Dataset analysis and preprocessing for Llama

The datasets need to be analysed and prepared for fine-tuning the model. This includes removing invalid data and formatting into the required format with special tokens the Llama 2 model requires. Dataset acquired from HuggingFace

- PubMedQA - <https://huggingface.co/datasets/qiaojin/PubMedQA>

```
1 # format input from datasets into prompt format required by llama model
2 def format_llama_prompt(user_message, model_answer):
3     prompt = '<s>[INST]' # special token - commence instruct
4     prompt += user_message.strip()
5     prompt += ' [/INST]' # special token - end instruct
6     prompt += model_answer.strip()
7     prompt += '</s>'# special token - end
8     return prompt
9
10 def format_llama_prompt_with_context(user_message, model_answer, context):
11     prompt = '<s>[INST] Given this context: ' # special token - commence instruct
12     prompt += context['contexts'][0].strip()
13     prompt += ' Question: '
14     prompt += user_message.strip()
15     prompt += ' [/INST]' # special token - end instruct
16     prompt += model_answer.strip()
17     prompt += '</s>'# special token - end
18     return prompt
19
20 # load pub_med_qa dataset
21 #pub_med_qa = pd.read_parquet("hf://datasets/qiaojin/PubMedQA/pqa_unlabeled/train-00000-of-00001.parquet")
22 pub_med_qa = load_dataset("qiaojin/PubMedQA", "pqa_artificial")
```

→ README.md: 100% 5.19k/5.19k [00:00<00:00, 275kB/s]

train-00000-of-00001.parquet: 100% 233M/233M [00:01<00:00, 226MB/s]

Generating train split: 100% 211269/211269 [00:03<00:00, 52551.39 examples/s]

1 pub_med_qa

→ DatasetDict({
 train: Dataset({
 features: ['pubid', 'question', 'context', 'long_answer', 'final_decision'],
 num_rows: 211269
 })
})

1 pub_med_qa['train'][0]

→ {'pubid': 25429730,
 'question': 'Are group 2 innate lymphoid cells (ILC2s) increased in chronic rhinosinusitis with nasal polyps or eosinophilia?',
 'context': {'contexts': ['Chronic rhinosinusitis (CRS) is a heterogeneous disease with an uncertain pathogenesis. Group 2 innate lymphoid cells (ILC2s) represent a recently discovered cell population which has been implicated in driving Th2 inflammation in CRS; however, their relationship with clinical disease characteristics has yet to be investigated.', 'The aim of this study was to identify ILC2s in sinus mucosa in patients with CRS and controls and compare ILC2s across characteristics of disease.', 'A cross-sectional study of patients with CRS undergoing endoscopic sinus surgery was conducted. Sinus mucosal biopsies were obtained during surgery and control tissue from patients undergoing pituitary tumour resection through transphenoidal approach. ILC2s were identified as CD45(+) Lin(-) CD127(+) CD4(-) CD8(-) CRTH2(CD294)(+) CD161(+) cells in single cell suspensions through flow cytometry. ILC2 frequencies, measured as a percentage of CD45(+) cells, were compared across CRS phenotype, endotype, inflammatory CRS subtype and other disease characteristics including blood eosinophils, serum IgE, asthma status and nasal symptom score.'],
 '35 patients (40% female, age 48 ± 17 years) including 13 with eosinophilic CRS (eCRS), 13 with non-eCRS and 9 controls were recruited. ILC2 frequencies were associated with the presence of nasal polyps (P = 0.002) as well as high tissue eosinophilia (P = 0.004) and eosinophil-dominant CRS (P = 0.001) (Mann-Whitney U). They were also associated with increased blood eosinophilia (P = 0.005). There were no significant associations found between ILC2s and serum total IgE and allergic disease. In the CRS with nasal polyps (CRSwNP) population, ILC2s were increased in patients with co-existing asthma (P = 0.03). ILC2s were also correlated with

```
worsening nasal symptom score in CRS ( $P = 0.04$ ).'],
  'labels': ['BACKGROUND', 'OBJECTIVE', 'METHODS', 'RESULTS'],
  'meshes': ['Adult',
    'Aged',
    'Antigens, Surface',
    'Case-Control Studies',
    'Chronic Disease',
    'Eosinophilia',
    'Female',
    'Humans',
    'Hypersensitivity',
    'Immunity, Innate',
    'Immunoglobulin E',
    'Immunophenotyping',
    'Leukocyte Count',
    'Lymphocyte Subsets',
    'Male',
    'Middle Aged',
    'Nasal Mucosa',
    'Nasal Polyps',
    'Neutrophil Infiltration',
    'Patient Outcome Assessment',
    'Rhinitis',
    'Sinusitis',
    'Young Adult']},
  'long_answer': 'As ILC2s are elevated in patients with CRSwNP, they may drive nasal polyp formation in CRS. ILC2s are also linked with high tissue and blood eosinophilia and have a potential role in the activation and survival of eosinophils during the Th2 immune response. The association of innate lymphoid cells in CRS provides insights into its pathogenesis.',
  'final_decision': 'yes'}
```

```
1 # create list of formatted prompts
2 pub_med_qa_as_prompt = []
3 for item in pub_med_qa['train']:
4   if item['final_decision'] == 'yes': # only get rows with yes decision
5     #pub_med_qa_as_prompt.append(format_llama_prompt(item['question'], item['long_answer']))
6     pub_med_qa_as_prompt.append(format_llama_prompt_with_context(item['question'], item['long_answer'], item['context']))
7

1 print(len(pub_med_qa_as_prompt))
2 print(pub_med_qa_as_prompt[0])
```

→ 196144
`<s>[INST] Given this context: Chronic rhinosinusitis (CRS) is a heterogeneous disease with an uncertain pathogenesis. Group 2 innate

```
1 pub_med_qa_as_prompt[200]
```

→ `<s>[INST] Given this context: To investigate the effect of the percentage of free prostate-specific antigen (%fPSA) on future prostate cancer risk. Question: Is low percentage of free prostate-specific antigen (PSA) a strong predictor of later detection of prostate cancer among Japanese men with serum levels of total PSA of 4.0 ng/mL or less? [/INST] A low %fPSA is a strong predictor of

```
1 # collect all formatted data into one dataset
2 prompt_dataset = pub_med_qa_as_prompt # + med_dialog_hcm_as_prompt + med_dialog_ic_as_prompt
3 print(len(prompt_dataset))
```

→ 196144

```
1 # shuffle the combined dataset to disperse the data
2 # random.seed(42)
3 # random.shuffle(prompt_dataset)
```

```
1 len(prompt_dataset)
```

→ 196144

```
1 # convert to huggingface dataset, create splits and push to hub
2 prompt_hf_dataset = Dataset.from_dict({'text': prompt_dataset})
```

```
1 ds_train = prompt_hf_dataset.train_test_split(test_size=0.2, seed=42)
```

```
1 ds_train
```

→ DatasetDict({
 train: Dataset({
 features: ['text'],
 num_rows: 156915
 })
 test: Dataset({
 features: ['text'],
 num_rows: 39229
 })})

```

        })
    })

1 ds_test = ds_train['test'].train_test_split(test_size=0.5, seed=42)
2 ds_test

↳ DatasetDict({
    train: Dataset({
        features: ['text'],
        num_rows: 19614
    })
    test: Dataset({
        features: ['text'],
        num_rows: 19615
    })
})

1 ds_splits = DatasetDict({
2     'train': ds_train['train'],
3     'validation': ds_test['train'],
4     'test': ds_test['test']
5 })

1 ds_splits

↳ DatasetDict({
    train: Dataset({
        features: ['text'],
        num_rows: 156915
    })
    validation: Dataset({
        features: ['text'],
        num_rows: 19614
    })
    test: Dataset({
        features: ['text'],
        num_rows: 19615
    })
})

1 # push to hugging face
2 ds_splits.push_to_hub(f"{hf_username}/{project_name}-pubmedqa-artifical-with-context", private=True)

↳ Uploading the dataset shards: 100% 1/1 [00:09<00:00, 9.05s/it]
Creating parquet from Arrow format: 100% 157/157 [00:04<00:00, 27.06ba/s]
Uploading the dataset shards: 100% 1/1 [00:02<00:00, 2.72s/it]
Creating parquet from Arrow format: 100% 20/20 [00:01<00:00, 22.05ba/s]
Uploading the dataset shards: 100% 1/1 [00:01<00:00, 1.29s/it]
Creating parquet from Arrow format: 100% 20/20 [00:00<00:00, 37.79ba/s]
CommitInfo(commit_url='https://huggingface.co/datasets/digitalblue/ed_medical-pubmedqa-artifical-with-
context/commit/d305407ce4ca80e249fea443677aac7451183561', commit_message='Upload dataset', commit_description='',
nid='d305407ce4ca80e249fea443677aac7451183561', nr_url=None)

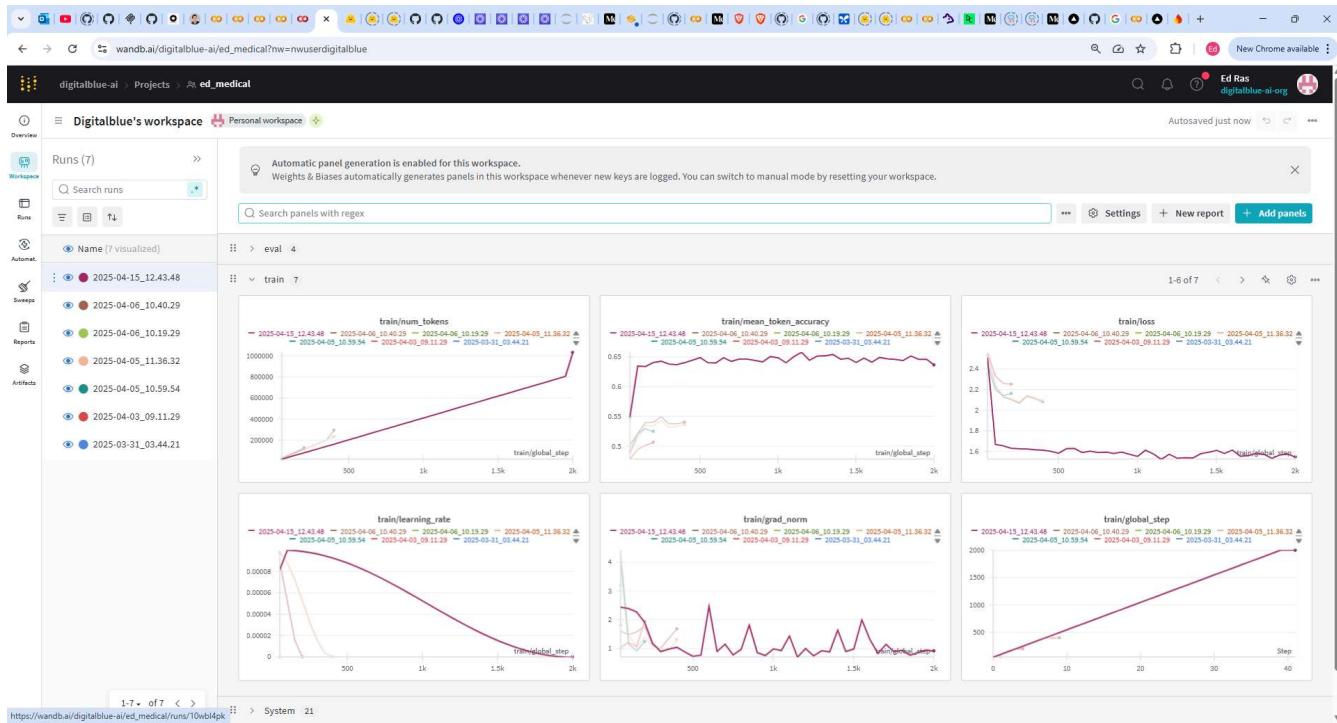
```

4. Model Training Pipeline

The model training pipeline is setup for fine-tuning. During this process additional data processing was required to filter out data with excessive token length that would exhaust GPU resources.

Initially four variants(A,B,C,D) of the model were fine-tuned on combination dataset, however focus is now on model variant E. It maintains integration with HuggingFace and Weight and Bias platform so models and runs are saved for later retrieval.

The models were fine-tuned on the **train** set, and evaluated with **validation** set.



```
1 # load pre-processed dataset created in last step from hugging face
2 #qa_dataset = load_dataset("digitalblue/ed_medical")
3 #qa_dataset = load_dataset("digitalblue/ed_medical-pubmedqa-artifical") # load only pubmedqa YES rows
4 qa_dataset = load_dataset("digitalblue/ed_medical-pubmedqa-artifical-with-context") # load only pubmedqa YES rows
5 qa_dataset
```

	README.md: 100%	532/532 [00:00<00:00, 35.4kB/s]
	train-00000-of-00001.parquet: 100%	63.8M/63.8M [00:01<00:00, 57.6MB/s]
	validation-00000-of-00001.parquet: 100%	7.96M/7.96M [00:00<00:00, 42.7MB/s]
	test-00000-of-00001.parquet: 100%	7.96M/7.96M [00:00<00:00, 28.0MB/s]
	Generating train split: 100%	156915/156915 [00:00<00:00, 228029.73 examples/s]
	Generating validation split: 100%	19614/19614 [00:00<00:00, 135985.73 examples/s]
	Generating test split: 100%	19615/19615 [00:00<00:00, 171856.41 examples/s]
	DatasetDict({ train: Dataset({ features: ['text'], num_rows: 156915 }) validation: Dataset({ features: ['text'], num_rows: 19614 }) test: Dataset({ features: ['text'], num_rows: 19615 }) })	

```
1 # getting token lengths from training dataset
2 train_length = []
3 train_token_length = []
4 for item in qa_dataset['train']:
5     if item['text']:
6         train_length.append(len(item['text']))
7         tokens = tokenizer.encode(item['text'])
8         train_token_length.append(len(tokens))
9     else:
10        print(item)
```

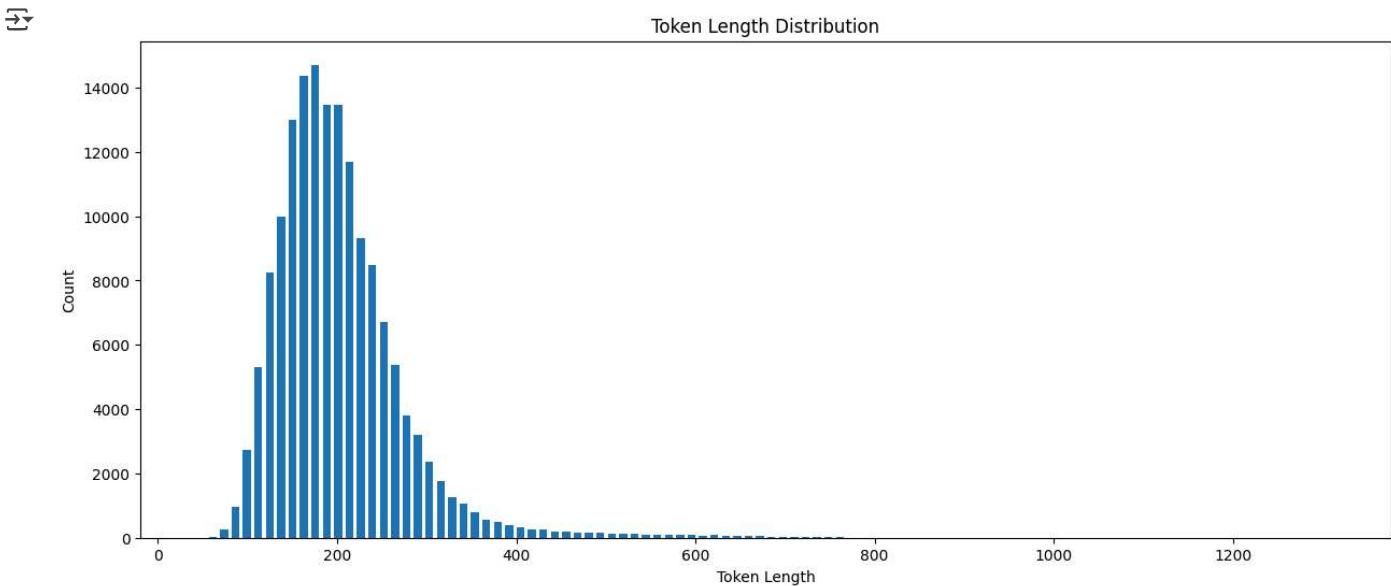
```
1 len(train_token_length)
```

```
 156915
```

```
1 # print the min, max and average character lengths and token count lengths for training data
2 print(min(train_length), max(train_length), sum(train_length)/len(train_length))
3 print(min(train_token_length), max(train_token_length), sum(train_token_length)/len(train_token_length))
```

→ 138 4677 786.0767421852595
42 1315 201.8653156167352

```
1 # plot the token lengths
2 plt.figure(figsize=(15, 6))
3 plt.hist(train_token_length, rwidth=0.7, bins=100)
4 plt.xlabel('Token Length')
5 plt.ylabel('Count')
6 plt.title('Token Length Distribution')
7 plt.show()
```



```
1 # avg is now approx 104, but will keep 230 tokens, therefore will create new datasets
2 # to filter out any rows above this token count
3 MAX_TOKENS = 250
4 DATASET_SIZE = 1000
5 def get_filtered_dataset(dataset, max_tokens, size):
6     count = 0
7     train_size = round(size * 0.9) # 90% for train
8     valid_size = round(size * 0.1) # 10% for eval
9     print(f"Train size: {train_size}, Eval size: {valid_size}")
10
11    filtered_train_dataset = []
12    for item in dataset['train']:
13        if item['text']:
14            tokens = tokenizer.encode(item['text'])
15            if len(tokens) < max_tokens:
16                filtered_train_dataset.append(item)
17                count += 1
18            if count >= train_size:
19                break
20
21    count = 0
22    filtered_eval_dataset = []
23    for item in dataset['validation']:
24        if item['text']:
25            tokens = tokenizer.encode(item['text'])
26            if len(tokens) < max_tokens:
27                filtered_eval_dataset.append(item)
28
29            count += 1
30            if count >= valid_size:
31                break
32
33    return filtered_train_dataset, filtered_eval_dataset
```

```
1 qa_train, qa_val = get_filtered_dataset(qa_dataset, MAX_TOKENS, DATASET_SIZE)
2 qa_train = Dataset.from_list(qa_train)
```

```

3 qa_val = Dataset.from_list(qa_val)
4 print(qa_train)
5 print(qa_val)

→ Train size: 900, Eval size: 100
Dataset({
    features: ['text'],
    num_rows: 900
})
Dataset({
    features: ['text'],
    num_rows: 100
})

1 qa_val[10]

→ {'text': '<s>[INST] Given this context: More than one half of lower respiratory cultures are negative for ventilator-associated pneumonia (VAP) and final reporting requires 72 hours to 96 hours. A previous retrospective study concluded that preliminary bronchoalveolar lavage (BAL) culture results (pBAL), reported at approximately 24 hours, accurately predicted final BAL culture results (fBAL). Our objective was to verify the predictive value of pBALs for fBALs, and evaluate the use of insignificant (1-99,999 cfu/mL) pBALs for rapid discontinuation of empirical antibiotics. Question: Does utility of preliminary bronchoalveolar lavage result in suspected ventilator-associated pneumonia? [/INST] Preliminary BALs were highly predictive for final results. Using insignificant pBALs appears to be a safe strategy for promptly discontinuing empirical antibiotics in trauma patients with suspected VAP. </s>'}

1 # set constants
2 MAX_SEQUENCE_LENGTH = 230 # calculated from avg token lenght in dataset
3
4 # Run name for saving the model in the hub
5 RUN_NAME = f"{datetime.now():%Y-%m-%d_%H.%M.%S}"
6 PROJECT_RUN_NAME = f"{project_name}-{RUN_NAME}"
7 HUB_MODEL_NAME = f"{hf_username}/{PROJECT_RUN_NAME}"
8
9 # qllora hyper params
10 LORA_R = 16 # Reduce LoRA rank (lower = less memory) , initial = 32
11 LORA_ALPHA = 32 # Lower alpha , initial = 64
12 TARGET_MODULES = ["q_proj", "k_proj", "v_proj", "o_proj"]
13 LORA_DROPOUT = 0.05
14 QUANT_4_BIT = True
15
16 # training hyper params
17 EPOCHS = 3
18 BATCH_SIZE = 4 # Reduce batch size , initial = 4
19 GRADIENT_ACCUMULATION_STEPS = 8 # Simulate batch size 8, initial = 1
20 LEARNING_RATE = 2e-4
21 LR_SCHEDULER_TYPE = 'cosine'
22 WARMUP_RATIO = 0.03
23 OPTIMIZER = "paged_adamw_32bit"
24 STEPS = 50
25 SAVE_STEPS = 2000
26 LOG_TO_WANDB = True

1 response_template = "[/INST]"
2 collator = DataCollatorForCompletionOnlyLM(response_template, tokenizer=tokenizer)

1 lora_params = LoraConfig(
2     r=LORA_R,
3     lora_alpha=LORA_ALPHA,
4     lora_dropout=LORA_DROPOUT,
5     target_modules=TARGET_MODULES,
6     bias="none",
7     task_type="CAUSAL_LM"
8 )

1 training_params = SFTConfig(
2     output_dir=PROJECT_RUN_NAME,
3     num_train_epochs=EPOCHS,
4     per_device_train_batch_size=BATCH_SIZE,
5     per_device_eval_batch_size=1,
6     #eval_strategy="no",
7     gradient_accumulation_steps=GRADIENT_ACCUMULATION_STEPS,
8     optim=OPTIMIZER,
9     save_steps=SAVE_STEPS,
10    save_total_limit=10,
11    logging_steps=STEPS,
12    learning_rate=LEARNING_RATE,
13    weight_decay=0.001,
14    fp16=False,
15    bf16=True,
16    max_grad_norm=0.3,

```

```

17     max_steps=-1,
18     warmup_ratio=WARMUP_RATIO,
19     group_by_length=True,
20     lr_scheduler_type=LR_SCHEDULER_TYPE,
21     report_to="wandb" if True else None,
22     run_name=RUN_NAME,
23     max_seq_length=MAX_SEQUENCE_LENGTH,
24     dataset_text_field="text",
25     save_strategy="steps",
26     hub_strategy="every_save",
27     push_to_hub=True,
28     hub_model_id=HUB_MODEL_NAME,
29     hub_private_repo=True
30     #neftune_noise_alpha=5 # using NEFTune as describe in SFT Trainer docs for increased conversational quality
31 )

```

```
1 !nvidia-smi
```

→ Sun May 11 04:10:46 2025

NVIDIA-SMI 550.54.15			Driver Version: 550.54.15		CUDA Version: 12.4		
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	Tesla T4	Off	00000000:00:04.0	Off	0%	Default	0
N/A	68C	P0	32W / 70W	6652MiB / 15360MiB			N/A

Processes:				
GPU	GI	CI	PID	Type
ID	ID			Process name
				GPU Memory Usage

```

1 # clear gpu memory
2 torch.cuda.empty_cache()
3 torch.cuda.reset_max_memory_allocated()

```

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

```

1 # Prepare the model for k-bit training (important for 4-bit)
2 model = prepare_model_for_kbit_training(model)

```

```
1 model = get_peft_model(model, lora_params)
```

```

1 trainer = SFTTrainer(
2     model=model,
3     train_dataset=qa_train,
4     eval_dataset=qa_val,
5     peft_config=lora_params,
6     args=training_params
7 )
8

```

→ Converting train dataset to ChatML: 100% 900/900 [00:00<00:00, 12383.29 examples/s]

Adding EOS to train dataset: 100% 900/900 [00:00<00:00, 12465.28 examples/s]

Tokenizing train dataset: 100% 900/900 [00:00<00:00, 1677.64 examples/s]

Truncating train dataset: 100% 900/900 [00:00<00:00, 32253.95 examples/s]

Converting eval dataset to ChatML: 100% 100/100 [00:00<00:00, 4262.20 examples/s]

Adding EOS to eval dataset: 100% 100/100 [00:00<00:00, 4069.10 examples/s]

Tokenizing eval dataset: 100% 100/100 [00:00<00:00, 1147.91 examples/s]

Truncating eval dataset: 100% 100/100 [00:00<00:00, 5246.68 examples/s]

No label_names provided for model class `PeftModelForCausalLM`. Since `PeftModel` hides base models input arguments, if label_names

Double-click (or enter) to edit

```
1 #torch.cuda.empty_cache()
2 trainer.train()
3 trainer.model.push_to_hub(PROJECT_RUN_NAME, private=True)
4 print(f"Model saved to HuggingFace as: {PROJECT_RUN_NAME}")
5 model.save_pretrained(PROJECT_RUN_NAME)
6
```

Changes to your `wandb` environment variables will be ignored because your `wandb` session has already started. For more information on how to modify your settings with `wandb.init()` arguments, please refer to [the W&B docs](#).

Tracking run with wandb version 0.19.10

Run data is saved locally in /content/wandb/run-20250511_041757-hei1h23j

Syncing run [2025-05-11 04.05.52](#) to [Weights & Biases \(docs\)](#)

View project at https://wandb.ai/digitalblue-ai/ed_medical

View run at https://wandb.ai/digitalblue-ai/ed_medical/runs/hei1h23j

wandb: **WARNING** The `get_url` method is deprecated and will be removed in a future release. Please use `run.url` instead.

`use_cache=True` is incompatible with gradient checkpointing. Setting `use_cache=False`.

```
/usr/local/lib/python3.11/dist-packages/torch/_dynamo/eval_frame.py:745: UserWarning: torch.utils.checkpoint: the use_reentrant parameter is deprecated and will be removed in a future release. Please use use_reentrant=False instead.
```

[84/84 2:46:10, Epoch 2/3]

Step Training Loss

Step	Training Loss
50	1.693400

wandb: Adding directory to artifact (./ed_medical-2025-05-11_04.05.52/checkpoint-84)... Done. 0.6s

No label_names provided for model class `PeftModelForCausalLM`. Since `PeftModel` hides base models input arguments, if label_names are provided, they will be passed to the base model. If label_names are not provided, the base model's input arguments will be used.

README.md: 100%

1.52k/1.52k [00:00<00:00, 93.3kB/s]

No files have been modified since last commit. Skipping to prevent empty commit.

WARNING:huggingface_hub.hf_api:No files have been modified since last commit. Skipping to prevent empty commit.

Model saved to HuggingFace as: ed_medical-2025-05-11_04.05.52

```
1 results = trainer.evaluate()
2 print(results)
3
```

[100/100 02:50]

```
{'eval loss': 1.4443223476409912, 'eval runtime': 172.7306, 'eval samples per second': 0.579, 'eval steps per second': 0.579}
```

```
1 wandb.finish()
```

**Run history:**

```
eval/loss          -
eval/runtime       -
eval/samples_per_second -
eval/steps_per_second -
train/epoch        [redacted]
train/global_step  [redacted]
train/grad_norm    -
train/learning_rate -
train/loss         -
train/mean_token_accuracy [redacted]
train/num_tokens   [redacted]
```

Run summary:

eval/loss	1.44432
eval/runtime	172.7306
eval/samples_per_second	0.579
eval/steps_per_second	0.579
total_flos	1.8483999762579456e+16
train/epoch	2.92444
train/global_step	84
train/grad_norm	0.26786
train/learning_rate	8e-05
train/loss	1.6934
train/mean_token_accuracy	0.66989
train/num_tokens	480088
train_loss	1.57578
train_runtime	10115.7149
train_samples_per_second	0.267
train_steps_per_second	0.008

View run [2025-05-11_04.05.52](https://wandb.ai/digitalblue-ai/ed_medical/runs/he1h23i) at: https://wandb.ai/digitalblue-ai/ed_medical/runs/he1h23i

View project at: https://wandb.ai/digitalblue-ai/ed_medical

Synced 5 W&B file(s), 0 media file(s), 17 artifact file(s) and 0 other file(s)

Find logs at: ./wandb/run-20250511_041757-he1h23i/logs

```
1 !zip -r ed_medical-2025-05-11_04.05.52.zip ed_medical-2025-05-11_04.05.52
```

```
[redacted] adding: ed_medical-2025-05-11_04.05.52/ (stored 0%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/ (stored 0%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/rng_state.pth (deflated 25%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/adapter_model.safetensors (deflated 7%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/training_args.bin (deflated 51%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/optimizer.pt (deflated 9%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/trainer_state.json (deflated 55%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/special_tokens_map.json (deflated 74%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/tokenizer.model (deflated 55%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/scheduler.pt (deflated 57%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/tokenizer.json (deflated 85%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/tokenizer_config.json (deflated 66%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/adapter_config.json (deflated 55%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/checkpoint-84/README.md (deflated 66%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/adapter_model.safetensors (deflated 7%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/training_args.bin (deflated 51%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/special_tokens_map.json (deflated 74%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/tokenizer.model (deflated 55%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/tokenizer.json (deflated 85%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/tokenizer_config.json (deflated 66%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/adapter_config.json (deflated 55%)
[redacted] adding: ed_medical-2025-05-11_04.05.52/README.md (deflated 46%)
```

```
1 !ls -la
```

```
[redacted] total 242888
drwxr-xr-x 1 root root      4096 May 11 07:10 .
drwxr-xr-x 1 root root      4096 May 11 03:49 ..
drwxr-xr-x 4 root root     4096 May  8 13:38 .config
drwx----- 6 root root     4096 May 11 03:50 drive
drwxr-xr-x 3 root root     4096 May 11 07:06 ed_medical-2025-05-11_04.05.52
-rw-r--r-- 1 root root 248680821 May 11 07:10 ed_medical-2025-05-11_04.05.52.zip
drwxr-xr-x 1 root root     4096 May  8 13:38 sample_data
drwxr-xr-x 3 root root     4096 May 11 04:17 wandb
```

```
1 !cp ed_medical-2025-05-11_04.05.52.zip /content/drive/MyDrive/
```

```
1 # stop notebook and disconnect GPU after finishing above steps as this process can take several hours
2 # from google.colab import runtime
3 # runtime.unassign()
```

⌄ 5. Evaluate the models - generate evaluation data

At this step the fine-tuned model created in the previous step can be re-loaded without need to re-run step 4. Generated responses will be collected, in addition to the base model to use as a benchmark.

The prompts to generate the responses are 50 samples from the **test** set, which was not used for fine-tuning.

The generated response for each model are then saved to HuggingFace hub for later analysis.

```
1 # define models saved to huggingface hub from previous step
2 MODEL_A = "digitalblue/ed_medical-2025-04-03_09.11.29" # trained on 200 rows
3 MODEL_B = "digitalblue/ed_medical-2025-04-05_10.59.54" # trained on 800 rows
4 MODEL_C = "digitalblue/ed_medical-2025-04-05_11.36.32" # trained on 1600 rows
5 MODEL_D = "digitalblue/ed_medical-2025-04-06_10.40.29" # trained on 1600 rows with NEFTune
6 MODEL_E = "digitalblue/ed_medical-2025-04-15_12.43.48" # trained on 8000 rows of pubmedqa only w/ NEFTune
7 MODEL_F = "digitalblue/ed_medical-2025-05-05_23.51.25" # trained on 9000 row, lower learning rate, 3 epochs
8 MODEL_G = "digitalblue/ed_medical-2025-05-08_01.23.37" # trained on 9000 rows, NEFTune false
9 MODEL_H = "digitalblue/ed_medical-2025-05-11_04.05.52" # trained on 900 row, with context, 3 epochs
```

```
1 # NOTE: Only using model E, F, G, H
```

```
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_base = AutoModelForCausalLM.from_pretrained(
2     model_name,
3     quantization_config=quant_config,
4     device_map="auto")
```

→ Loading checkpoint shards: 100% 2/2 [01:28<00:00, 40.27s/it]

```
1 # initialise tokenizer and pipeline
2 tokenizer = AutoTokenizer.from_pretrained(model_name, token=hf_token)
3 tokenizer.pad_token = tokenizer.eos_token
4 tokenizer.padding_side = "right"
```

```
1 #qa_test = load_dataset("digitalblue/ed_medical", split="test[:50]") # first 50 from test set
2 qa_test = load_dataset("digitalblue/ed_medical-pubmedqa-artificial", split="test[:50]") # load only pubmedqa YES rows, test set
```

```
1 # split into question and answer list
2 qa_list = []
3 for item in qa_test:
4     prompt = item['text']
5     question = re.search(r'\[INST\] (.*) \[/INST\]', prompt).group(1)
6     response = re.search(r'\[INST\] (.*) \</s\>', prompt).group(1)
7     qa_list.append([question, response])
```

```
1 qa_list[20][0] # question
```

→ 'Do n-glycans of human amniotic fluid transferrin stimulate progesterone production in human first trimester trophoblast cells in v

```
1 prompt1 = "You are a helpful personalised medical assistant. In your response do not include any personal names and keep your answer"
2 #prompt2 = "You are a helpful personalised medical assistant. Provide a friendly personal response to the patients medical question
3
4 def get_model_responses(model_x, qa_list):
5     model_responses = []
6     for item in qa_list:
7         question = item[0]
8         # print(item)
9         # print(question)
10        # print("----")
11        messages = [
```

```

12     {"role": "system", "content": prompt1},
13     {"role": "user", "content": question}
14 ]
15 inputs = tokenizer.apply_chat_template(messages, return_tensors="pt").to("cuda")
16 outputs = model_x.generate(inputs, max_new_tokens=250, temperature=0.1)
17 response = tokenizer.decode(outputs[0], skip_special_tokens=False)
18 answer = response.split('[/INST] ')[1] # get text after /INST token
19 answer = answer.replace('</s>', '') # remove trailing special token if present
20 model_responses.append(answer)
21 return model_responses

```

```

1 print(model_base.lm_head.weight[0, :10])
2 print(model_g.lm_head.weight[0, :10])

→ tensor([-0.0036,  0.0027, -0.0074, -0.0100,  0.0063,  0.0049, -0.0065,  0.0012,
         -0.0009,  0.0107], device='cuda:0', dtype=torch.float16)
tensor([-0.0036,  0.0027, -0.0074, -0.0100,  0.0063,  0.0049, -0.0065,  0.0012,
         -0.0009,  0.0107], device='cuda:0', dtype=torch.float16)

```

```
1 model_base_responses_pubmedqa = get_model_responses(model_base, qa_list)
```

```
1 prompt_base_pubmedqa_dataset = Dataset.from_dict({"text": model_base_responses_pubmedqa})
```

```
1 prompt_base_pubmedqa_dataset.push_to_hub(f"hf_username}/model_base_responses_pubmedqa", private=True)
```

→ Uploading the dataset shards: 100% 1/1 [00:01<00:00, 1.36s/it]

Creating parquet from Arrow format: 100% 1/1 [00:00<00:00, 110.81ba/s]

CommitInfo(commit_url='https://huggingface.co/datasets/digitalblue/model_base_responses_pubmedqa/commit/e7d5da1e7221601b91e115833b9; commit_message='Upload dataset', commit_description='', oid='e7d5da1e7221601b91e115833b91cd02cf1c0bd4', pr_url=None, reno_url=RenonUrl('https://huggingface.co/datasets/digitalblue/model_base_responses_pubmedqa').endpoint='https://huggingface.co').

```
1 prompt_base_pubmedqa_dataset[1]
```

→ {'text': 'Dose interval between 2 and 4 hours is recommended to achieve the highest bioavailability. '}

```
1 model_e = PeftModel.from_pretrained(model_base, MODEL_E)
2 model_f = PeftModel.from_pretrained(model_base, MODEL_F)
```

```
1 model_e_responses = get_model_responses(model_e, qa_list)
```

```
1 model_f_responses = get_model_responses(model_f, qa_list)
```

```
1 prompt_e_dataset = Dataset.from_dict({"text": model_e_responses})
2 prompt_f_dataset = Dataset.from_dict({"text": model_f_responses})
```

```
1 prompt_e_dataset.push_to_hub(f"hf_username}/model_e_responses", private=True)
2 prompt_f_dataset.push_to_hub(f"hf_username}/model_f_responses", private=True)
```

→ Uploading the dataset shards: 100% 1/1 [00:01<00:00, 1.37s/it]

Creating parquet from Arrow format: 100% 1/1 [00:00<00:00, 112.82ba/s]

Uploading the dataset shards: 100% 1/1 [00:01<00:00, 1.46s/it]

Creating parquet from Arrow format: 100% 1/1 [00:00<00:00, 117.08ba/s]

CommitInfo(commit_url='https://huggingface.co/datasets/digitalblue/model_f_responses/commit/8ac2484b0737fdde8a0f5f87d6471242afca9fd6; commit_message='Upload dataset', commit_description='', oid='8ac2484b0737fdde8a0f5f87d6471242afca9fd6', pr_url=None, reno_url=RenonUrl('https://huggingface.co/datasets/digitalblue/model_f_responses').endpoint='https://huggingface.co').

```
1 model_g = PeftModel.from_pretrained(model_base, MODEL_G)
```

→ adapter_config.json: 100% 850/850 [00:00<00:00, 38.3kB/s]

adapter_model.safetensors: 100% 793M/793M [00:34<00:00, 21.9MB/s]

```

1 model_g_responses = get_model_responses(model_g, qa_list)
2 prompt_g_dataset = Dataset.from_dict({"text": model_g_responses})
3 prompt_g_dataset.push_to_hub(f"hf_username}/model_g_responses", private=True)

```

```

Uploading the dataset shards: 100%                                         1/1 [00:01<00:00, 1.42s/it]
Creating parquet from Arrow format: 100%                                     1/1 [00:00<00:00, 78.17ba/s]
CommitInfo(commit_url='https://huggingface.co/datasets/digitalblue/model_g_responses/commit/18661dc8f0a4129ddbce079efe62700c8f909c5c', commit_message='Upload dataset', commit_description='', oid='18661dc8f0a4129ddbce079efe62700c8f909c5d', pr_url=None, renn_url=Renrollr('https://huggingface.co/datasets/digitalblue/model_g_responses'), endpoint='https://huggingface.co')
```

```

1 model_h = PeftModel.from_pretrained(model_base, MODEL_H)
2 model_h_responses = get_model_responses(model_h, qa_list)
3 prompt_h_dataset = Dataset.from_dict({"text": model_h_responses})
4 prompt_h_dataset.push_to_hub(f"hf_username/{model_h_responses}", private=True)
```

Uploading the dataset shards: 100% 1/1 [00:01<00:00, 1.43s/it]

Creating parquet from Arrow format: 100% 1/1 [00:00<00:00, 39.14ba/s]

```

CommitInfo(commit_url='https://huggingface.co/datasets/digitalblue/model_h_responses/commit/06fce8700b9d76ff82f9434c13399375f50ce198', commit_message='Upload dataset', commit_description='', oid='06fce8700b9d76ff82f9434c13399375f50ce198', pr_url=None, renn_url=Renrollr('https://huggingface.co/datasets/digitalblue/model_h_responses'), endpoint='https://huggingface.co').
```

```

1 # stop notebook and disconnect GPU after finishing above steps as this process can take several hours
2 from google.colab import runtime
3 runtime.unassign()
```

6. Evaluate fine tuned models - get metrics from generated responses

In this final step, the generated response can be reloaded from HuggingFace to evaluate and calculate metrics against the expected responses in **test** dataset.

```

1 qa_test = load_dataset("digitalblue/ed_medical-pubmedqa-artifical", split="test[:50]") # load only pubmedqa YES rows, test set

1 qa_list = []
2 for item in qa_test:
3   prompt = item['text']
4   question = re.search(r'\[INST\] (.*) \[/INST\]', prompt).group(1)
5   response = re.search(r'\[INST\] (.*) \</s\>', prompt).group(1)
6   qa_list.append([question, response])

1 qa_list[20][1] # index 1 is dataset real response
The results suggest that amnion-transferrin and especially its N-glycans modulate the endocrine function of trophoblasts in cultur
```

```

1 qa_references = [item[1] for item in qa_list] # get list of reference responses

1 # load saved response data to evaluate
2 model_base_pubmedqa_dataset = load_dataset("digitalblue/model_base_responses_pubmedqa")
3 model_e_dataset = load_dataset("digitalblue/model_e_responses")
4 model_f_dataset = load_dataset("digitalblue/model_f_responses")
5 model_g_dataset = load_dataset("digitalblue/model_g_responses")
6 model_h_dataset = load_dataset("digitalblue/model_h_responses")
```

```

 README.md: 100%                                         267/267 [00:00<00:00, 17.3kB/s]
 train-00000-of-00001.parquet: 100%                      11.1k/11.1k [00:00<00:00, 1.12MB/s]
 Generating train split: 100%                           100/100 [00:00<00:00, 1707.56 examples/s]
 README.md: 100%                                         267/267 [00:00<00:00, 27.1kB/s]
 train-00000-of-00001.parquet: 100%                      11.1k/11.1k [00:00<00:00, 1.11MB/s]
 Generating train split: 100%                           100/100 [00:00<00:00, 5452.25 examples/s]
 README.md: 100%                                         267/267 [00:00<00:00, 28.1kB/s]
 train-00000-of-00001.parquet: 100%                      11.3k/11.3k [00:00<00:00, 1.23MB/s]
 Generating train split: 100%                           100/100 [00:00<00:00, 4327.95 examples/s]
 README.md: 100%                                         266/266 [00:00<00:00, 27.9kB/s]
 train-00000-of-00001.parquet: 100%                      25.4k/25.4k [00:00<00:00, 2.31MB/s]
 Generating train split: 100%                           50/50 [00:00<00:00, 2815.42 examples/s]
 README.md: 100%                                         263/263 [00:00<00:00, 18.5kB/s]
 train-00000-of-00001.parquet: 100%                      5.44k/5.44k [00:00<00:00, 611kB/s]
 Generating train split: 100%                           50/50 [00:00<00:00, 2617.35 examples/s]

```

```

1 # convert to list
2 base_pubmedqa_response_list = model_base_pubmedqa_dataset['train']['text'][:50]
3 e_response_list = model_e_dataset['train']['text'][:50]
4 f_response_list = model_f_dataset['train']['text'][:50]
5 g_response_list = model_g_dataset['train']['text']
6 h_response_list = model_h_dataset['train']['text']

1 def calc_rouge_metric(references, predictions):
2     rouge = evaluate.load('rouge')
3     results = rouge.compute(
4         predictions=predictions,
5         references=references
6     )
7     return results
8
9 def calc_bleu_metric(references, predictions):
10    bleu = evaluate.load('bleu')
11    results = bleu.compute(
12        predictions=predictions,
13        references=references
14    )
15    return results
16
17 def calc_meteor_metric(references, predictions):
18    meteor = evaluate.load('meteor')
19    results = meteor.compute(
20        predictions=predictions,
21        references=[[ref] for ref in references]
22    )
23    return results
24
25 def calc_bert_score(references, predictions):
26    bert_scorer = BERTScorer(model_type='bert-base-uncased')
27    P = []
28    R = []
29    F1 = []
30
31    for i in range(len(predictions)):
32        p_i, r_i, f1_i = bert_scorer.score([predictions[i]], [references[i]])
33        P.append(p_i)
34        R.append(r_i)
35        F1.append(f1_i)
36    results = {
37        'P': np.mean(P),
38        'R': np.mean(R),
39        'F1': np.mean(F1)
40    }
41    return results

```

```

1 rouge_base_pubmedqa = calc_rouge_metric(qa_references, base_pubmedqa_response_list)
2 rouge_e = calc_rouge_metric(qa_references, e_response_list)
3 rouge_f = calc_rouge_metric(qa_references, f_response_list)
4 rouge_g = calc_rouge_metric(qa_references, g_response_list)
5 rouge_h = calc_rouge_metric(qa_references, h_response_list)

```

Downloading builder script: 100% 6.27k/6.27k [00:00<00:00, 640kB/s]

```

1 print(rouge_base_pubmedqa)
2 print(rouge_e)
3 print(rouge_f)
4 print(rouge_g)
5 print(rouge_h)

```

```

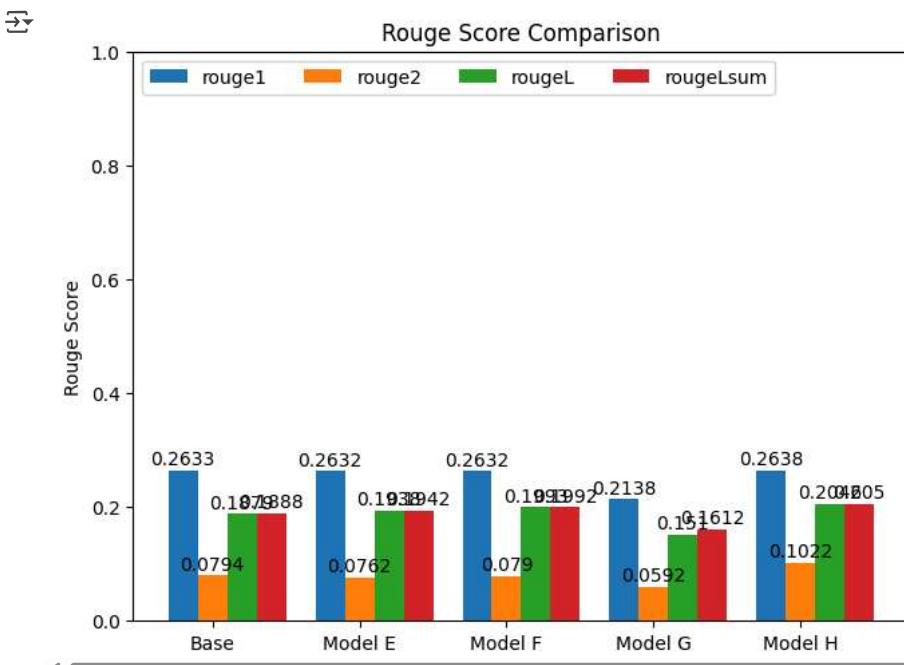
[2]: {'rouge1': np.float64(0.26328473964307153), 'rouge2': np.float64(0.07937912538217656), 'rougeL': np.float64(0.18786841702419282), 'rougeLsum': np.float64(0.19380980070682116), 'rougeLsum': np.float64(0.19380980070682116), 'rouge1': np.float64(0.263218379753081), 'rouge2': np.float64(0.07624340373229432), 'rougeL': np.float64(0.19380980070682116), 'rougeLsum': np.float64(0.19380980070682116), 'rouge1': np.float64(0.263197307016511), 'rouge2': np.float64(0.07898295649206828), 'rougeL': np.float64(0.19925342907701712), 'rougeLsum': np.float64(0.19925342907701712), 'rouge1': np.float64(0.2138327372000946), 'rouge2': np.float64(0.059210807794686446), 'rougeL': np.float64(0.15098541311199265), 'rougeLsum': np.float64(0.15098541311199265), 'rouge1': np.float64(0.2637618500720388), 'rouge2': np.float64(0.10223797199683642), 'rougeL': np.float64(0.2046431306067485), 'rougeLsum': np.float64(0.2046431306067485)}

```

```

1 # plot rouge scores
2 rouge_dict = {
3     'rouge1': (rouge_base_pubmedqa['rouge1'], rouge_e['rouge1'], rouge_f['rouge1'], rouge_g['rouge1'], rouge_h['rouge1']),
4     'rouge2': (rouge_base_pubmedqa['rouge2'], rouge_e['rouge2'], rouge_f['rouge2'], rouge_g['rouge2'], rouge_h['rouge2']),
5     'rougeL': (rouge_base_pubmedqa['rougeL'], rouge_e['rougeL'], rouge_f['rougeL'], rouge_g['rougeL'], rouge_h['rougeL']),
6     'rougeLsum': (rouge_base_pubmedqa['rougeLsum'], rouge_e['rougeLsum'], rouge_f['rougeLsum'], rouge_g['rougeLsum'], rouge_h['rougeLsum'])
7 }
8 model_labels = ('Base', 'Model E', 'Model F', 'Model G', 'Model H')
9
10 x = np.arange(len(model_labels))
11 width = 0.2
12 multiplier = 0
13
14 fig, ax = plt.subplots(layout='constrained')
15
16 for attribute, measurement in rouge_dict.items():
17     offset = width * multiplier
18     msr = np.round(measurement, 4) # round to 4 decimal places for plot
19     rects = ax.bar(x + offset, msr, width, label=attribute)
20     ax.bar_label(rects, padding=1)
21     multiplier += 1
22
23 ax.set_ylabel('Rouge Score')
24 ax.set_title('Rouge Score Comparison')
25 ax.set_xticks(x + width, model_labels)
26 ax.legend(loc='upper left', ncols=4)
27 ax.set_ylim(0, 1)
28
29 plt.show()

```



```

1 bleu_base_pubmedqa = calc_bleu_metric(qa_references, base_pubmedqa_response_list)
2 bleu_e = calc_bleu_metric(qa_references, e_response_list)

```

```
3 bleu_f = calc_bleu_metric(qa_references, f_response_list)
4 bleu_g = calc_bleu_metric(qa_references, g_response_list)
5 bleu_h = calc_bleu_metric(qa_references, h_response_list)
```

→ Downloading builder script: 100% 5.94k/5.94k [00:00<00:00, 106kB/s]

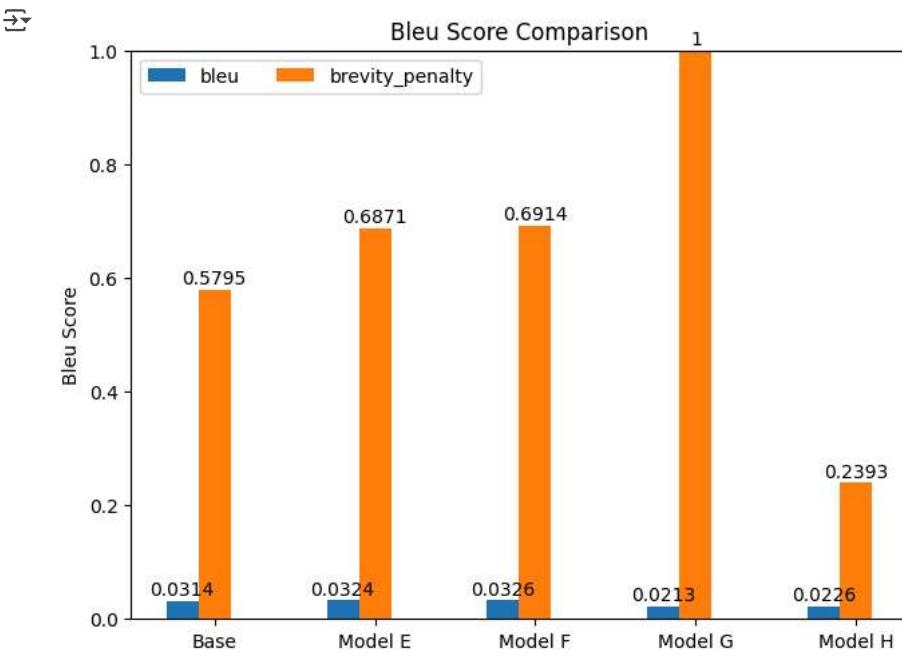
Downloading extra modules: 4.07k/? [00:00<00:00, 99.1kB/s]

Downloading extra modules: 100% 3.34k/3.34k [00:00<00:00, 356kB/s]

```
1 print(bleu_base_pubmedqa)
2 print(bleu_e)
3 print(bleu_f)
4 print(bleu_g)
5 print(bleu_h)

→ {'bleu': 0.03138988119399103, 'precisions': [0.3380281690140845, 0.07929176289453425, 0.029623698959167333, 0.010842368640533779], 'bleu': 0.032379811967089916, 'precisions': [0.3113456464379947, 0.07094133697135062, 0.025423728813559324, 0.008784773060029283], 'bleu': 0.0326166306491595, 'precisions': [0.3112278397898884, 0.07196198234894773, 0.025298664792691498, 0.00873998543335761], 'bleu': 0.021269429662351482, 'precisions': [0.1367226770070372, 0.032402719513959205, 0.01049103890426927, 0.004403346543372964], 'bleu': 0.022637121921975846, 'precisions': [0.40326340326340326, 0.12128712871287128, 0.052770448548812667, 0.031029619181946404],
```

```
1 # plot bleu scores
2 bleu_dict = {
3     'bleu': (bleu_base_pubmedqa['bleu'], bleu_e['bleu'], bleu_f['bleu'], bleu_g['bleu'], bleu_h['bleu']),
4     'brevity_penalty': (bleu_base_pubmedqa['brevity_penalty'], bleu_e['brevity_penalty'], bleu_f['brevity_penalty'], bleu_g['brevity'],
5 } 6
7 model_labels = ('Base', 'Model E', 'Model F', 'Model G', 'Model H')
8
9 x = np.arange(len(model_labels))
10 width = 0.2
11 multiplier = 0
12
13 fig, ax = plt.subplots(layout='constrained')
14
15 for attribute, measurement in bleu_dict.items():
16     offset = width * multiplier
17     msr = np.round(measurement, 4) # round to 4 decimal places for plot
18     rects = ax.bar(x + offset, msr, width, label=attribute)
19     ax.bar_label(rects, padding=1)
20     multiplier += 1
21
22 ax.set_ylabel('Bleu Score')
23 ax.set_title('Bleu Score Comparison')
24 ax.set_xticks(x + width, model_labels)
25 ax.legend(loc='upper left', ncols=4)
26 ax.set_ylim(0, 1)
27
28 plt.show()
```



```
1 meteor_base_pubmedqa = calc_meteor_metric(qa_references, base_pubmedqa_response_list)
2 meteor_e = calc_meteor_metric(qa_references, e_response_list)
```

```
3 meteor_f = calc_meteor_metric(qa_references, f_response_list)
4 meteor_g = calc_meteor_metric(qa_references, g_response_list)
5 meteor_h = calc_meteor_metric(qa_references, h_response_list)
```

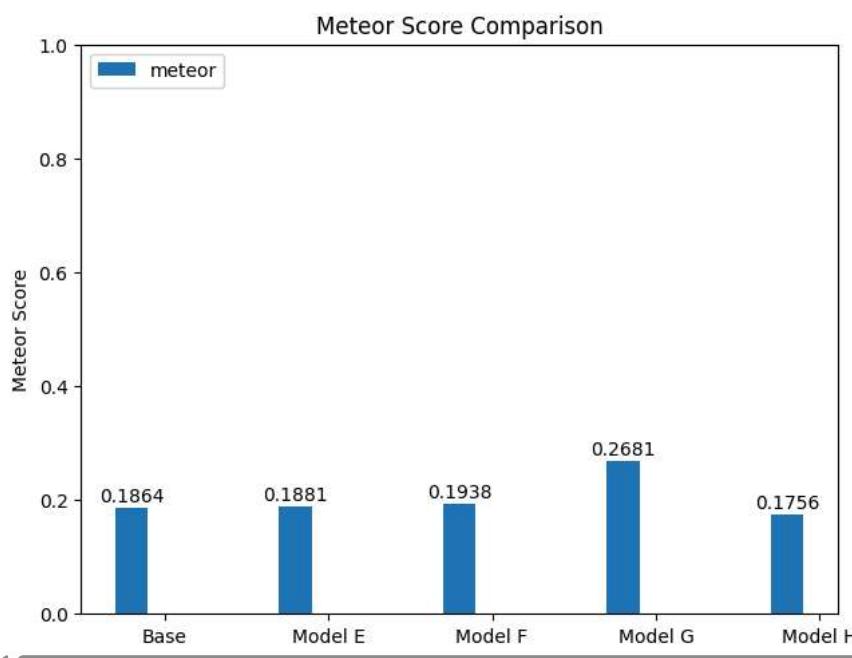
→ Downloading builder script: 100% 7.02k/7.02k [00:00<00:00, 481kB/s]

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]  Unzipping tokenizers/punkt_tab.zip.
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]  Package wordnet is already up-to-date!
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]  Package punkt_tab is already up-to-date!
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data]  Package omw-1.4 is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]  Package wordnet is already up-to-date!
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]  Package punkt_tab is already up-to-date!
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data]  Package omw-1.4 is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]  Package wordnet is already up-to-date!
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]  Package punkt_tab is already up-to-date!
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data]  Package omw-1.4 is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]  Package wordnet is already up-to-date!
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]  Package punkt_tab is already up-to-date!
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data]  Package omw-1.4 is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]  Package wordnet is already up-to-date!
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]  Package punkt_tab is already up-to-date!
[nltk_data] Downloading package omw-1.4 to /root/nltk_data...
[nltk_data]  Package omw-1.4 is already up-to-date!
```

```
1 print(meteor_base_pubmedqa)
2 print(meteor_e)
3 print(meteor_f)
4 print(meteor_g)
5 print(meteor_h)
```

→ {'meteor': np.float64(0.18639960354363067)}
{'meteor': np.float64(0.1881447693596208)}
{'meteor': np.float64(0.19384352470409308)}
{'meteor': np.float64(0.26814099277140213)}
{'meteor': np.float64(0.17556055146287494)}

```
1 # plot meteor scores
2 meteor_dict = {
3     'meteor': (meteor_base_pubmedqa['meteor'], meteor_e['meteor'], meteor_f['meteor'], meteor_g['meteor'], meteor_h['meteor'])
4 }
5 model_labels = ('Base', 'Model E', 'Model F', 'Model G', 'Model H')
6
7 x = np.arange(len(model_labels))
8 width = 0.2
9 multiplier = 0
10
11 fig, ax = plt.subplots(layout='constrained')
12
13 for attribute, measurement in meteor_dict.items():
14     offset = width * multiplier
15     msr = np.round(measurement, 4) # round to 4 decimal places for plot
16     rects = ax.bar(x + offset, msr, width, label=attribute)
17     ax.bar_label(rects, padding=1)
18     multiplier += 1
19
20 ax.set_ylabel('Meteor Score')
21 ax.set_title('Meteor Score Comparison')
22 ax.set_xticks(x + width, model_labels)
23 ax.legend(loc='upper left', ncols=4)
24 ax.set_ylim(0, 1)
25
26 plt.show()
```



```

1 bert_score_base_pubmedqa = calc_bert_score(qa_references, base_pubmedqa_response_list)
2 bert_score_e = calc_bert_score(qa_references, e_response_list)
3 bert_score_f = calc_bert_score(qa_references, f_response_list)
4 bert_score_g = calc_bert_score(qa_references, g_response_list)
5 bert_score_h = calc_bert_score(qa_references, h_response_list)

```

```

→ tokenizer_config.json: 100% 48.0/48.0 [00:00<00:00, 1.14kB/s]
config.json: 100% 570/570 [00:00<00:00, 16.0kB/s]
vocab.txt: 100% 232k/232k [00:00<00:00, 19.6MB/s]
tokenizer.json: 100% 466k/466k [00:00<00:00, 2.15MB/s]

Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better p
WARNING:huggingface_hub.file_download:Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back
model.safetensors: 100% 440M/440M [00:01<00:00, 273MB/s]

```

```

1 print(bert_score_base_pubmedqa)
2 print(bert_score_e)
3 print(bert_score_f)
4 print(bert_score_g)
5 print(bert_score_h)

→ {'P': np.float32(0.6467173), 'R': np.float32(0.58598), 'F1': np.float32(0.6117929)}
{'P': np.float32(0.64296514), 'R': np.float32(0.585245), 'F1': np.float32(0.61016136)}
{'P': np.float32(0.64098334), 'R': np.float32(0.58719736), 'F1': np.float32(0.6097659)}
{'P': np.float32(0.5433609), 'R': np.float32(0.6343017), 'F1': np.float32(0.58330375)}
{'P': np.float32(0.6551912), 'R': np.float32(0.56881934), 'F1': np.float32(0.60660475)}

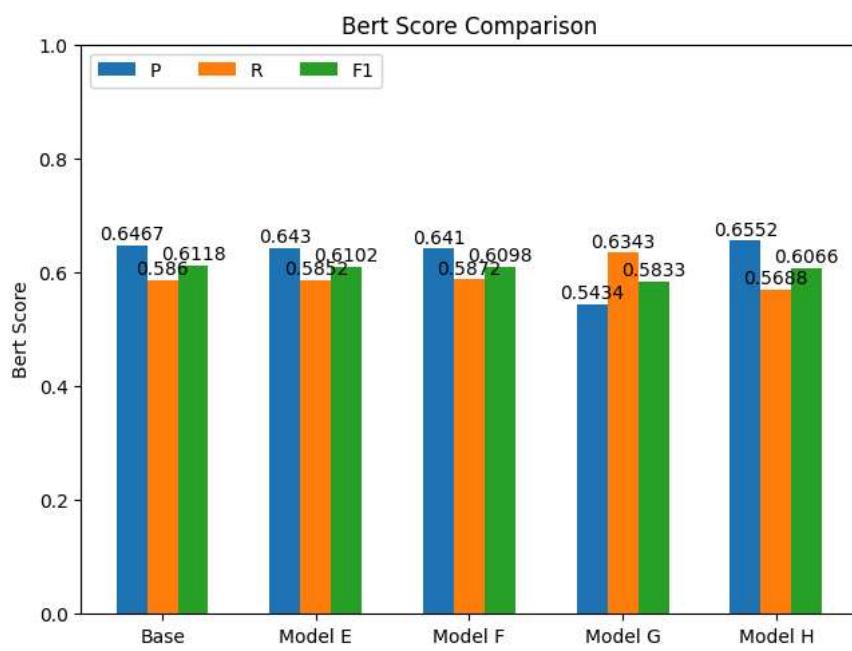
1 # plot bert scores
2 bert_score_dict = {
3     'P': (bert_score_base_pubmedqa['P'], bert_score_e['P'], bert_score_f['P'], bert_score_g['P'], bert_score_h['P']),
4     'R': (bert_score_base_pubmedqa['R'], bert_score_e['R'], bert_score_f['R'], bert_score_g['R'], bert_score_h['R']),
5     'F1': (bert_score_base_pubmedqa['F1'], bert_score_e['F1'], bert_score_f['F1'], bert_score_g['F1'], bert_score_h['F1'])
6 }
7 model_labels = ('Base', 'Model E', 'Model F', 'Model G', 'Model H')
8
9 x = np.arange(len(model_labels))
10 width = 0.2
11 multiplier = 0
12
13 fig, ax = plt.subplots(layout='constrained')
14
15 for attribute, measurement in bert_score_dict.items():
16     offset = width * multiplier
17     msr = np.round(measurement, 4) # round to 4 decimal places for plot
18     rects = ax.bar(x + offset, msr, width, label=attribute)
19     ax.bar_label(rects, padding=1)
20     multiplier += 1
21
22 ax.set_ylabel('Bert Score')
23 ax.set_title('Bert Score Comparison')

```

```

24 ax.set_xticks(x + width, model_labels)
25 ax.legend(loc='upper left', ncols=4)
26 ax.set_xlim(0, 1)
27
28 plt.show()

```



```

1 # human evaluation
2 def print_response(QA_NUM):
3     for i in range(QA_NUM):
4         print('-----')
5         print('\nTest question: ' + qa_list[i][0])
6         print('\nTest answer: ' + qa_list[i][1])
7         print('\nBase: ' + base_pubmedqa_response_list[i])
8         print('\nModel E: ' + e_response_list[i])
9         print('\nModel F: ' + f_response_list[i])
10        print('\nModel G: ' + g_response_list[i])
11        print('\nModel H: ' + h_response_list[i])
12
13 print_response(5)

```

Test question: Does simvastatin attenuate hepatic sensitization to lipopolysaccharide after partial hepatectomy?

Test answer: These novel findings demonstrate that simvastatin protects the remnant liver against endotoxemic injury following major hepatectomy. The results suggest that simvastatin attenuates the increase in the expression of CD14 and TNF-alpha after partial hepatectomy.

Base: The results suggest that simvastatin attenuates the increase in mRNA expression of TNF-alpha and IL-1beta in the liver after partial hepatectomy.

Model E: The results suggest that simvastatin attenuates the increase in mRNA expression of TNF-alpha and IL-1beta in the liver after partial hepatectomy.

Model F: The results suggest that simvastatin attenuates the increase in mRNA expression of TNF-alpha and IL-1beta in the liver after partial hepatectomy.

Model G: Yes, simvastatin has been shown to attenuate hepatic sensitization to lipopolysaccharide (LPS) after partial hepatectomy. Hepatic sensitization to LPS is a common complication after partial hepatectomy, leading to increased inflammation and impaired liver function. Studies have shown that simvastatin can reduce the expression of toll-like receptor 4 (TLR4) and myeloid differentiation factor 88 (MDRF88).

Overall, the available evidence suggests that simvastatin may be beneficial in preventing hepatic sensitization after partial hepatectomy.

Model H: Simvastatin attenuates hepatic sensitization to lipopolysaccharide after partial hepatectomy.

Test question: Does mizoribine require individual dosing due to variation of bioavailability?

Test answer: Individual dosing is required to optimize C(max) in childhood-onset glomerular disease patients. The safe dosage of mizoribine is approximately 10 mg/m²/day, divided into two equal doses given 12 hours apart. Dose interval between 2 and 4 hours is recommended to achieve the highest bioavailability.

Base: Dose interval between 2 and 4 hours is recommended to achieve the highest bioavailability.

Model E: Dose interval between 2 and 4 hours is recommended to achieve the highest bioavailability.

Model F: Dear patient,

Model G: Yes, mizoribine can exhibit significant inter-individual variability in bioavailability, which may require individualized dosing.

Model H: Mizoribine is a drug that requires individual dosing due to variation of bioavailability.

Test question: Does cold reperfusion before rewarming reduce neurological events after deep hypothermic circulatory arrest?

Test answer: sDHCA remains a safe and easy tool for cerebral protection when DHCA duration is expected to be less than 30 min. When

Base: This study suggests that cold reperfusion followed by rewarming may reduce the incidence of neurological events after DHCA.

Model E: This study suggests that cold reperfusion followed by rewarming may reduce neurological events after DHCA.

Model F: Cold reperfusion is associated with a significant reduction in neurological events after deep hypothermic circulatory arr

Model G: Cold reperfusion before rewarming has been shown to reduce the incidence of neurological events after deep hypothermic c

Cold reperfusion, which involves rewarming the body at a slow rate, has been found to reduce the risk of neurological events such

However, it is important to note that the evidence for the benefit of cold reperfusion is not yet definitive, and more research is

Model H: The results of this study suggest that cold reperfusion before rewarming does not reduce the incidence of neurological ev

▼ 7. Evaluation Summary

Evaluation metrics do not show improvement over base model. The Llama 2 model does not disclose its training dataset, however as PubMedQA is public data, there exists the possibility that it may have already been trained on this dataset.

Alternatively, the training hyperparameters used in training and/or the volume of data used may have been insufficient for the model to successfully learn the domain.

Next steps:

- Continue to experiment with further hyperparameter tuning
- Increase volume of training data
- Adjust system prompt, and model parameters on inference
- Explore different dataset

▼ 8. Model packaging

This section explored how to bundle the fine-tuned model with adapters and including tokenizer, into an LLM that can be downloaded and used standalone. The intention was to investigate how it could be deployed to a cloud provider, and used on edge device or local PC via Ollama.

```
1 # clear gpu memory
2 torch.cuda.empty_cache()
3 torch.cuda.reset_max_memory_allocated()

[2]: /usr/local/lib/python3.11/dist-packages/torch/cuda/memory.py:391: FutureWarning: torch.cuda.reset_max_memory_allocated now calls torch._C._cudaResetMaxMemoryAllocated. To silence this warning, use warnings.warn()
      warnings.warn(
      └────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────────................................................................
```

```
↳ config.json: 100% 614/614 [00:00<00:00, 45.7kB/s]
model.safetensors.index.json: 100% 26.8k/26.8k [00:00<00:00, 2.64MB/s]
Fetching 2 files: 100% 2/2 [01:16<00:00, 76.74s/it]
model-00001-of-00002.safetensors: 100% 9.98G/9.98G [01:16<00:00, 87.7MB/s]
model-00002-of-00002.safetensors: 100% 3.50G/3.50G [00:38<00:00, 141MB/s]
Loading checkpoint shards: 100% 2/2 [00:58<00:00, 26.71s/it]
generation_config.json: 100% 188/188 [00:00<00:00, 9.34kB/s]
WARNING:accelerate.hip_modeling:Some parameters are on the meta device because they were offloaded to the cpu and disk.
```

```
1 # initialise tokenizer and pipeline
2 tokenizer = AutoTokenizer.from_pretrained(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, 84.0kB/s]
tokenizer.model: 100% 500k/500k [00:00<00:00, 25.5MB/s]
tokenizer.json: 100% 1.84M/1.84M [00:00<00:00, 4.13MB/s]
special_tokens_map.json: 100% 414/414 [00:00<00:00, 45.4kB/s]
```

```
1 MODEL_A = "digitalblue/ed_medical-2025-04-03_09.11.29" # trained on 200 rows
2 MODEL_B = "digitalblue/ed_medical-2025-04-05_10.59.54" # trained on 800 rows
3 MODEL_C = "digitalblue/ed_medical-2025-04-05_11.36.32" # trained on 1600 rows
4 MODEL_D = "digitalblue/ed_medical-2025-04-06_10.40.29" # trained on 1600 rows with NEFTune
5 MODEL_E = "digitalblue/ed_medical-2025-04-15_12.43.48" # trained on 8000 rows of pubmedqa only w/ NEFTune
6 MODEL_F = "digitalblue/ed_medical-2025-05-05_23.51.25" # trained on 9000 row, lower learning rate, 3 epochs
```

```
1 model_e = PeftModel.from_pretrained(model, MODEL_E)
```

```
↳ adapter_config.json: 100% 754/754 [00:00<00:00, 56.4kB/s]
adapter_model.safetensors: 100% 134M/134M [00:00<00:00, 247MB/s]
```

```
1 print(model.lm_head.weight[0, :10])
2 print(model_e.lm_head.weight[0, :10])
3 print(model_f.lm_head.weight[0, :10])
```

```
↳ tensor([-0.0115,  0.0007, -0.0078, -0.0107,  0.0119,  0.0019,  0.0006, -0.0050,
         -0.0156,  0.0156], device='cuda:0', dtype=torch.float16,
         grad_fn=<SliceBackward0>)
tensor([-0.0115,  0.0007, -0.0078, -0.0107,  0.0119,  0.0019,  0.0006, -0.0050,
         -0.0156,  0.0156], device='cuda:0', dtype=torch.float16,
         grad_fn=<SliceBackward0>)
tensor([-0.0115,  0.0007, -0.0078, -0.0107,  0.0119,  0.0019,  0.0006, -0.0050,
         -0.0156,  0.0156], device='cuda:0', dtype=torch.float16,
         grad_fn=<SliceBackward0>)
```

```
1 model_e.print_trainable_parameters()
2 model_f.print_trainable_parameters()
```

```
↳ trainable params: 262,144,000 || all params: 7,008,948,224 || trainable%: 3.7401
trainable params: 262,144,000 || all params: 7,008,948,224 || trainable%: 3.7401
```

```
1 messages = [
2     {"role": "system", "content": "You are a helpful personalised medical assistant. In your response do not include any personal names or details."}
3     {"role": "user", "content": "Does methylphenidate improve the quality of life of children and adolescents with ADHD and difficulty concentrating?"}
4 ]
```

```
1 og_output = model_e.generate(tokenizer.apply_chat_template(messages, return_tensors="pt").to("cuda"), max_new_tokens=250)
```

```
↳ The attention mask and the pad token id were not set. As a consequence, you may observe unexpected behavior. Please pass your input Setting `pad_token_id` to `eos_token_id` for open-end generation.
The attention mask is not set and cannot be inferred from input because pad token is same as eos token. As a consequence, you may obtain None for the attention mask.
```

```
1 tokenizer.decode(og_output[0], skip_special_tokens=False)
```

```
↳ '<s> [INST] <> You are a helpful personalised medical assistant. In your response do not include any personal names and keep your answer professional and concise.<> <>\n\nDoes methylphenidate improve the quality of life of children and adolescents with ADHD and difficulty concentrating?'
```

```
1 base_output = model.generate(tokenizer.apply_chat_template(messages, return_tensors="pt").to("cuda"), max_new_tokens=250)
2 tokenizer.decode(base_output[0], skip_special_tokens=False)

→ The attention mask and the pad token id were not set. As a consequence, you may observe unexpected behavior. Please pass your input
Setting `pad_token_id` to `eos_token_id`:2 for open-end generation.
'<s> [INST] <>SYS>\nYou are a helpful personalised medical assistant. In your response do not include any personal names and keep
your answer professional and concise.\n</>SYS>\n\nDoes methyphenidate improve the quality of life of children and adolescents wit
1 merge_model_e = model_e.merge_and_unload()

→ /usr/local/lib/python3.11/dist-packages/peft/tuners/lora/bnb.py:351: UserWarning: Merge lora module to 4-bit linear may get differer
warnings.warn(
```



```
1 merge_model_f = model_f.merge_and_unload()

1 merge_model_f
```

```
→ LlamaForCausalLM(
    (model): LlamaModel(
Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.
```