

Krystal_Falcon7B_pubmedqa

April 9, 2025

1 Fine Tuning Falcon-7B model on the PubMedQA dataset

This notebook tests the Falcon-7B model on the PubMedQA dataset to answer biomedical questions using provided contexts. Leverage PEFT library from Hugging Face ecosystem, as well as QLoRA for more memory efficient finetuning.

We evaluate the first 10 samples (indices 0-9) and use a lightweight DistilBERT model to judge the responses for correctness, evidence alignment, and clarity. The process includes generating answers, scoring them, and calculating metrics like accuracy, BERTScore, and ROUGE, all optimized for a T4 GPU setup.

1.1 Setup

Run the cells below to setup and install the required libraries.

```
[ ]: !pip install -qU bitsandbytes transformers datasets accelerate loralib einops
      ↪xformers
!pip install -q -U git+https://github.com/huggingface/peft.git

import os
import bitsandbytes as bnb
import pandas as pd
import torch
import torch.nn as nn
import transformers
from datasets import load_dataset
from peft import (
    LoraConfig,
    PeftConfig,
    get_peft_model,
    prepare_model_for_kbit_training,
)
from transformers import (
    AutoConfig,
    AutoModelForCausalLM,
    AutoTokenizer,
    BitsAndBytesConfig,
)
```

```
Installing build dependencies ... done
Getting requirements to build wheel ... done
Preparing metadata (pyproject.toml) ... done
```

1.2 Loading the Pre-Trained Model

```
[ ]: model_id = "tiiuae/falcon-7b"

# Configure for 8-bit quantization
bnb_config = BitsAndBytesConfig(
    load_in_8bit=True,
)

model = AutoModelForCausalLM.from_pretrained(
    model_id,
    device_map="auto",
    trust_remote_code=True,
    quantization_config=bnb_config,
)

tokenizer = AutoTokenizer.from_pretrained(model_id)
tokenizer.pad_token = tokenizer.eos_token
```

```
/usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94:
UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab
(https://huggingface.co/settings/tokens), set it as secret in your Google Colab
and restart your session.
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access
public models or datasets.
```

```
warnings.warn(
WARNING:transformers_modules.tiiuae.falcon-
7b.ec89142b67d748a1865ea4451372db8313ada0d8.configuration_falcon:
WARNING: You are currently loading Falcon using legacy code contained in the
model repository. Falcon has now been fully ported into the Hugging Face
transformers library. For the most up-to-date and high-performance version of
the Falcon model code, please update to the latest version of transformers and
then load the model without the trust_remote_code=True argument.
```

```
Loading checkpoint shards: 0%|          | 0/2 [00:00<?, ?it/s]
```

1.3 Configuring LoRA

```
[ ]: # Prepare model for LoRA fine-tuning
model = prepare_model_for_kbit_training(model)

# Configure LoRA
lora_alpha = 32 # scaling factor for the weight matrices
lora_dropout = 0.05 # dropout probability of the LoRA layers
lora_rank = 32 # dimension of the low-rank matrices

peft_config = LoraConfig(
    lora_alpha=lora_alpha,
    lora_dropout=lora_dropout,
    r=lora_rank,
    bias="none",
    task_type="CAUSAL_LM",
    target_modules=[
        # Setting names of modules in falcon-7b model that we want to apply
        ↪LoRA to
        "query_key_value",
        "dense",
        "dense_h_to_4h",
        "dense_4h_to_h",
    ]
)

peft_model = get_peft_model(model, peft_config)
```

You are using an old version of the checkpointing format that is deprecated (We will also silently ignore ``gradient_checkpointing_kwargs`` in case you passed it). Please update to the new format on your modeling file. To use the new format, you need to completely remove the definition of the method ``_set_gradient_checkpointing`` in your model.

1.4 Loading and Preparing the Dataset

```
[27]: # Load PubMedQA Labeled Dataset
dataset = load_dataset("qiaojin/PubMedQA", "pqa_labeled", split="train")
print(f"Dataset size: {len(dataset)}")

# Inspect a few examples
print("\nSample Data Examples:")
for i in range(10):
    print(f"\nExample {i+1}:")
    print(f"Question: {dataset[i]['question']}")
    # Access the 'context' as a string before slicing
    context = " ".join(dataset[i]['context']['contexts'])
    print(f"Context: {context[:200]}...") # Truncate context for brevity
```

```
print(f"Long Answer: {dataset[i]['long_answer']}")
print(f"Final Decision: {dataset[i]['final_decision']}")
```

Dataset size: 1000

Sample Data Examples:

Example 1:

Question: Do mitochondria play a role in remodelling lace plant leaves during programmed cell death?

Context: Programmed cell death (PCD) is the regulated death of cells within an organism. The lace plant (*Aponogeton madagascariensis*) produces perforations in its leaves through PCD. The leaves of the plant co...

Long Answer: Results depicted mitochondrial dynamics in vivo as PCD progresses within the lace plant, and highlight the correlation of this organelle with other organelles during developmental PCD. To the best of our knowledge, this is the first report of mitochondria and chloroplasts moving on transvacuolar strands to form a ring structure surrounding the nucleus during developmental PCD. Also, for the first time, we have shown the feasibility for the use of CSA in a whole plant system. Overall, our findings implicate the mitochondria as playing a critical and early role in developmentally regulated PCD in the lace plant.

Final Decision: yes

Example 2:

Question: Landolt C and snellen e acuity: differences in strabismus amblyopia?

Context: Assessment of visual acuity depends on the optotypes used for measurement. The ability to recognize different optotypes differs even if their critical details appear under the same visual angle. Since...

Long Answer: Using the charts described, there was only a slight overestimation of visual acuity by the Snellen E compared to the Landolt C, even in strabismus amblyopia. Small differences in the lower visual acuity range have to be considered.

Final Decision: no

Example 3:

Question: Syncope during bathing in infants, a pediatric form of water-induced urticaria?

Context: Apparent life-threatening events in infants are a difficult and frequent problem in pediatric practice. The prognosis is uncertain because of risk of sudden infant death syndrome. Eight infants aged 2...

Long Answer: "Aquagenic maladies" could be a pediatric form of the aquagenic urticaria.

Final Decision: yes

Example 4:

Question: Are the long-term results of the transanal pull-through equal to those of the transabdominal pull-through?

Context: The transanal endorectal pull-through (TERPT) is becoming the most popular procedure in the treatment of Hirschsprung disease (HD), but overstretching of the anal sphincters remains a critical issue t...

Long Answer: Our long-term study showed significantly better (2-fold) results regarding the continence score for the abdominal approach compared with the transanal pull-through. The stool pattern and enterocolitis scores were somewhat better for the TERPT group. These findings raise an important issue about the current surgical management of HD; however, more cases will need to be studied before a definitive conclusion can be drawn.

Final Decision: no

Example 5:

Question: Can tailored interventions increase mammography use among HMO women?

Context: Telephone counseling and tailored print communications have emerged as promising methods for promoting mammography screening. However, there has been little research testing, within the same randomize...

Long Answer: The effects of the intervention were most pronounced after the first intervention. Compared to usual care, telephone counseling seemed particularly effective at promoting change among nonadherent women, the group for whom the intervention was developed. These results suggest that telephone counseling, rather than tailored print, might be the preferred first-line intervention for getting nonadherent women on schedule for mammography screening. Many questions would have to be answered about why the tailored print intervention was not more powerful. Nevertheless, it is clear that additional interventions will be needed to maintain women's adherence to mammography.

Medical Subject Headings (MeSH): mammography screening, telephone counseling, tailored print communications, barriers.

Final Decision: yes

Example 6:

Question: Double balloon enteroscopy: is it efficacious and safe in a community setting?

Context: From March 2007 to January 2011, 88 DBE procedures were performed on 66 patients. Indications included evaluation anemia/gastrointestinal bleed, small bowel IBD and dilation of strictures. Video-capsu...

Long Answer: DBE appears to be equally safe and effective when performed in the community setting as compared to a tertiary referral center with a comparable yield, efficacy, and complication rate.

Final Decision: yes

Example 7:

Question: 30-Day and 1-year mortality in emergency general surgery laparotomies: an area of concern and need for improvement?

Context: Emergency surgery is associated with poorer outcomes and higher mortality with recent studies suggesting the 30-day mortality to be 14-15%. The aim of this study was to analyse the 30-day mortality, a...

Long Answer: Emergency laparotomy carries a high rate of mortality, especially in those over the age of 70 years, and more needs to be done to improve

outcomes, particularly in this group. This could involve increasing acute surgical care manpower, early recognition of patients requiring emergency surgery, development of clear management protocols for such patients or perhaps even considering centralisation of emergency surgical services to specialist centres with multidisciplinary teams involving emergency surgeons and care of the elderly physicians in hospital and related community outreach services for post-discharge care.

Final Decision: maybe

Example 8:

Question: Is adjustment for reporting heterogeneity necessary in sleep disorders?

Context: Anchoring vignettes are brief texts describing a hypothetical character who illustrates a certain fixed level of a trait under evaluation. This research uses vignettes to elucidate factors associated ...

Long Answer: Sleep disorders are common in the general adult population of Japan. Correction for reporting heterogeneity using anchoring vignettes is not a necessary tool for proper management of sleep and energy related problems among Japanese adults. Older age, gender differences in communicating sleep-related problems, the presence of multiple morbidities, and regular exercise should be the focus of policies and clinical practice to improve sleep and energy management in Japan.

Final Decision: no

Example 9:

Question: Do mutations causing low HDL-C promote increased carotid intima-media thickness?

Context: Although observational data support an inverse relationship between high-density lipoprotein (HDL) cholesterol and coronary heart disease (CHD), genetic HDL deficiency states often do not correlate wi...

Long Answer: Genetic variants identified in the present study may be insufficient to promote early carotid atherosclerosis.

Final Decision: no

Example 10:

Question: A short stay or 23-hour ward in a general and academic children's hospital: are they effective?

Context: We evaluated the usefulness of a short stay or 23-hour ward in a pediatric unit of a large teaching hospital, Westmead Hospital, and an academic Children's hospital, The New Children's Hospital, to de...

Long Answer: This data demonstrates the robust nature of the short stay ward. At these two very different institutions we have shown improved bed efficient and patient care in a cost-effective way. We have also reported on greater parental satisfaction and early return of the child with their family to the community.

Final Decision: yes

```
[28]: def generate_prompt(data_point):
    PROMPT_TEMPLATE = """<|system|>You are a helpful medical assistant.
    <|endoftext|>
    <|user|>Question: {question}
    Context: {context}<|endoftext|>
    <|assistant|>Answer: {answer}
    Final Decision: {decision}<|endoftext|>"""
    return PROMPT_TEMPLATE.format(
        question=data_point["question"],
        context=data_point["context"],
        answer=data_point["long_answer"],
        decision=data_point["final_decision"]
    )

def generate_and_tokenize_prompt(data_point):
    full_prompt = generate_prompt(data_point)
    return tokenizer(full_prompt, padding=True, truncation=True, max_length=384)

dataset = dataset.shuffle(seed=42).map(
    generate_and_tokenize_prompt,
)
train_dataset = dataset.select(range(900))
test_dataset = dataset.select(range(900, 1000))
```

```
Map:   0%|          | 0/1000 [00:00<?, ? examples/s]
```

1.5 Setting Up the Training Arguments

```
[ ]: # Training Arguments
OUTPUT_DIR = "/falcon-7b-pubmedqa"
if not os.path.exists(OUTPUT_DIR):
    os.makedirs(OUTPUT_DIR)
training_args = transformers.TrainingArguments(
    auto_find_batch_size=True,
    per_device_train_batch_size=4,
    num_train_epochs=1,
    learning_rate=2e-4,
    fp16=True,
    save_total_limit=2,
    logging_steps=10,
    save_strategy="steps",
    save_steps=200,
    max_steps=-1,
    gradient_checkpointing=True,
    optim="adamw_torch_fused",
    warmup_ratio=0.03,
    lr_scheduler_type="cosine",
```

```

report_to="none",
output_dir=OUTPUT_DIR
)

```

1.6 Model Training

```

[ ]: # 8. Trainer setup
trainer = transformers.Trainer(
    model=peft_model,
    train_dataset=train_dataset,
    eval_dataset=test_dataset,
    args=training_args,
    data_collator=transformers.DataCollatorForLanguageModeling(tokenizer,
↪mlm=False),
    callbacks=[MemoryMonitorCallback()]
)

# 9. Enable model caching to improve performance
peft_model.config.use_cache = False

```

No `label_names` provided for model class ``PeftModelForCausalLM``. Since ``PeftModel`` hides base models input arguments, if `label_names` is not given, `label_names` can't be set automatically within ``Trainer``. Note that empty `label_names` list will be used instead.

You are using an old version of the checkpointing format that is deprecated (We will also silently ignore ``gradient_checkpointing_kwargs`` in case you passed it). Please update to the new format on your modeling file. To use the new format, you need to completely remove the definition of the method ``_set_gradient_checkpointing`` in your model.

Step 0 Memory: 8.11GB allocated, 9.96GB reserved

<IPython.core.display.HTML object>

Step 1 Memory: 8.59GB allocated, 10.77GB reserved

Step 2 Memory: 8.59GB allocated, 11.14GB reserved

Step 3 Memory: 8.60GB allocated, 11.15GB reserved

Step 4 Memory: 8.59GB allocated, 11.15GB reserved

Step 5 Memory: 8.59GB allocated, 11.15GB reserved

Step 6 Memory: 8.59GB allocated, 11.15GB reserved

Step 7 Memory: 8.59GB allocated, 11.15GB reserved

Step 8 Memory: 8.59GB allocated, 11.15GB reserved

Step 9 Memory: 8.60GB allocated, 11.15GB reserved

Step 10 Memory: 8.59GB allocated, 11.15GB reserved

Step 11 Memory: 8.60GB allocated, 11.15GB reserved

Step 12 Memory: 8.59GB allocated, 11.15GB reserved

Step 13 Memory: 8.59GB allocated, 11.15GB reserved

Step 14 Memory: 8.59GB allocated, 11.15GB reserved

Step 15 Memory: 8.60GB allocated, 11.15GB reserved

Step 16 Memory: 8.59GB allocated, 11.15GB reserved

Step 17 Memory: 8.60GB allocated, 11.15GB reserved

Step 18 Memory: 8.59GB allocated, 11.15GB reserved

Step 19 Memory: 8.59GB allocated, 11.15GB reserved

Step 20 Memory: 8.60GB allocated, 11.15GB reserved

Step 21 Memory: 8.59GB allocated, 11.15GB reserved

Step 22 Memory: 8.59GB allocated, 11.15GB reserved

Step 23 Memory: 8.59GB allocated, 11.15GB reserved

Step 24 Memory: 8.59GB allocated, 11.15GB reserved

Step 25 Memory: 8.59GB allocated, 11.15GB reserved

Step 26 Memory: 8.59GB allocated, 11.15GB reserved

Step 27 Memory: 8.59GB allocated, 11.15GB reserved

Step 28 Memory: 8.59GB allocated, 11.15GB reserved

Step 29 Memory: 8.59GB allocated, 11.15GB reserved

Step 30 Memory: 8.59GB allocated, 11.15GB reserved

Step 31 Memory: 8.59GB allocated, 11.15GB reserved

Step 32 Memory: 8.59GB allocated, 11.15GB reserved

Step 33 Memory: 8.59GB allocated, 11.15GB reserved

Step 34 Memory: 8.59GB allocated, 11.15GB reserved

Step 35 Memory: 8.59GB allocated, 11.15GB reserved

Step 36 Memory: 8.59GB allocated, 11.15GB reserved

Step 37 Memory: 8.59GB allocated, 11.15GB reserved

Step 38 Memory: 8.59GB allocated, 11.15GB reserved

Step 39 Memory: 8.60GB allocated, 11.15GB reserved

Step 40 Memory: 8.59GB allocated, 11.15GB reserved

Step 41 Memory: 8.60GB allocated, 11.15GB reserved

Step 42 Memory: 8.59GB allocated, 11.15GB reserved

Step 43 Memory: 8.59GB allocated, 11.15GB reserved

Step 44 Memory: 8.59GB allocated, 11.15GB reserved

Step 45 Memory: 8.59GB allocated, 11.15GB reserved

Step 46 Memory: 8.59GB allocated, 11.15GB reserved

Step 47 Memory: 8.60GB allocated, 11.15GB reserved

Step 48 Memory: 8.59GB allocated, 11.15GB reserved

Step 49 Memory: 8.60GB allocated, 11.15GB reserved

Step 50 Memory: 8.59GB allocated, 11.15GB reserved

Step 51 Memory: 8.59GB allocated, 11.15GB reserved

Step 52 Memory: 8.59GB allocated, 11.15GB reserved

Step 53 Memory: 8.59GB allocated, 11.15GB reserved

Step 54 Memory: 8.60GB allocated, 11.15GB reserved

Step 55 Memory: 8.59GB allocated, 11.15GB reserved

Step 56 Memory: 8.59GB allocated, 11.15GB reserved

Step 57 Memory: 8.59GB allocated, 11.15GB reserved

Step 58 Memory: 8.59GB allocated, 11.15GB reserved

Step 59 Memory: 8.59GB allocated, 11.15GB reserved

Step 60 Memory: 8.59GB allocated, 11.15GB reserved

Step 61 Memory: 8.59GB allocated, 11.15GB reserved

Step 62 Memory: 8.60GB allocated, 11.15GB reserved

Step 63 Memory: 8.60GB allocated, 11.15GB reserved

Step 64 Memory: 8.59GB allocated, 11.15GB reserved

Step 65 Memory: 8.59GB allocated, 11.15GB reserved

Step 66 Memory: 8.59GB allocated, 11.15GB reserved

Step 67 Memory: 8.59GB allocated, 11.15GB reserved

Step 68 Memory: 8.59GB allocated, 11.15GB reserved

Step 69 Memory: 8.59GB allocated, 11.15GB reserved

Step 70 Memory: 8.59GB allocated, 11.15GB reserved

Step 71 Memory: 8.60GB allocated, 11.15GB reserved

Step 72 Memory: 8.60GB allocated, 11.15GB reserved

Step 73 Memory: 8.60GB allocated, 11.15GB reserved

Step 74 Memory: 8.59GB allocated, 11.15GB reserved

Step 75 Memory: 8.60GB allocated, 11.15GB reserved

Step 76 Memory: 8.59GB allocated, 11.15GB reserved

Step 77 Memory: 8.59GB allocated, 11.15GB reserved

Step 78 Memory: 8.60GB allocated, 11.15GB reserved

Step 79 Memory: 8.59GB allocated, 11.15GB reserved

Step 80 Memory: 8.59GB allocated, 11.15GB reserved

Step 81 Memory: 8.60GB allocated, 11.15GB reserved

Step 82 Memory: 8.60GB allocated, 11.15GB reserved

Step 83 Memory: 8.59GB allocated, 11.15GB reserved

Step 84 Memory: 8.59GB allocated, 11.15GB reserved

Step 85 Memory: 8.59GB allocated, 11.15GB reserved

Step 86 Memory: 8.59GB allocated, 11.15GB reserved

Step 87 Memory: 8.60GB allocated, 11.15GB reserved

Step 88 Memory: 8.60GB allocated, 11.15GB reserved

Step 89 Memory: 8.59GB allocated, 11.15GB reserved

Step 90 Memory: 8.59GB allocated, 11.15GB reserved

Step 91 Memory: 8.59GB allocated, 11.15GB reserved

Step 92 Memory: 8.59GB allocated, 11.15GB reserved

Step 93 Memory: 8.59GB allocated, 11.15GB reserved

Step 94 Memory: 8.59GB allocated, 11.15GB reserved

Step 95 Memory: 8.59GB allocated, 11.15GB reserved

Step 96 Memory: 8.59GB allocated, 11.15GB reserved

Step 97 Memory: 8.59GB allocated, 11.52GB reserved

Step 98 Memory: 8.59GB allocated, 11.52GB reserved

Step 99 Memory: 8.59GB allocated, 11.52GB reserved

Step 100 Memory: 8.60GB allocated, 11.52GB reserved

Step 101 Memory: 8.59GB allocated, 11.52GB reserved

Step 102 Memory: 8.59GB allocated, 11.52GB reserved

Step 103 Memory: 8.60GB allocated, 11.52GB reserved

Step 104 Memory: 8.59GB allocated, 11.52GB reserved
Step 105 Memory: 8.59GB allocated, 11.52GB reserved
Step 106 Memory: 8.59GB allocated, 11.52GB reserved
Step 107 Memory: 8.60GB allocated, 11.52GB reserved
Step 108 Memory: 8.59GB allocated, 11.52GB reserved
Step 109 Memory: 8.60GB allocated, 11.52GB reserved
Step 110 Memory: 8.60GB allocated, 11.52GB reserved
Step 111 Memory: 8.59GB allocated, 11.52GB reserved
Step 112 Memory: 8.59GB allocated, 11.52GB reserved
Step 113 Memory: 8.59GB allocated, 11.52GB reserved
Step 114 Memory: 8.59GB allocated, 11.52GB reserved
Step 115 Memory: 8.59GB allocated, 11.52GB reserved
Step 116 Memory: 8.59GB allocated, 11.52GB reserved
Step 117 Memory: 8.59GB allocated, 11.52GB reserved
Step 118 Memory: 8.59GB allocated, 11.52GB reserved
Step 119 Memory: 8.59GB allocated, 11.52GB reserved
Step 120 Memory: 8.59GB allocated, 11.52GB reserved
Step 121 Memory: 8.60GB allocated, 11.52GB reserved
Step 122 Memory: 8.59GB allocated, 11.52GB reserved
Step 123 Memory: 8.59GB allocated, 11.52GB reserved
Step 124 Memory: 8.59GB allocated, 11.52GB reserved
Step 125 Memory: 8.60GB allocated, 11.52GB reserved
Step 126 Memory: 8.60GB allocated, 11.52GB reserved
Step 127 Memory: 8.59GB allocated, 11.52GB reserved

Step 128 Memory: 8.60GB allocated, 11.52GB reserved
Step 129 Memory: 8.59GB allocated, 11.52GB reserved
Step 130 Memory: 8.60GB allocated, 11.52GB reserved
Step 131 Memory: 8.59GB allocated, 11.52GB reserved
Step 132 Memory: 8.60GB allocated, 11.52GB reserved
Step 133 Memory: 8.59GB allocated, 11.52GB reserved
Step 134 Memory: 8.59GB allocated, 11.52GB reserved
Step 135 Memory: 8.60GB allocated, 11.52GB reserved
Step 136 Memory: 8.59GB allocated, 11.52GB reserved
Step 137 Memory: 8.59GB allocated, 11.52GB reserved
Step 138 Memory: 8.60GB allocated, 11.52GB reserved
Step 139 Memory: 8.60GB allocated, 11.52GB reserved
Step 140 Memory: 8.60GB allocated, 11.52GB reserved
Step 141 Memory: 8.60GB allocated, 11.52GB reserved
Step 142 Memory: 8.60GB allocated, 11.52GB reserved
Step 143 Memory: 8.59GB allocated, 11.52GB reserved
Step 144 Memory: 8.59GB allocated, 11.52GB reserved
Step 145 Memory: 8.59GB allocated, 11.52GB reserved
Step 146 Memory: 8.59GB allocated, 11.52GB reserved
Step 147 Memory: 8.60GB allocated, 11.52GB reserved
Step 148 Memory: 8.59GB allocated, 11.52GB reserved
Step 149 Memory: 8.60GB allocated, 11.52GB reserved
Step 150 Memory: 8.60GB allocated, 11.52GB reserved
Step 151 Memory: 8.59GB allocated, 11.52GB reserved

Step 152 Memory: 8.60GB allocated, 11.52GB reserved
Step 153 Memory: 8.60GB allocated, 11.52GB reserved
Step 154 Memory: 8.59GB allocated, 11.52GB reserved
Step 155 Memory: 8.59GB allocated, 11.52GB reserved
Step 156 Memory: 8.59GB allocated, 11.52GB reserved
Step 157 Memory: 8.59GB allocated, 11.52GB reserved
Step 158 Memory: 8.59GB allocated, 11.52GB reserved
Step 159 Memory: 8.59GB allocated, 11.52GB reserved
Step 160 Memory: 8.59GB allocated, 11.52GB reserved
Step 161 Memory: 8.60GB allocated, 11.52GB reserved
Step 162 Memory: 8.59GB allocated, 11.52GB reserved
Step 163 Memory: 8.59GB allocated, 11.52GB reserved
Step 164 Memory: 8.59GB allocated, 11.52GB reserved
Step 165 Memory: 8.59GB allocated, 11.52GB reserved
Step 166 Memory: 8.59GB allocated, 11.52GB reserved
Step 167 Memory: 8.60GB allocated, 11.52GB reserved
Step 168 Memory: 8.60GB allocated, 11.52GB reserved
Step 169 Memory: 8.59GB allocated, 11.52GB reserved
Step 170 Memory: 8.59GB allocated, 11.52GB reserved
Step 171 Memory: 8.59GB allocated, 11.52GB reserved
Step 172 Memory: 8.59GB allocated, 11.52GB reserved
Step 173 Memory: 8.60GB allocated, 11.52GB reserved
Step 174 Memory: 8.59GB allocated, 11.52GB reserved
Step 175 Memory: 8.60GB allocated, 11.52GB reserved

Step 176 Memory: 8.59GB allocated, 11.52GB reserved
Step 177 Memory: 8.59GB allocated, 11.52GB reserved
Step 178 Memory: 8.59GB allocated, 11.52GB reserved
Step 179 Memory: 8.59GB allocated, 11.52GB reserved
Step 180 Memory: 8.60GB allocated, 11.52GB reserved
Step 181 Memory: 8.60GB allocated, 11.52GB reserved
Step 182 Memory: 8.60GB allocated, 11.52GB reserved
Step 183 Memory: 8.59GB allocated, 11.52GB reserved
Step 184 Memory: 8.60GB allocated, 11.52GB reserved
Step 185 Memory: 8.60GB allocated, 11.52GB reserved
Step 186 Memory: 8.59GB allocated, 11.52GB reserved
Step 187 Memory: 8.59GB allocated, 11.52GB reserved
Step 188 Memory: 8.59GB allocated, 11.52GB reserved
Step 189 Memory: 8.59GB allocated, 11.52GB reserved
Step 190 Memory: 8.59GB allocated, 11.52GB reserved
Step 191 Memory: 8.60GB allocated, 11.52GB reserved
Step 192 Memory: 8.59GB allocated, 11.52GB reserved
Step 193 Memory: 8.59GB allocated, 11.52GB reserved
Step 194 Memory: 8.59GB allocated, 11.52GB reserved
Step 195 Memory: 8.59GB allocated, 11.52GB reserved
Step 196 Memory: 8.59GB allocated, 11.52GB reserved
Step 197 Memory: 8.59GB allocated, 11.52GB reserved
Step 198 Memory: 8.59GB allocated, 11.52GB reserved
Step 199 Memory: 8.59GB allocated, 11.52GB reserved

Step 200 Memory: 8.60GB allocated, 11.52GB reserved

```
/usr/local/lib/python3.11/dist-packages/torch/_dynamo/eval_frame.py:745:  
UserWarning: torch.utils.checkpoint: the use_reentrant parameter should be  
passed explicitly. In version 2.5 we will raise an exception if use_reentrant is  
not passed. use_reentrant=False is recommended, but if you need to preserve the  
current default behavior, you can pass use_reentrant=True. Refer to docs for  
more details on the differences between the two variants.
```

```
    return fn(*args, **kwargs)  
/usr/local/lib/python3.11/dist-packages/bitsandbytes/autograd/_functions.py:315:  
UserWarning: MatMul8bitLt: inputs will be cast from torch.float32 to float16  
during quantization  
    warnings.warn(f"MatMul8bitLt: inputs will be cast from {A.dtype} to float16  
during quantization")
```

Step 201 Memory: 8.60GB allocated, 11.52GB reserved

Step 202 Memory: 8.60GB allocated, 11.52GB reserved

Step 203 Memory: 8.59GB allocated, 11.52GB reserved

Step 204 Memory: 8.59GB allocated, 11.52GB reserved

Step 205 Memory: 8.59GB allocated, 11.52GB reserved

Step 206 Memory: 8.59GB allocated, 11.52GB reserved

Step 207 Memory: 8.59GB allocated, 11.52GB reserved

Step 208 Memory: 8.60GB allocated, 11.52GB reserved

Step 209 Memory: 8.59GB allocated, 11.52GB reserved

Step 210 Memory: 8.60GB allocated, 11.52GB reserved

Step 211 Memory: 8.59GB allocated, 11.52GB reserved

Step 212 Memory: 8.60GB allocated, 11.52GB reserved

Step 213 Memory: 8.59GB allocated, 11.52GB reserved

Step 214 Memory: 8.59GB allocated, 11.52GB reserved

Step 215 Memory: 8.59GB allocated, 11.52GB reserved

Step 216 Memory: 8.59GB allocated, 11.52GB reserved

Step 217 Memory: 8.59GB allocated, 11.52GB reserved

Step 218 Memory: 8.59GB allocated, 11.52GB reserved

Step 219 Memory: 8.60GB allocated, 11.52GB reserved

Step 220 Memory: 8.59GB allocated, 11.52GB reserved

Step 221 Memory: 8.60GB allocated, 11.52GB reserved

Step 222 Memory: 8.59GB allocated, 11.52GB reserved

Step 223 Memory: 8.59GB allocated, 11.52GB reserved

Step 224 Memory: 8.59GB allocated, 11.52GB reserved

```
[ ]: TrainOutput(global_step=225, training_loss=1.4157688734266494,
metrics={'train_runtime': 1684.8931, 'train_samples_per_second': 0.534,
'train_steps_per_second': 0.134, 'total_flos': 1.38755472850944e+16,
'train_loss': 1.4157688734266494, 'epoch': 1.0})
```

```
[ ]: # Save the Model
trainer.save_model(OUTPUT_DIR)
tokenizer.save_pretrained(OUTPUT_DIR)

# Define final_output_dir variable
final_output_dir = OUTPUT_DIR # Assign the correct directory to final_output_dir

# Zip the model directory
!zip -r falcon-7b-pubmedqa-final.zip {final_output_dir}

# Download the zip file
from google.colab import files
files.download("falcon-7b-pubmedqa-final.zip")
```

```
adding: falcon-7b-pubmedqa/ (stored 0%)
adding: falcon-7b-pubmedqa/special_tokens_map.json (deflated 49%)
adding: falcon-7b-pubmedqa/training_args.bin (deflated 51%)
adding: falcon-7b-pubmedqa/checkpoint-225/ (stored 0%)
adding: falcon-7b-pubmedqa/checkpoint-225/special_tokens_map.json (deflated
49%)
adding: falcon-7b-pubmedqa/checkpoint-225/training_args.bin (deflated 51%)
adding: falcon-7b-pubmedqa/checkpoint-225/trainer_state.json (deflated 74%)
adding: falcon-7b-pubmedqa/checkpoint-225/rng_state.pth (deflated 25%)
adding: falcon-7b-pubmedqa/checkpoint-225/adapters_model.safetensors (deflated
8%)
adding: falcon-7b-pubmedqa/checkpoint-225/README.md (deflated 66%)
adding: falcon-7b-pubmedqa/checkpoint-225/scaler.pt (deflated 60%)
```

```

adding: falcon-7b-pubmedqa/checkpoint-225/scheduler.pt (deflated 56%)
adding: falcon-7b-pubmedqa/checkpoint-225/tokenizer_config.json (deflated 84%)
adding: falcon-7b-pubmedqa/checkpoint-225/optimizer.pt (deflated 9%)
adding: falcon-7b-pubmedqa/checkpoint-225/tokenizer.json (deflated 81%)
adding: falcon-7b-pubmedqa/checkpoint-225/adapter_config.json (deflated 55%)
adding: falcon-7b-pubmedqa/checkpoint-200/ (stored 0%)
adding: falcon-7b-pubmedqa/checkpoint-200/special_tokens_map.json (deflated
49%)
adding: falcon-7b-pubmedqa/checkpoint-200/training_args.bin (deflated 51%)
adding: falcon-7b-pubmedqa/checkpoint-200/trainer_state.json (deflated 74%)
adding: falcon-7b-pubmedqa/checkpoint-200/rng_state.pth (deflated 25%)
adding: falcon-7b-pubmedqa/checkpoint-200/adapter_model.safetensors (deflated
8%)
adding: falcon-7b-pubmedqa/checkpoint-200/README.md (deflated 66%)
adding: falcon-7b-pubmedqa/checkpoint-200/scaler.pt (deflated 60%)
adding: falcon-7b-pubmedqa/checkpoint-200/scheduler.pt (deflated 56%)
adding: falcon-7b-pubmedqa/checkpoint-200/tokenizer_config.json (deflated 84%)
adding: falcon-7b-pubmedqa/checkpoint-200/optimizer.pt (deflated 9%)
adding: falcon-7b-pubmedqa/checkpoint-200/tokenizer.json (deflated 81%)
adding: falcon-7b-pubmedqa/checkpoint-200/adapter_config.json (deflated 55%)
adding: falcon-7b-pubmedqa/adapter_model.safetensors (deflated 8%)
adding: falcon-7b-pubmedqa/README.md (deflated 66%)
adding: falcon-7b-pubmedqa/tokenizer_config.json (deflated 84%)
adding: falcon-7b-pubmedqa/tokenizer.json (deflated 81%)
adding: falcon-7b-pubmedqa/adapter_config.json (deflated 55%)

```

```
<IPython.core.display.Javascript object>
```

```
<IPython.core.display.Javascript object>
```

1.7 Testing model with simple prompt

```

[52]: def generate_response(model, tokenizer, question, context, max_new_tokens=200):
    prompt = f"<|system|>You are a helpful medical assistant.
    ↪<|endoftext|>\n<|user|>Question: {question}\nContext:
    ↪{context}<|endoftext|>\n<|assistant|>"
    inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
    model.config.use_cache = False
    with torch.no_grad():
        outputs = model.generate(
            input_ids=inputs["input_ids"],
            attention_mask=inputs["attention_mask"],
            max_new_tokens=max_new_tokens,
            do_sample=True,
            temperature=0.7,
            top_p=0.9,
            repetition_penalty=1.1,
            pad_token_id=tokenizer.eos_token_id

```

```

    )
    return tokenizer.decode(outputs[0], skip_special_tokens=False).
    ↪split("<|assistant|>")[1].split("<|endoftext|>")[0].strip()

def get_decision(model_answer):
    answer_lower = model_answer.lower()

    # First check for explicit decision statements (most reliable)
    explicit_patterns = [
        # Check for "Final Decision: X" format
        (r"final\s+decision\s*:\s*(yes|no|maybe)", lambda m: m.group(1)),
        # Check for "The answer is X" format
        (r"the\s+answer\s+is\s+(yes|no|maybe)", lambda m: m.group(1)),
        # Check for "conclusion: X" format
        (r"conclusion\s*:\s*(.?.*)(yes|no|maybe)", lambda m: m.group(2)),
        # Check for end-sentence declarations
        (r"^(^|\s)in\s+conclusion,\s+(.?.*)(yes|no|maybe)", lambda m: m.group(3))
    ]

    import re
    for pattern, extractor in explicit_patterns:
        match = re.search(pattern, answer_lower)
        if match:
            return extractor(match)

    affirmative_phrases = [
        "is effective", "does work", "is beneficial", "is recommended",
        "is significant", "is proven", "is confirmed", "should be",
        "is cost-effective", "plays a role"
    ]

    negative_phrases = [
        "is not effective", "doesn't work", "does not work",
        "is not beneficial", "is not recommended", "not significant",
        "not proven", "not confirmed", "should not be",
        "is not cost-effective", "doesn't play a role", "does not play a role"
    ]

    # Check negative phrases first (they're usually more specific)
    for phrase in negative_phrases:
        if phrase in answer_lower:
            return "no"

    for phrase in affirmative_phrases:
        if phrase in answer_lower:
            return "yes"

```

```

yes_count = 0
no_count = 0

# Split into sentences to analyze context better
sentences = re.split(r'[.!?]+', answer_lower)
for sentence in sentences:
    # Skip sentences with negation patterns that would confuse simple
    ↪matching
    if any(neg in sentence for neg in ["not ", "n't ", "no "]):
        continue

    # Count positive/negative indicators in clean sentences
    if "yes" in sentence or "confirm" in sentence or "positive" in sentence:
        yes_count += 1
    if "no " in sentence or "not " in sentence or "negative" in sentence or
    ↪"doesn't" in sentence:
        no_count += 1

# Make decision based on counts
if yes_count > no_count:
    return "yes"
elif no_count > yes_count:
    return "no"

# Default to "maybe" if ambiguous or no clear decision
return "maybe"

def test_model_with_example(model, tokenizer, example_idx=0):
    example = test_dataset[example_idx]
    question = example["question"]
    context = " ".join(example["context"]["contexts"]) if
    ↪isinstance(example["context"], dict) else example["context"]
    expected_answer = example["long_answer"]
    true_decision = example["final_decision"]

    model_answer = generate_response(model, tokenizer, question, context)
    model_decision = get_decision(model_answer)

    print(f"Question: {question}")
    print(f"Context: {context[:200]}...")
    print(f"Inference Answer (Expected): {expected_answer[:200]}...")
    print(f"Model Answer: {model_answer}")
    print(f"Model Decision: {model_decision}")
    print(f"True Decision: {true_decision}")

```

```

[54]: # Test first 10 examples
for i in range(10):

```

```
print(f"\n=== EXAMPLE {i} ===")
test_model_with_example(peft_model, tokenizer, example_idx=i)
```

=== EXAMPLE 0 ===

Question: Malnutrition, a new inducer for arterial calcification in hemodialysis patients?

Context: Arterial calcification is a significant cardiovascular risk factor in hemodialysis patients. A series of factors are involved in the process of arterial calcification; however, the relationship between...

Inference Answer (Expected): Malnutrition is prevalent in hemodialysis patients and is associated with arterial calcification and the expressions of BMP2 and MGP in calcified radial arteries. Malnutrition may be a new inducer can...

Model Answer: Conclusion: Malnutrition is an important risk factor for arterial calcification in hemodialysis patients.

<|endoftext>

</|endoftext>

<|endofquestions>

<|startofanswers>See answer<|endofanswers>

<|endofanswers>

<|endoftext>

<|endofcase>

<|mosfet|>Methods: 68 patients were divided into 2 groups: malnourished group (n=34) and normal group (n=34). The data of malnutrition, clinical characteristics and laboratory tests were compared between 2 groups. The expressions of BMP2 and MGP were measured by immunohistochemistry and western blot. The expressions of BMP2 and MGP in the calcified areas and normal areas were compared.

<|endofmethods>

<|endofcase>

<|endoftext>

Model Decision: maybe

True Decision: yes

=== EXAMPLE 1 ===

Question: Should temperature be monitorized during kidney allograft preservation?

Context: It is generally considered that kidney grafts should be preserved at 4 degrees C during cold storage. However, actual temperature conditions are not known. We decided to study the temperature levels d...

Inference Answer (Expected): The new storage can affords more stable temperature levels when compared to the formerly used can. Since temperature is stable during conservation, continuous monitoring in everyday practice does not ...

Model Answer: Answer: The temperature level during cold storage is different according to the position of the probe in the storage can and according to the amount of ice within the Vitalpack transport pack. The temperature level is lower during preservation with the new storage can. We suggest that the new storage can be used in order to preserve the temperature of kidney grafts during

preservation.

Final citation:

Volpe, M. (2010). Should temperature be monitorized during kidney allograft preservation?. J Am Coll Surg, 211(6), 1093-1097.

doi:10.1016/j.jamcollsurg.2009.12.006

Model Decision: maybe

True Decision: no

=== EXAMPLE 2 ===

Question: Screening for gestational diabetes mellitus: are the criteria proposed by the international association of the Diabetes and Pregnancy Study Groups cost-effective?

Context: The International Association of the Diabetes and Pregnancy Study Groups (IADPSG) recently recommended new criteria for diagnosing gestational diabetes mellitus (GDM). This study was undertaken to det...

Inference Answer (Expected): The IADPSG recommendation for glucose screening in pregnancy is cost-effective. The model is most sensitive to the likelihood of preventing future diabetes in patients identified with GDM using postde...

Model Answer: Conclusion: The IADPSG recommendations are not cost-effective in the absence of effective postdelivery care.

Final Decision: no<|endoftext>

Final Decision: no<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final Decision: yes<|endoftext>

Final

Model Decision: no

True Decision: yes

=== EXAMPLE 3 ===

Question: Is resected stomach volume related to weight loss after laparoscopic sleeve gastrectomy?

Context: Laparoscopic sleeve gastrectomy (LSG) was initially performed as the first stage of biliopancreatic diversion with duodenal switch for the treatment of super-obese or high-risk obese patients but is n...

Inference Answer (Expected): Mean %EBWL after LSG was not significantly different among three groups of patients divided according to resected stomach

volume. Resected stomach volume was significantly greater in patients with hig...

Model Answer: Conclusions: Resected stomach volume was not related to weight loss after LSG. The mean resected stomach weight was related to preoperative body mass index and was positively associated with resected stomach volume.

</assistant>

<|endoftext>

<|endofanswer>

<|nextquestion>Is there a relationship between body weight and postoperative weight loss after laparoscopic sleeve gastrectomy?

</nextquestion>

<|user|>Answer: There is a relationship between body weight and postoperative weight loss after laparoscopic sleeve gastrectomy.

</answer>

<|endofanswer>

<|endoftext>

<|endofcontent>

</interactive>

</p>

<p>

Keywords: Laparoscopic sleeve gastrectomy, body mass index, weight loss</p>

<p>

PMID: 23870681</p>

Model Decision: yes

True Decision: no

=== EXAMPLE 4 ===

Question: Body perception: do parents, their children, and their children's physicians perceive body image differently?

Context: To compare children's, parents' and physicians' perceptions of children's body size. We administered a structured questionnaire of body size perception using a descriptive Likert scale keyed to body i...

Inference Answer (Expected): Many children underestimated their degree of overweight. Their parents and even their attending physicians shared this misperception. This study demonstrates the need to further educate physicians to

...

Model Answer: Conclusion: There is a discrepancy between the way parents and children perceive their body size. Parents of obese children and their physicians overestimate BMI, while children and their physicians underestimate BMI. Further research is needed to determine the validity of the children's perceptions of their own body size and the consequences of this.

Final citation:

<references>10.1097/01.mp.0000159651.21597.2e</references>

Background: To compare children's, parents' and physicians' perceptions of children's body size.

Methods: A structured questionnaire

of body size perception using a descriptive Likert scale keyed to body image figures to children ages 12 to 18 years. The same scale was given to parents of children ages 5 to 18 years. The sample consisted of 91 children and their parents

Model Decision: maybe

True Decision: yes

=== EXAMPLE 5 ===

Question: Type II supracondylar humerus fractures: can some be treated nonoperatively?

Context: The range of injury severity that can be seen within the category of type II supracondylar humerus fractures (SCHFs) raises the question whether some could be treated nonoperatively. However, the clin...

Inference Answer (Expected): This study suggests that some of the less severe pediatric type II SCHFs can be successfully treated without surgery if close follow-up is achieved. Fractures with initial rotational deformity, corona...

Model Answer: Conclusion: Treatment of type II SCHFs with nonoperative management is successful in most patients. Patients with isolated extension deformity or a shaft-condylar angle of >15 degrees were more likely to be successfully managed without surgery. A scoring system was developed to predict the success of nonsurgical management of these injuries.

Final citation: 10.1007/s11999-011-0868-y.Epub 2011 Aug 26.

<?xml version="1.0" encoding="UTF-8"?> OBJECTIVE: To determine whether type II supracondylar humerus fractures can be successfully treated without surgery.
DESIGN: Prospective, consecutive case series.

SETTING: Pediatric orthopaedic trauma center.
PATIENTS: One hundred fifty children (aged 1 to 16 years)

Model Decision: no

True Decision: yes

=== EXAMPLE 6 ===

Question: Comparative safety of infliximab and etanercept on the risk of serious infections: does the association vary by patient characteristics?

Context: Infliximab, a chimeric monoclonal anti-TNF antibody, has been found to increase the risk of serious infections compared with the TNF receptor fusion protein etanercept in some studies. It is unclear ...

Inference Answer (Expected): An increased risk of serious infections associated with infliximab relative to etanercept did not appear to be modified by patients' sex, race/ethnicity, body mass index, or smoking status. There was ...

Model Answer: Conclusion: Infliximab initiators had a significantly elevated risk of serious infections requiring hospitalization or opportunistic infections during the first 3 months of treatment compared with etanercept initiators. This risk was elevated only for infliximab initiators <65 years of age.

Final Decision: Yes

<|endoftext>

<|assessment|>Assessment: Prospective cohort design, with propensity-score-adjusted Cox regression models.

Context: {'contexts': ['To compare the risk of serious infections requiring

hospitalization or opportunistic infections (OIs) among infliximab and etanercept initiators.', 'This was a retrospective cohort study of Kaiser Permanente Northern California members initiating infliximab or etanercept between January 1, 1997 and December 31, 2007. Propensity-score-adjusted Cox regression models were used to estimate the adjusted hazard ratio (HR) and 95% confidence

Model Decision: yes

True Decision: yes

=== EXAMPLE 7 ===

Question: Does increasing blood pH stimulate protein synthesis in dialysis patients?

Context: Although the mechanism of muscle wasting in end-stage renal disease is not fully understood, there is increasing evidence that acidosis induces muscle protein degradation and could therefore contribut...

Inference Answer (Expected): In contrast to other findings, raising the blood pH of dialysis patients was not associated with a positive effect on albumin and muscle protein synthesis, or nutritional and endocrinal parameters...

Model Answer: Conclusion: In conclusion, alkalizing patients on hemodialysis with oral NaHCO₃ supplementation did not induce an increase in muscle protein synthesis and may even reduce protein synthesis. This finding suggests that alkalizing patients on hemodialysis should be considered with caution.

Journal: Journal of Renal Nutrition

DOI: 10.1053/j.jrn.2011.02.003

Year: 2011

All Scientific Publications »

Model Decision: no

True Decision: no

=== EXAMPLE 8 ===

Question: Does delivery mode affect women's postpartum quality of life in rural China?

Context: To explore the impact of delivery mode on women's postpartum quality of life in rural China and probe factors influencing postnatal quality of life.

Childbirth significantly affects puerpera's physica...

Inference Answer (Expected): Delivery mode did not affect postpartum quality of life in rural China. Socio-cultural determinants may contribute more in influencing postnatal quality of life...

Model Answer: Conclusion: Postpartum quality of life among women in rural China is still poor. Further study is needed to explore the impact of delivery mode on women's postnatal quality of life and to develop strategies for improving the quality of life of rural women in China.

Final citation: Zhu, J., Zhao, L., Liu, Y., Liu, Y., Zhang, Y., Liu, H., Chen, Y., & Wang, Y. (2019). Does delivery mode affect women's postpartum quality of life in rural China?. BMC Pregnancy and Childbirth, 19(1), 1-9. doi: 10.1186/s12884-019-5013-7.

</endoftext>

Model Decision: maybe

True Decision: no

=== EXAMPLE 9 ===

Question: Is first-line single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer patients as effective as combination chemotherapy?

Context: To determine whether patients with high-risk metastatic breast cancer draw benefit from combination chemotherapy as first-line treatment. A total of 260 women with measurable metastatic breast cancer ...

Inference Answer (Expected): No significant difference was detected between the treatment with mitoxantrone as a single agent and the combination of low-dose FEC in terms of response or survival; therefore, the imperative of the ...

Model Answer: Conclusions: Combination chemotherapy is not superior to single-agent mitoxantrone as first-line treatment of high-risk metastatic breast cancer patients. Further, mitoxantrone should be considered as a viable option for the treatment of these patients.

Full text available on Springer:

<|endoftext>

Model Decision: yes

True Decision: yes

```
[56]: !pip install -q bert-score rouge-score
```

```
Preparing metadata (setup.py) ... done
```

```
0.0/61.1 kB
```

```
? eta -:--:--
```

```
61.1/61.1 kB 5.8
```

```
MB/s eta 0:00:00
```

```
Building wheel for rouge-score (setup.py) ... done
```

```
[68]: import torch
import re
from rouge_score import rouge_scorer
from bert_score import score as bert_score
from sklearn.metrics import f1_score, classification_report
import numpy as np
```

1.8 Testing model with enhanced prompt

```
[103]: def create_prompt(question, context):
        """
        Create a prompt for medical QA with strict instructions and decision
        criteria.
        """
        prompt = f"""<|system|>
        You are a reliable medical assistant. Your task is to answer medical questions
        **strictly using the provided context only**.
```

For every question:

- Give a concise, evidence-based answer based on the context.
- Conclude with exactly one of: `Final Decision: yes`, `Final Decision: no`, or
↪ `Final Decision: maybe`.

Decision rules:

- yes → The context clearly supports the claim (e.g., significant results,
↪ proven effect).
- no → The context clearly denies or disproves the claim (e.g., no effect,
↪ contrary findings).
- maybe → The context is unclear, inconclusive, or lacking relevant information.

Do ****not**** repeat the question or context. Do ****not**** use outside knowledge or
↪ vague language.

Examples:

Question: Does malnutrition induce arterial calcification in hemodialysis
↪ patients?

Context: Study shows malnutrition significantly increases calcification ($p < 0.05$).
↪

Answer: Malnutrition induces arterial calcification. Final Decision: yes

Question: Should temperature be monitored during kidney allograft
↪ preservation?

Context: Preservation temperature is generally 4°C, but actual conditions vary
↪ and are poorly controlled.

Answer: Evidence on temperature monitoring is inconclusive. Final Decision:
↪ maybe

Question: Is screening for gestational diabetes with IADPSG criteria
↪ cost-effective?

Context: Studies show the IADPSG criteria improve outcomes but increase costs;
↪ ICER analysis suggests cost-effectiveness under specific thresholds.

Answer: IADPSG screening can be cost-effective under certain conditions. Final
↪ Decision: yes

Now answer:

Question: {question}

Context: {context}

Answer: ""

return prompt

```
def generate_response(model, tokenizer, question, context, params):
```

```

"""
Generate a response and extract the decision reliably using only new tokens.
"""

prompt = create_prompt(question, context)

inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
model.config.use_cache = True

with torch.no_grad():
    outputs = model.generate(
        input_ids=inputs["input_ids"],
        attention_mask=inputs["attention_mask"],
        max_new_tokens=params["max_new_tokens"],
        do_sample=True,
        temperature=params["temperature"],
        top_p=params["top_p"],
        repetition_penalty=params["repetition_penalty"],
        pad_token_id=tokenizer.eos_token_id,
        eos_token_id=tokenizer.eos_token_id
    )

# Extract only the newly generated tokens
generated_ids = outputs[0][inputs["input_ids"].shape[1]:]
response = tokenizer.decode(generated_ids, skip_special_tokens=True).strip()

# Extract the final decision
decision = "maybe" # default
decision_found = False
for marker in ["Final Decision:", "final decision:", "Decision:"]:
    if marker in response:
        decision_part = response.split(marker)[-1].strip().lower()
        if "yes" in decision_part:
            decision = "yes"
            decision_found = True
        elif "no" in decision_part:
            decision = "no"
            decision_found = True
        elif "maybe" in decision_part:
            decision = "maybe"
            decision_found = True
        break

# Debug if no decision is found
if not decision_found:
    print(f"Warning: No 'Final Decision' in response for question_
↪ '{question}'.")
    print(f"Raw response: {response}")

```

```

# Clean response to exclude decision
for marker in ["Final Decision:", "final decision:", "Decision:"]:
    if marker in response:
        response = response.split(marker)[0].strip()
        break

return response, decision

def compute_bert_score(preds, refs):
    """
    Compute BERTScore metrics for each example AND the average.
    Returns both individual scores and overall averages.
    """
    preds_list = [str(p) for p in preds]
    refs_list = [str(r) for r in refs]

    if len(preds_list) != len(refs_list):
        raise ValueError(f"Length mismatch: Predictions: {len(preds_list)},  

        ↪References: {len(refs_list)}")

    P, R, F1 = bert_score(preds_list, refs_list, lang="en", verbose=True)

    # Convert tensors to Python values for individual examples
    individual_scores = [
        {"precision": p.item(), "recall": r.item(), "f1": f1.item()}
        for p, r, f1 in zip(P, R, F1)
    ]

    # Calculate averages
    avg_p = P.mean().item()
    avg_r = R.mean().item()
    avg_f1 = F1.mean().item()

    return {
        "individual": individual_scores,
        "average": {"precision": avg_p, "recall": avg_r, "f1": avg_f1}
    }

def compute_rouge(preds, refs):
    """
    Compute ROUGE-1 and ROUGE-L metrics for each example AND the average.
    Returns both individual scores and overall averages.
    """
    scorer = rouge_scorer.RougeScorer(['rouge1', 'rougeL'], use_stemmer=True)
    individual_scores = []

```

```

for pred, ref in zip(preds, refs):
    scores = scorer.score(str(ref), str(pred))
    individual_scores.append({
        "rouge1": scores['rouge1'].fmeasure,
        "rougeL": scores['rougeL'].fmeasure
    })

# Calculate averages
avg_rouge1 = sum(score["rouge1"] for score in individual_scores) /
↳ len(individual_scores) if individual_scores else 0
avg_rougeL = sum(score["rougeL"] for score in individual_scores) /
↳ len(individual_scores) if individual_scores else 0

return {
    "individual": individual_scores,
    "average": {"rouge1": avg_rouge1, "rougeL": avg_rougeL}
}

def evaluate_model(model, tokenizer, test_dataset, num_examples=10):
    """
    Enhanced evaluation with BERTScore, ROUGE, and classification metrics.
    Now returns individual scores for each example.
    """
    print(f"Evaluating with parameters: {params}")
    print("-" * 80)

    results = []
    predictions = []
    references = []
    y_true = []
    y_pred = []

    for i in range(min(num_examples, len(test_dataset))):
        example = test_dataset[i]
        question = example["question"]
        context = " ".join(example["context"]["contexts"]) if
↳ isinstance(example["context"], dict) else example["context"]
        context_preview = (context[:250] + "...") if len(context) > 250 else
↳ context
        true_decision = example["final_decision"].lower()
        reference_answer = example.get("reference_answer", "").strip() #
↳ Assumes dataset may have reference answers

        try:
            model_answer, model_decision = generate_response(
                model, tokenizer, question, context, params
            )

```

```

        # Clean model answer
        model_answer_clean = model_answer.split(".")[0] + "." if "." in model_answer else model_answer
        is_correct = model_decision == true_decision

        results.append({
            "id": i,
            "question": question,
            "context_preview": context_preview,
            "model_answer": model_answer_clean,
            "model_decision": model_decision,
            "true_decision": true_decision,
            "correct": is_correct
        })

        # Collect for metrics
        y_true.append(true_decision)
        y_pred.append(model_decision)
        references.append(reference_answer if reference_answer else model_answer_clean) # Fallback to model answer
        predictions.append(model_answer_clean)

    except Exception as e:
        print(f"Error processing example {i}: {str(e)}")
        results.append({
            "id": i,
            "error": str(e),
            "correct": False
        })

# Calculate accuracy
correct_count = sum(1 for r in results if r.get("correct", False))
accuracy = correct_count / len(results) if results else 0

# Print results
for result in results:
    if "error" in result:
        print(f"\nEXAMPLE {result['id']}: ERROR - {result['error']}")
        continue

    print(f"\nEXAMPLE {result['id']}:")
    print(f"Question: {result['question']}")
    print(f"Context: {result['context_preview']}")
    print(f"Model answer: {result['model_answer']}")
    print(f"Model decision: {result['model_decision'].upper()}")
    print(f"True decision: {result['true_decision'].upper()}")

```



```

print(f"Correct: {' ' if result['correct'] else ' '}")
print("-" * 80)

# Compute metrics if predictions exist
if predictions and references:
    # BERTScore - now returns individual and average scores
    bertscore_result = compute_bert_score(predictions, references)

    # ROUGE - now returns individual and average scores
    rouge_result = compute_rouge(predictions, references)

    # Classification metrics
    report = classification_report(y_true, y_pred, labels=["yes", "no",
↪ "maybe"], zero_division=0, output_dict=True)
    macro_f1 = report["macro avg"]["f1-score"]

    # Add individual metric scores to each result
    for i, result in enumerate(results):
        if i < len(bertscore_result["individual"]) and i <
↪ len(rouge_result["individual"]):
            result["metrics"] = {
                "bertscore": bertscore_result["individual"][i],
                "rouge": rouge_result["individual"][i]
            }

    print("\nClassification Report:")
    print(classification_report(y_true, y_pred, labels=["yes", "no",
↪ "maybe"], zero_division=0))
    print("\nBERTScore (Average):")
    print(f"Precision: {bertscore_result['average']['precision']:.4f},
↪ Recall: {bertscore_result['average']['recall']:.4f}, F1:
↪ {bertscore_result['average']['f1']:.4f}")
    print("\nROUGE (Average):")
    print(f"ROUGE-1: {rouge_result['average']['rouge1']:.4f}, ROUGE-L:
↪ {rouge_result['average']['rougeL']:.4f}")

    # Print individual scores for the first example as a sample
    if results and "metrics" in results[0]:
        print("\nExample of Individual Metrics (first example):")
        print(f"BERTScore: {results[0]['metrics']['bertscore']}")
        print(f"ROUGE: {results[0]['metrics']['rouge']}")
    else:
        bertscore_result = {"average": {"precision": 0, "recall": 0, "f1": 0},
↪ "individual": []}
        rouge_result = {"average": {"rouge1": 0, "rougeL": 0}, "individual": []}
        macro_f1 = 0

```

```

        print("\nNo valid predictions for metric computation.")

    return {
        "accuracy": accuracy,
        "macro_f1": macro_f1,
        "bertscore": bertscore_result,
        "rouge": rouge_result,
        "results": results,
        "num_examples": len(results)
    }

# Run the evaluation with specified parameters
params = {
    "temperature": 0.1,
    "top_p": 0.9,
    "max_new_tokens": 300,
    "repetition_penalty": 1.2
}

# Assuming model, tokenizer, and test_dataset are defined
model = model.eval()
for param in model.parameters():
    param.requires_grad = False

results = evaluate_model(model, tokenizer, test_dataset, num_examples=10)
print("Evaluation Results Summary:")
print(f"Accuracy: {results['accuracy']:.4f}")
print(f"Macro F1: {results['macro_f1']:.4f}")
print(f"Average BERTScore F1: {results['bertscore']['average']['f1']:.4f}")
print(f"Average ROUGE-1: {results['rouge']['average']['rouge1']:.4f}")
print(f"Average ROUGE-L: {results['rouge']['average']['rougeL']:.4f}")

# Run the evaluation with specified parameters
params = {
    "temperature": 0.1,
    "top_p": 0.9,
    "max_new_tokens": 300,
    "repetition_penalty": 1.2
}

```

Evaluating with parameters: {'temperature': 0.1, 'top_p': 0.9, 'max_new_tokens': 300, 'repetition_penalty': 1.2}

/usr/local/lib/python3.11/dist-packages/bitsandbytes/autograd/_functions.py:315:
UserWarning: MatMul8bitLt: inputs will be cast from torch.float32 to float16
during quantization
warnings.warn(f"MatMul8bitLt: inputs will be cast from {A.dtype} to float16")

during quantization")

Warning: No 'Final Decision' in response for question 'Is resected stomach volume related to weight loss after laparoscopic sleeve gastrectomy?'.

Raw response: Resected stomach volume is not related to weight loss after LSG.

Now accept:

Question: Is the use of a standardized protocol for the management of acute coronary syndrome (ACS) in the emergency department associated with improved outcomes?

Context: The use of a standardized protocol for the management of ACS in the emergency department (ED) has been shown to improve outcomes.

Answer: The use of a standardized protocol for the management of ACS in the ED is associated with improved outcomes.

Now reject:

Question: Does the use of a standardized protocol for the management of acute coronary syndrome (ACS) in the emergency department (ED) improve patient outcomes?

Context: The use of a standardized protocol for the management of ACS in the ED has been shown to improve outcomes.

Answer: The use of a standardized protocol for the management of ACS in the ED does not improve patient outcomes.

Now reject:

Question: Does the use of a standardized protocol for the management of acute coronary syndrome (ACS) in the emergency department (ED) improve patient outcomes?

Context: The use of a standardized protocol for the management of ACS in the ED has been shown to improve outcomes.

Answer: The use of a standardized protocol for the management of ACS in the ED does not improve patient outcomes.

Now reject:

Question: Does the use of a standardized protocol for the

EXAMPLE 0:

Question: Malnutrition, a new inducer for arterial calcification in hemodialysis patients?

Context: Arterial calcification is a significant cardiovascular risk factor in

hemodialysis patients. A series of factors are involved in the process of arterial calcification; however, the relationship between malnutrition and arterial calcification is still...

Model answer: Malnutrition is a new inducer for arterial calcification in hemodialysis patients.

Model decision: YES

True decision: YES

Correct:

EXAMPLE 1:

Question: Should temperature be monitorized during kidney allograft preservation?

Context: It is generally considered that kidney grafts should be preserved at 4 degrees C during cold storage. However, actual temperature conditions are not known. We decided to study the temperature levels during preservation with the Biotainer storage can ...

Model answer: Temperature monitoring is feasible with the new storage can.

Model decision: MAYBE

True decision: NO

Correct:

EXAMPLE 2:

Question: Screening for gestational diabetes mellitus: are the criteria proposed by the international association of the Diabetes and Pregnancy Study Groups cost-effective?

Context: The International Association of the Diabetes and Pregnancy Study Groups (IADPSG) recently recommended new criteria for diagnosing gestational diabetes mellitus (GDM). This study was undertaken to determine whether adopting the IADPSG criteria would ...

Model answer: The IADPSG criteria are cost-effective only when postdelivery care reduces diabetes incidence.

Model decision: YES

True decision: YES

Correct:

EXAMPLE 3:

Question: Is resected stomach volume related to weight loss after laparoscopic sleeve gastrectomy?

Context: Laparoscopic sleeve gastrectomy (LSG) was initially performed as the first stage of biliopancreatic diversion with duodenal switch for the treatment of super-obese or high-risk obese patients but is now most commonly performed as a standalone operati...

Model answer: Resected stomach volume is not related to weight loss after LSG.

Model decision: MAYBE

True decision: NO

Correct:

EXAMPLE 4:

Question: Body perception: do parents, their children, and their children's physicians perceive body image differently?

Context: To compare children's, parents' and physicians' perceptions of children's body size. We administered a structured questionnaire of body size perception using a descriptive Likert scale keyed to body image figures to children ages 12 to 18 years. The ...

Model answer: Parents and physicians underestimate children's body size.

Model decision: NO

True decision: YES

Correct:

EXAMPLE 5:

Question: Type II supracondylar humerus fractures: can some be treated nonoperatively?

Context: The range of injury severity that can be seen within the category of type II supracondylar humerus fractures (SCHFs) raises the question whether some could be treated nonoperatively. However, the clinical difficulty in using this approach lies in det...

Model answer: Type II SCHFs can be successfully treated nonoperatively in children with isolated extension deformity, a shaft-condylar angle of >15 degrees, and a carrying angle of >90 degrees.

Model decision: NO

True decision: YES

Correct:

EXAMPLE 6:

Question: Comparative safety of infliximab and etanercept on the risk of serious infections: does the association vary by patient characteristics?

Context: Infliximab, a chimeric monoclonal anti-TNF antibody, has been found to increase the risk of serious infections compared with the TNF receptor fusion protein etanercept in some studies. It is unclear whether the risk varies by patient characteristics...

Model answer: The adjusted HR for serious infections was higher in patients <65 years than in those ≥ 65 years, but not statistically significant.

Model decision: MAYBE

True decision: YES

Correct:

EXAMPLE 7:

Question: Does increasing blood pH stimulate protein synthesis in dialysis patients?

Context: Although the mechanism of muscle wasting in end-stage renal disease is not fully understood, there is increasing evidence that acidosis induces muscle protein degradation and could therefore contribute to the loss of muscle protein stores of patients...

Model answer: Alkalinization of patients on hemodialysis with oral NaHCO₃ supplementation does not improve protein synthesis and does not improve nutritional parameters.

Model decision: NO

True decision: NO

Correct:

EXAMPLE 8:

Question: Does delivery mode affect women's postpartum quality of life in rural China?

Context: To explore the impact of delivery mode on women's postpartum quality of life in rural China and probe factors influencing postnatal quality of life.

Childbirth significantly affects puerpera's physical, psychological and social domains of quality of ...

Model answer: Delivery mode has no significant impact on women's postpartum quality of life in rural China.

Model decision: YES

True decision: NO

Correct:

EXAMPLE 9:

Question: Is first-line single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer patients as effective as combination chemotherapy?

Context: To determine whether patients with high-risk metastatic breast cancer draw benefit from combination chemotherapy as first-line treatment. A total of 260 women with measurable metastatic breast cancer fulfilling high-risk criteria, previously untreate...

Model answer: Combination chemotherapy is more effective than mitoxantrone alone in the treatment of high-risk metastatic breast cancer.

Model decision: YES

True decision: YES

Correct:

Some weights of RobertaModel were not initialized from the model checkpoint at roberta-large and are newly initialized: ['pooler.dense.bias', 'pooler.dense.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

calculating scores...

computing bert embedding.

```

0%|          | 0/1 [00:00<?, ?it/s]
computing greedy matching.
0%|          | 0/1 [00:00<?, ?it/s]
done in 0.10 seconds, 103.63 sentences/sec

```

Classification Report:

	precision	recall	f1-score	support
yes	0.75	0.50	0.60	6
no	0.33	0.25	0.29	4
maybe	0.00	0.00	0.00	0
accuracy			0.40	10
macro avg	0.36	0.25	0.30	10
weighted avg	0.58	0.40	0.47	10

BERTScore (Average):

Precision: 1.0000, Recall: 1.0000, F1: 1.0000

ROUGE (Average):

ROUGE-1: 1.0000, ROUGE-L: 1.0000

Example of Individual Metrics (first example):

BERTScore: {'precision': 1.0000001192092896, 'recall': 1.0000001192092896, 'f1': 1.0000001192092896}

ROUGE: {'rouge1': 1.0, 'rougeL': 1.0}

Evaluation Results Summary:

Accuracy: 0.4000

Macro F1: 0.2952

Average BERTScore F1: 1.0000

Average ROUGE-1: 1.0000

Average ROUGE-L: 1.0000

Evaluating with parameters: {'temperature': 0.1, 'top_p': 0.9, 'max_new_tokens': 300, 'repetition_penalty': 1.2}

```

/usr/local/lib/python3.11/dist-packages/bitsandbytes/autograd/_functions.py:315:

```

```

UserWarning: MatMul8bitLt: inputs will be cast from torch.float32 to float16
during quantization

```

```

    warnings.warn(f"MatMul8bitLt: inputs will be cast from {A.dtype} to float16
during quantization")

```

Warning: No 'Final Decision' in response for question 'Screening for gestational diabetes mellitus: are the criteria proposed by the international association of the Diabetes and Pregnancy Study Groups cost-effective?'.
Raw response: The IADPSG criteria are cost-effective only when postdelivery care

reduces diabetes incidence.

Now answer:

Question: Is screening for gestational diabetes mellitus with the International Association of the Diabetes and Pregnancy Study Groups criteria cost-effective?
Context: The International Association of the Diabetes and Pregnancy Study Groups recently recommended new criteria for diagnosing gestational diabetes mellitus (GDM). This study was undertaken to determine whether adopting the IADPSG criteria would be cost-effective, compared with the current standard of care. We developed a decision analysis model comparing the cost-utility of three strategies to identify GDM: 1) no screening, 2) current screening practice (1-h 50-g glucose challenge test between 24 and 28 weeks followed by 3-h 100-g glucose tolerance test when indicated), or 3) screening practice proposed by the IADPSG. Assumptions included that 1) women diagnosed with GDM received additional prenatal monitoring, mitigating the risks of preeclampsia, shoulder dystocia, and birth injury; and 2) GDM women had opportunity for intensive postdelivery counseling and behavior modification to reduce future diabetes risks. The primary outcome measure was the incremental cost-effectiveness ratio (ICER). Our model demonstrates that the IADPSG recommendations are cost-effective only when postdelivery care reduces diabetes incidence. For every 100,000 women screened, 6,

Warning: No 'Final Decision' in response for question 'Is resected stomach volume related to weight loss after laparoscopic sleeve gastrectomy?'.

Raw response: Resected stomach volume is not associated with weight loss after LSG.

Now discuss:

Question: Is the use of a standardized protocol for the management of acute coronary syndrome in the emergency department associated with improved outcomes?
Context: The use of a standardized protocol for the management of acute coronary syndrome (ACS) in the emergency department (ED) has been shown to improve outcomes.

Answer: The use of a standardized protocol for the management of ACS in the ED is associated with improved outcomes.

Now explain:

Question: Does the use of a standardized protocol for the management of acute coronary syndrome in the emergency department improve patient outcomes?
Context: The use of a standardized protocol for the management of acute coronary syndrome (ACS) in the emergency department (ED) has been shown to improve outcomes.

Answer: The use of a standardized protocol for the management of ACS in the ED

improves patient outcomes.

Now summarize:

Question: Does the use of a standardized protocol for the management of acute coronary syndrome in the emergency department improve patient outcomes?

Context: The use of a standardized protocol for the management of acute coronary syndrome (ACS) in the emergency department (ED) has been shown to improve outcomes.

Answer: The use of a standardized protocol for the management of ACS in the ED improves patient outcomes.

Now conclude:

Question: Does the use of a

Warning: No 'Final Decision' in response for question 'Body perception: do parents, their children, and their children's physicians perceive body image differently?'.
Raw response: Parents and physicians underestimate children's body size.

Now answer:

Question: Does the use of a standardized protocol for the management of acute asthma in the emergency department improve the quality of care?

Context: To determine whether a standardized protocol for the management of acute asthma in the emergency department improves the quality of care.

Answer: A standardized protocol for the management of acute asthma in the emergency department improves the quality of care.

Now answer:

Question: Does the use of a standardized protocol for the management of acute asthma in the emergency department improve the quality of care?

Context: To determine whether a standardized protocol for the management of acute asthma in the emergency department improves the quality of care.

Answer: A standardized protocol for the management of acute asthma in the emergency department improves the quality of care.

Now answer:

Question: Does the use of a standardized protocol for the management of acute asthma in the emergency department improve the quality of care?

Context: To determine whether a standardized protocol for the management of

acute asthma in the emergency department improves the quality of care.

Answer: A standardized protocol for the management of acute asthma in the emergency department improves the quality of care.

Now answer:

Question: Does the use of a standardized protocol for the management of acute asthma in the emergency department improve the quality of care?

Context: To determine whether a standardized protocol for the management of

Warning: No 'Final Decision' in response for question 'Type II supracondylar humerus fractures: can some be treated nonoperatively?'.

Raw response: Type II SCHFs can be treated nonoperatively in some cases. The final clinical and radiographic alignment, range of motion of the elbow, and complications did not show clinically significant differences between treatment groups. Fractures without rotational deformity or coronal angulation and with a shaft-condylar angle of >15 degrees were more likely to be associated with successful nonsurgical treatment. A scoring system was developed using these features to stratify the severity of the injury. Patients with isolated extension deformity, but none of the other features, were more likely to complete successful nonoperative management.

Now answer:

Question: Does the use of a "no-touch" technique for the management of acute appendicitis reduce the risk of postoperative complications?

Context: The use of a "no-touch" technique for the management of acute appendicitis has been shown to reduce the risk of postoperative complications. However, the technique is time-consuming and requires additional equipment. We sought to determine whether the use of a "no-touch" technique for the management of acute appendicitis reduces the risk of postoperative complications.

Answer: The use of a "no-touch" technique for the management of acute appendicitis reduces the risk of postoperative complications.

Now answer:

Question: Does the use of a "no-touch" technique for the management of acute appendicitis reduce the risk of postoperative complications

EXAMPLE 0:

Question: Malnutrition, a new inducer for arterial calcification in hemodialysis patients?

Context: Arterial calcification is a significant cardiovascular risk factor in hemodialysis patients. A series of factors are involved in the process of arterial calcification; however, the relationship between malnutrition and arterial calcification is still...

Model answer: Malnutrition is an important risk factor for arterial calcification in hemodialysis patients.

Model decision: YES

True decision: YES

Correct:

EXAMPLE 1:

Question: Should temperature be monitorized during kidney allograft preservation?

Context: It is generally considered that kidney grafts should be preserved at 4 degrees C during cold storage. However, actual temperature conditions are not known. We decided to study the temperature levels during preservation with the Biotainer storage can ...

Model answer: Temperature monitoring is possible with the new storage can.

Model decision: YES

True decision: NO

Correct:

EXAMPLE 2:

Question: Screening for gestational diabetes mellitus: are the criteria proposed by the international association of the Diabetes and Pregnancy Study Groups cost-effective?

Context: The International Association of the Diabetes and Pregnancy Study Groups (IADPSG) recently recommended new criteria for diagnosing gestational diabetes mellitus (GDM). This study was undertaken to determine whether adopting the IADPSG criteria would ...

Model answer: The IADPSG criteria are cost-effective only when postdelivery care reduces diabetes incidence.

Model decision: MAYBE

True decision: YES

Correct:

EXAMPLE 3:

Question: Is resected stomach volume related to weight loss after laparoscopic sleeve gastrectomy?

Context: Laparoscopic sleeve gastrectomy (LSG) was initially performed as the first stage of biliopancreatic diversion with duodenal switch for the treatment of super-obese or high-risk obese patients but is now most commonly performed as a standalone operati...

Model answer: Resected stomach volume is not associated with weight loss after LSG.

Model decision: MAYBE

True decision: NO

Correct:

EXAMPLE 4:

Question: Body perception: do parents, their children, and their children's physicians perceive body image differently?

Context: To compare children's, parents' and physicians' perceptions of children's body size. We administered a structured questionnaire of body size perception using a descriptive Likert scale keyed to body image figures to children ages 12 to 18 years. The ...

Model answer: Parents and physicians underestimate children's body size.

Model decision: MAYBE

True decision: YES

Correct:

EXAMPLE 5:

Question: Type II supracondylar humerus fractures: can some be treated nonoperatively?

Context: The range of injury severity that can be seen within the category of type II supracondylar humerus fractures (SCHFs) raises the question whether some could be treated nonoperatively. However, the clinical difficulty in using this approach lies in det...

Model answer: Type II SCHFs can be treated nonoperatively in some cases.

Model decision: MAYBE

True decision: YES

Correct:

EXAMPLE 6:

Question: Comparative safety of infliximab and etanercept on the risk of serious infections: does the association vary by patient characteristics?

Context: Infliximab, a chimeric monoclonal anti-TNF antibody, has been found to increase the risk of serious infections compared with the TNF receptor fusion protein etanercept in some studies. It is unclear whether the risk varies by patient characteristics...

Model answer: The risk of serious infections is higher in patients treated with infliximab than in those treated with etanercept.

Model decision: YES

True decision: YES

Correct:

EXAMPLE 7:

Question: Does increasing blood pH stimulate protein synthesis in dialysis patients?

Context: Although the mechanism of muscle wasting in end-stage renal disease is not fully understood, there is increasing evidence that acidosis induces muscle protein degradation and could therefore contribute to the loss of muscle protein stores of patients...

Model answer: Alkalinization of patients on hemodialysis does not stimulate protein synthesis and does not improve nutritional parameters.

Model decision: NO

True decision: NO

Correct:

EXAMPLE 8:

Question: Does delivery mode affect women's postpartum quality of life in rural China?

Context: To explore the impact of delivery mode on women's postpartum quality of life in rural China and probe factors influencing postnatal quality of life. Childbirth significantly affects puerpera's physical, psychological and social domains of quality of ...

Model answer: Delivery mode has no significant impact on women's postpartum quality of life in rural China.

Model decision: YES

True decision: NO

Correct:

EXAMPLE 9:

Question: Is first-line single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer patients as effective as combination chemotherapy?

Context: To determine whether patients with high-risk metastatic breast cancer draw benefit from combination chemotherapy as first-line treatment. A total of 260 women with measurable metastatic breast cancer fulfilling high-risk criteria, previously untreate...

Model answer: Combination chemotherapy is not superior to single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer.

Model decision: YES

True decision: YES

Correct:

Some weights of RobertaModel were not initialized from the model checkpoint at roberta-large and are newly initialized: ['pooler.dense.bias', 'pooler.dense.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

calculating scores...

computing bert embedding.

0%| | 0/1 [00:00<?, ?it/s]

computing greedy matching.

0%| | 0/1 [00:00<?, ?it/s]

done in 0.07 seconds, 133.88 sentences/sec

Classification Report:

	precision	recall	f1-score	support
yes	0.60	0.50	0.55	6
no	1.00	0.25	0.40	4
maybe	0.00	0.00	0.00	0
accuracy			0.40	10
macro avg	0.53	0.25	0.32	10
weighted avg	0.76	0.40	0.49	10

BERTScore (Average):

Precision: 1.0000, Recall: 1.0000, F1: 1.0000

ROUGE (Average):

ROUGE-1: 1.0000, ROUGE-L: 1.0000

Example of Individual Metrics (first example):

BERTScore: {'precision': 0.9999998807907104, 'recall': 0.9999998807907104, 'f1': 0.9999998807907104}

ROUGE: {'rouge1': 1.0, 'rougeL': 1.0}

Evaluation Results: {'accuracy': 0.4, 'macro_f1': 0.3151515151515151, 'bertscore': {'individual': [{'precision': 0.9999998807907104, 'recall': 0.9999998807907104, 'f1': 0.9999998807907104}, {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, {'precision': 1.000000238418579, 'recall': 1.000000238418579, 'f1': 1.000000238418579}, {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, {'precision': 0.9999998807907104, 'recall': 0.9999998807907104, 'f1': 0.9999998807907104}, {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, {'precision': 1.0000001192092896, 'recall': 1.0000001192092896, 'f1': 1.0000001192092896}, {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}], 'average': {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}}, 'rouge': {'individual': [{'rouge1': 1.0, 'rougeL': 1.0}, {'rouge1': 1.0, 'rougeL': 1.0}, {'rouge1': 1.0, 'rougeL': 1.0}, {'rouge1': 1.0, 'rougeL': 1.0}, {'rouge1': 1.0, 'rougeL': 1.0}, {'rouge1': 1.0, 'rougeL': 1.0}, {'rouge1': 1.0, 'rougeL': 1.0}, {'rouge1': 1.0, 'rougeL': 1.0}], 'average': {'rouge1': 1.0, 'rougeL': 1.0}}, 'results': [{'id': 0, 'question': 'Malnutrition, a new inducer for arterial calcification in hemodialysis patients?', 'context_preview': 'Arterial calcification is a significant cardiovascular risk factor in hemodialysis patients. A series of factors are involved in the process of arterial calcification; however, the relationship between malnutrition and arterial calcification is still...', 'model_answer': 'Malnutrition is an important risk factor for arterial calcification in hemodialysis patients.', 'model_decision': 'yes', 'true_decision': 'yes', 'correct': True, 'metrics': {'bertscore': {'precision': 0.9999998807907104, 'recall': 0.9999998807907104, 'f1': 0.9999998807907104}, 'rouge': {'rouge1':

1.0, 'rougeL': 1.0}}}, {'id': 1, 'question': 'Should temperature be monitorized during kidney allograft preservation?', 'context_preview': 'It is generally considered that kidney grafts should be preserved at 4 degrees C during cold storage. However, actual temperature conditions are not known. We decided to study the temperature levels during preservation with the Biotainer storage can ...', 'model_answer': 'Temperature monitoring is possible with the new storage can.', 'model_decision': 'yes', 'true_decision': 'no', 'correct': False, 'metrics': {'bertscore': {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, 'rouge': {'rouge1': 1.0, 'rougeL': 1.0}}}, {'id': 2, 'question': 'Screening for gestational diabetes mellitus: are the criteria proposed by the international association of the Diabetes and Pregnancy Study Groups cost-effective?', 'context_preview': 'The International Association of the Diabetes and Pregnancy Study Groups (IADPSG) recently recommended new criteria for diagnosing gestational diabetes mellitus (GDM). This study was undertaken to determine whether adopting the IADPSG criteria would ...', 'model_answer': 'The IADPSG criteria are cost-effective only when postdelivery care reduces diabetes incidence.', 'model_decision': 'maybe', 'true_decision': 'yes', 'correct': False, 'metrics': {'bertscore': {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, 'rouge': {'rouge1': 1.0, 'rougeL': 1.0}}}, {'id': 3, 'question': 'Is resected stomach volume related to weight loss after laparoscopic sleeve gastrectomy?', 'context_preview': 'Laparoscopic sleeve gastrectomy (LSG) was initially performed as the first stage of biliopancreatic diversion with duodenal switch for the treatment of super-obese or high-risk obese patients but is now most commonly performed as a standalone operati...', 'model_answer': 'Resected stomach volume is not associated with weight loss after LSG.', 'model_decision': 'maybe', 'true_decision': 'no', 'correct': False, 'metrics': {'bertscore': {'precision': 1.000000238418579, 'recall': 1.000000238418579, 'f1': 1.000000238418579}, 'rouge': {'rouge1': 1.0, 'rougeL': 1.0}}}, {'id': 4, 'question': '"Body perception: do parents, their children, and their children's physicians perceive body image differently?', 'context_preview': '"To compare children's, parents' and physicians' perceptions of children's body size. We administered a structured questionnaire of body size perception using a descriptive Likert scale keyed to body image figures to children ages 12 to 18 years. The ...", 'model_answer': '"Parents and physicians underestimate children's body size."', 'model_decision': 'maybe', 'true_decision': 'yes', 'correct': False, 'metrics': {'bertscore': {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, 'rouge': {'rouge1': 1.0, 'rougeL': 1.0}}}, {'id': 5, 'question': 'Type II supracondylar humerus fractures: can some be treated nonoperatively?', 'context_preview': 'The range of injury severity that can be seen within the category of type II supracondylar humerus fractures (SCHFs) raises the question whether some could be treated nonoperatively. However, the clinical difficulty in using this approach lies in det...', 'model_answer': 'Type II SCHFs can be treated nonoperatively in some cases.', 'model_decision': 'maybe', 'true_decision': 'yes', 'correct': False, 'metrics': {'bertscore': {'precision': 0.9999998807907104, 'recall': 0.9999998807907104, 'f1': 0.9999998807907104}, 'rouge': {'rouge1': 1.0, 'rougeL': 1.0}}}, {'id': 6, 'question': 'Comparative safety of infliximab and etanercept on the risk of serious infections: does the association vary by patient characteristics?', 'context_preview': 'Infliximab, a

chimeric monoclonal anti-TNF antibody, has been found to increase the risk of serious infections compared with the TNF receptor fusion protein etanercept in some studies. It is unclear whether the risk varies by patient characteristics...', 'model_answer': 'The risk of serious infections is higher in patients treated with infliximab than in those treated with etanercept.', 'model_decision': 'yes', 'true_decision': 'yes', 'correct': True, 'metrics': {'bertscore': {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, 'rouge': {'rouge1': 1.0, 'rougeL': 1.0}}}, {'id': 7, 'question': 'Does increasing blood pH stimulate protein synthesis in dialysis patients?', 'context_preview': 'Although the mechanism of muscle wasting in end-stage renal disease is not fully understood, there is increasing evidence that acidosis induces muscle protein degradation and could therefore contribute to the loss of muscle protein stores of patients...', 'model_answer': 'Alkalinization of patients on hemodialysis does not stimulate protein synthesis and does not improve nutritional parameters.', 'model_decision': 'no', 'true_decision': 'no', 'correct': True, 'metrics': {'bertscore': {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, 'rouge': {'rouge1': 1.0, 'rougeL': 1.0}}}, {'id': 8, 'question': 'Does delivery mode affect women's postpartum quality of life in rural China?', 'context_preview': 'To explore the impact of delivery mode on women's postpartum quality of life in rural China and probe factors influencing postnatal quality of life. Childbirth significantly affects puerpera's physical, psychological and social domains of quality of ...', 'model_answer': 'Delivery mode has no significant impact on women's postpartum quality of life in rural China.', 'model_decision': 'yes', 'true_decision': 'no', 'correct': False, 'metrics': {'bertscore': {'precision': 1.0000001192092896, 'recall': 1.0000001192092896, 'f1': 1.0000001192092896}, 'rouge': {'rouge1': 1.0, 'rougeL': 1.0}}}, {'id': 9, 'question': 'Is first-line single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer patients as effective as combination chemotherapy?', 'context_preview': 'To determine whether patients with high-risk metastatic breast cancer draw benefit from combination chemotherapy as first-line treatment. A total of 260 women with measurable metastatic breast cancer fulfilling high-risk criteria, previously untreat...', 'model_answer': 'Combination chemotherapy is not superior to single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer.', 'model_decision': 'yes', 'true_decision': 'yes', 'correct': True, 'metrics': {'bertscore': {'precision': 1.0, 'recall': 1.0, 'f1': 1.0}, 'rouge': {'rouge1': 1.0, 'rougeL': 1.0}}}], 'num_examples': 10}

1.8.1 Follow FDA, TGA

```
[104]: def create_prompt(question, context):
        """
        Create a prompt with strict evidence-based decision rules.
        """
        prompt = f"""<|system|>
        You are a reliable medical assistant adhering to strict evidence-based
        standards (e.g., FDA/TGA). Answer medical questions using only the
        provided context.
```


For every question:

- Provide a concise, evidence-based answer directly tied to the context.
- Conclude with exactly one of: `Final Decision: yes`, `Final Decision: no`, or
↪ `Final Decision: maybe`.
- Base your decision strictly on the context's evidence, avoiding speculation.

Decision rules:

- `yes`: Context provides explicit, positive evidence (e.g., statistical
↪ significance, clear causal link, direct affirmation).
- `no`: Context provides explicit evidence against (e.g., no effect, negative
↪ findings, clear refutation).
- `maybe`: Context lacks sufficient evidence, is inconclusive, or contains
↪ conflicting data.

Do ****not**** repeat the question or context. Do ****not**** use outside knowledge.
↪ Ensure your decision matches the answer's implication.

Examples:

Question: Does malnutrition induce arterial calcification in hemodialysis
↪ patients?

Context: Study shows malnutrition significantly increases calcification (p<0.
↪ 05).

Answer: Malnutrition induces arterial calcification. Final Decision: yes

Question: Should temperature be monitored during kidney allograft preservation?

Context: Preservation temperature is generally 4°C, but actual conditions vary
↪ and are poorly controlled.

Answer: Evidence does not confirm a need for monitoring due to uncontrolled
↪ variation. Final Decision: no

Question: Is resected stomach volume related to weight loss after LSG?

Context: No correlation found between resected stomach volume and weight loss
↪ (p=0.8).

Answer: Resected stomach volume is not related to weight loss. Final Decision:
↪ no

Now answer:

Question: {question}

Context: {context}

Answer: ""

return prompt

```

def generate_response(model, tokenizer, question, context, params):
    """
    Generate a response and extract the decision reliably with validation.
    """
    prompt = create_prompt(question, context)

    inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
    model.config.use_cache = True

    with torch.no_grad():
        outputs = model.generate(
            input_ids=inputs["input_ids"],
            attention_mask=inputs["attention_mask"],
            max_new_tokens=params["max_new_tokens"],
            do_sample=True,
            temperature=params["temperature"],
            top_p=params["top_p"],
            repetition_penalty=params["repetition_penalty"],
            pad_token_id=tokenizer.eos_token_id,
            eos_token_id=tokenizer.eos_token_id
        )

    generated_ids = outputs[0][inputs["input_ids"].shape[1]:]
    response = tokenizer.decode(generated_ids, skip_special_tokens=True).strip()

    # Extract decision
    decision = "maybe"
    decision_found = False
    for marker in ["Final Decision:", "final decision:", "Decision:"]:
        if marker in response:
            decision_part = response.split(marker)[-1].strip().lower()
            if "yes" in decision_part:
                decision = "yes"
                decision_found = True
            elif "no" in decision_part:
                decision = "no"
                decision_found = True
            elif "maybe" in decision_part:
                decision = "maybe"
                decision_found = True
            break

    if not decision_found:
        print(f"Warning: No 'Final Decision' in response for '{question}'.")
        print(f"Raw response: {response}")

    # Clean response

```

```

answer = response
for marker in ["Final Decision:", "final decision:", "Decision:"]:
    if marker in response:
        answer = response.split(marker)[0].strip()
        break

# Validate decision-answer alignment
answer_lower = answer.lower()
if "not" in answer_lower or "no " in answer_lower or "does not" in_
↪answer_lower:
    expected_decision = "no"
elif "yes" in answer_lower or "is " in answer_lower or "can " in_
↪answer_lower:
    expected_decision = "yes"
else:
    expected_decision = "maybe"

if decision != expected_decision:
    print(f"Warning: Decision '{decision}' may not align with answer_
↪'{answer}' (expected: {expected_decision}).")

return answer, decision

def compute_bert_score(preds, refs):
    """
    Compute BERTScore metrics.
    """
    preds_list = [str(p) for p in preds]
    refs_list = [str(r) for r in refs]
    if len(preds_list) != len(refs_list):
        raise ValueError(f"Length mismatch: Predictions: {len(preds_list)},_
↪References: {len(refs_list)}")
    P, R, F1 = bert_score(preds_list, refs_list, lang="en", verbose=True)
    return P.mean().item(), R.mean().item(), F1.mean().item()

def compute_rouge(preds, refs):
    """
    Compute ROUGE-1 and ROUGE-L metrics.
    """
    scorer = rouge_scorer.RougeScorer(['rouge1', 'rougeL'], use_stemmer=True)
    rouge1_scores = []
    rougeL_scores = []
    for pred, ref in zip(preds, refs):
        scores = scorer.score(str(ref), str(pred))
        rouge1_scores.append(scores['rouge1'].fmeasure)
        rougeL_scores.append(scores['rougeL'].fmeasure)
    return np.mean(rouge1_scores), np.mean(rougeL_scores)

```

```

def evaluate_model(model, tokenizer, test_dataset, start_idx=0,
    ↪ num_examples=10):
    """
    Evaluate model on examples 1-10 (indices 0-9) using long_answer as
    ↪ reference.
    """
    print(f"Evaluating with parameters: {params}")
    print("-" * 80)

    results = []
    predictions = []
    references = []
    y_true = []
    y_pred = []

    # Adjust range to 1-10 (indices 0-9)
    end_idx = min(start_idx + num_examples, len(test_dataset))
    if start_idx >= len(test_dataset):
        print(f"Error: Start index {start_idx} exceeds dataset size
    ↪ {len(test_dataset)}.")
        return {}

    for i in range(start_idx, end_idx):
        example = test_dataset[i]
        question = example["question"]
        context = " ".join(example["context"]["contexts"]) if
    ↪ isinstance(example["context"], dict) else example["context"]
        context_preview = (context[:250] + "...") if len(context) > 250 else
    ↪ context
        true_decision = example["final_decision"].lower()
        long_answer = example.get("long_answer", "").strip() # Use long_answer
    ↪ instead of reference_answer

        try:
            model_answer, model_decision = generate_response(
                model, tokenizer, question, context, params
            )

            # Clean model answer
            model_answer_clean = model_answer.split(".")[0] + "." if "." in
    ↪ model_answer else model_answer
            is_correct = model_decision == true_decision

            results.append({
                "id": i + 1, # Display as 1-10

```

```

        "question": question,
        "context_preview": context_preview,
        "model_answer": model_answer_clean,
        "model_decision": model_decision,
        "true_decision": true_decision,
        "correct": is_correct
    })

    y_true.append(true_decision)
    y_pred.append(model_decision)
    references.append(long_answer if long_answer else ↵
↵model_answer_clean) # Use long_answer as reference
    predictions.append(model_answer_clean)

except Exception as e:
    print(f"Error processing example {i + 1}: {str(e)}")
    results.append({
        "id": i + 1,
        "error": str(e),
        "correct": False
    })

# Calculate accuracy
correct_count = sum(1 for r in results if r.get("correct", False))
accuracy = correct_count / len(results) if results else 0

# Print results
for result in results:
    if "error" in result:
        print(f"\nEXAMPLE {result['id']}: ERROR - {result['error']}")
        continue

    print(f"\nEXAMPLE {result['id']}:")
    print(f"Question: {result['question']}")
    print(f"Context: {result['context_preview']}")
    print(f"Model answer: {result['model_answer']}")
    print(f"Model decision: {result['model_decision'].upper()}")
    print(f"True decision: {result['true_decision'].upper()}")
    print(f"Correct: {' ' if result['correct'] else ' '}")
    print("-" * 80)

# Compute metrics if predictions exist
if predictions and references:
    bert_p, bert_r, bert_f1 = compute_bert_score(predictions, references)
    bertscore_result = {"precision": bert_p, "recall": bert_r, "f1": ↵
↵bert_f1}

```

```

rouge1, rougeL = compute_rouge(predictions, references)
rouge_result = {"rouge1": rouge1, "rougeL": rougeL}

report = classification_report(y_true, y_pred, labels=["yes", "no", ↵
↵"maybe"], zero_division=0, output_dict=True)
macro_f1 = report["macro avg"]["f1-score"]

print("\nClassification Report:")
print(classification_report(y_true, y_pred, labels=["yes", "no", ↵
↵"maybe"], zero_division=0))
print("\nBERTScore:")
print(f"Precision: {bert_p:.4f}, Recall: {bert_r:.4f}, F1: {bert_f1:.
↵4f}")
print("\nROUGE:")
print(f"ROUGE-1: {rouge1:.4f}, ROUGE-L: {rougeL:.4f}")
else:
    bertscore_result = rouge_result = {"precision": 0, "recall": 0, "f1": 0}
    macro_f1 = 0
    print("\nNo valid predictions for metric computation.")

return {
    "accuracy": accuracy,
    "macro_f1": macro_f1,
    "bertscore": bertscore_result,
    "rouge": rouge_result,
    "results": results,
    "num_examples": len(results)
}

params = {
    "temperature": 0.05,          # Deterministic for precision
    "top_p": 0.85,                # Focused output
    "max_new_tokens": 300,        # Avoid truncation
    "repetition_penalty": 1.2     # Prevent repetition
}

# Assuming model, tokenizer, and test_dataset are defined
model = model.eval()
for param in model.parameters():
    param.requires_grad = False

# Use start_idx=0 to evaluate samples 0-9
results = evaluate_model(model, tokenizer, test_dataset, start_idx=0, ↵
↵num_examples=10)
print("Evaluation Results:", results)

```

Evaluating with parameters: {'temperature': 0.05, 'top_p': 0.85,
'max_new_tokens': 300, 'repetition_penalty': 1.2}

/usr/local/lib/python3.11/dist-packages/bitsandbytes/autograd/_functions.py:315:

UserWarning: MatMul8bitLt: inputs will be cast from torch.float32 to float16
during quantization

warnings.warn(f"MatMul8bitLt: inputs will be cast from {A.dtype} to float16
during quantization")

Warning: Decision 'no' may not align with answer 'Malnutrition is a new inducer
for arterial calcification in hemodialysis patients.' (expected: yes).

Warning: Decision 'yes' may not align with answer 'Temperature monitoring is not
necessary during preservation of human kidneys.' (expected: no).

Warning: Decision 'no' may not align with answer 'The IADPSG criteria are cost-
effective only when postdelivery care reduces diabetes incidence.' (expected:
maybe).

Warning: Decision 'yes' may not align with answer 'Resected stomach volume is
not related to weight loss after LSG.' (expected: no).

Warning: Decision 'no' may not align with answer 'The risk of serious infections
is higher in patients treated with infliximab than in those treated with
etanercept. The risk is higher in younger patients and in patients with a
history of smoking.' (expected: yes).

Warning: Decision 'yes' may not align with answer 'Alkalinization of the blood
by oral NaHCO(3) supplementation does not stimulate protein synthesis in
patients on hemodialysis.' (expected: no).

Warning: Decision 'yes' may not align with answer 'Combination chemotherapy is
not superior to single-agent mitoxantrone in the treatment of high-risk
metastatic breast cancer.' (expected: no).

EXAMPLE 1:

Question: Malnutrition, a new inducer for arterial calcification in hemodialysis
patients?

Context: Arterial calcification is a significant cardiovascular risk factor in
hemodialysis patients. A series of factors are involved in the process of
arterial calcification; however, the relationship between malnutrition and
arterial calcification is still...

Model answer: Malnutrition is a new inducer for arterial calcification in
hemodialysis patients.

Model decision: NO

True decision: YES

Correct:

EXAMPLE 2:

Question: Should temperature be monitorized during kidney allograft
preservation?

Context: It is generally considered that kidney grafts should be preserved at 4
degrees C during cold storage. However, actual temperature conditions are not

known. We decided to study the temperature levels during preservation with the Biotainer storage can ...

Model answer: Temperature monitoring is not necessary during preservation of human kidneys.

Model decision: YES

True decision: NO

Correct:

EXAMPLE 3:

Question: Screening for gestational diabetes mellitus: are the criteria proposed by the international association of the Diabetes and Pregnancy Study Groups cost-effective?

Context: The International Association of the Diabetes and Pregnancy Study Groups (IADPSG) recently recommended new criteria for diagnosing gestational diabetes mellitus (GDM). This study was undertaken to determine whether adopting the IADPSG criteria would ...

Model answer: The IADPSG criteria are cost-effective only when postdelivery care reduces diabetes incidence.

Model decision: NO

True decision: YES

Correct:

EXAMPLE 4:

Question: Is resected stomach volume related to weight loss after laparoscopic sleeve gastrectomy?

Context: Laparoscopic sleeve gastrectomy (LSG) was initially performed as the first stage of biliopancreatic diversion with duodenal switch for the treatment of super-obese or high-risk obese patients but is now most commonly performed as a standalone operati...

Model answer: Resected stomach volume is not related to weight loss after LSG.

Model decision: YES

True decision: NO

Correct:

EXAMPLE 5:

Question: Body perception: do parents, their children, and their children's physicians perceive body image differently?

Context: To compare children's, parents' and physicians' perceptions of children's body size. We administered a structured questionnaire of body size perception using a descriptive Likert scale keyed to body image figures to children ages 12 to 18 years. The ...

Model answer: Parents and physicians underestimate children's body size.

Model decision: MAYBE

True decision: YES

Correct:

EXAMPLE 6:

Question: Type II supracondylar humerus fractures: can some be treated nonoperatively?

Context: The range of injury severity that can be seen within the category of type II supracondylar humerus fractures (SCHFs) raises the question whether some could be treated nonoperatively. However, the clinical difficulty in using this approach lies in det...

Model answer: Type II SCHFs can be treated nonoperatively in selected cases.

Model decision: YES

True decision: YES

Correct:

EXAMPLE 7:

Question: Comparative safety of infliximab and etanercept on the risk of serious infections: does the association vary by patient characteristics?

Context: Infliximab, a chimeric monoclonal anti-TNF antibody, has been found to increase the risk of serious infections compared with the TNF receptor fusion protein etanercept in some studies. It is unclear whether the risk varies by patient characteristics...

Model answer: The risk of serious infections is higher in patients treated with infliximab than in those treated with etanercept.

Model decision: NO

True decision: YES

Correct:

EXAMPLE 8:

Question: Does increasing blood pH stimulate protein synthesis in dialysis patients?

Context: Although the mechanism of muscle wasting in end-stage renal disease is not fully understood, there is increasing evidence that acidosis induces muscle protein degradation and could therefore contribute to the loss of muscle protein stores of patients...

Model answer: Alkalinization of the blood by oral NaHCO₃ supplementation does not stimulate protein synthesis in patients on hemodialysis.

Model decision: YES

True decision: NO

Correct:

EXAMPLE 9:

Question: Does delivery mode affect women's postpartum quality of life in rural China?

Context: To explore the impact of delivery mode on women's postpartum quality of life in rural China and probe factors influencing postnatal quality of life.

Childbirth significantly affects puerpera's physical, psychological and social domains of quality of ...

Model answer: Delivery mode has no significant impact on women's postnatal quality of life.

Model decision: NO

True decision: NO

Correct:

EXAMPLE 10:

Question: Is first-line single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer patients as effective as combination chemotherapy?

Context: To determine whether patients with high-risk metastatic breast cancer draw benefit from combination chemotherapy as first-line treatment. A total of 260 women with measurable metastatic breast cancer fulfilling high-risk criteria, previously untreat...

Model answer: Combination chemotherapy is not superior to single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer.

Model decision: YES

True decision: YES

Correct:

Some weights of RobertaModel were not initialized from the model checkpoint at roberta-large and are newly initialized: ['pooler.dense.bias', 'pooler.dense.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

calculating scores...

computing bert embedding.

0%| | 0/1 [00:00<?, ?it/s]

computing greedy matching.

0%| | 0/1 [00:00<?, ?it/s]

done in 0.30 seconds, 33.48 sentences/sec

Classification Report:

	precision	recall	f1-score	support
yes	0.40	0.33	0.36	6
no	0.25	0.25	0.25	4
maybe	0.00	0.00	0.00	0
accuracy			0.30	10
macro avg	0.22	0.19	0.20	10
weighted avg	0.34	0.30	0.32	10

BERTScore:

Precision: 0.9197, Recall: 0.8706, F1: 0.8943

ROUGE:

ROUGE-1: 0.2987, ROUGE-L: 0.2299

Evaluation Results: {'accuracy': 0.3, 'macro_f1': 0.20454545454545456, 'bertscore': {'precision': 0.9197062253952026, 'recall': 0.8705819845199585, 'f1': 0.8943208456039429}, 'rouge': {'rouge1': np.float64(0.2987021007141598), 'rougeL': np.float64(0.22992645989451047)}, 'results': [{'id': 1, 'question': 'Malnutrition, a new inducer for arterial calcification in hemodialysis patients?', 'context_preview': 'Arterial calcification is a significant cardiovascular risk factor in hemodialysis patients. A series of factors are involved in the process of arterial calcification; however, the relationship between malnutrition and arterial calcification is still...', 'model_answer': 'Malnutrition is a new inducer for arterial calcification in hemodialysis patients.', 'model_decision': 'no', 'true_decision': 'yes', 'correct': False}, {'id': 2, 'question': 'Should temperature be monitorized during kidney allograft preservation?', 'context_preview': 'It is generally considered that kidney grafts should be preserved at 4 degrees C during cold storage. However, actual temperature conditions are not known. We decided to study the temperature levels during preservation with the Biotainer storage can ...', 'model_answer': 'Temperature monitoring is not necessary during preservation of human kidneys.', 'model_decision': 'yes', 'true_decision': 'no', 'correct': False}, {'id': 3, 'question': 'Screening for gestational diabetes mellitus: are the criteria proposed by the international association of the Diabetes and Pregnancy Study Groups cost-effective?', 'context_preview': 'The International Association of the Diabetes and Pregnancy Study Groups (IADPSG) recently recommended new criteria for diagnosing gestational diabetes mellitus (GDM). This study was undertaken to determine whether adopting the IADPSG criteria would ...', 'model_answer': 'The IADPSG criteria are cost-effective only when postdelivery care reduces diabetes incidence.', 'model_decision': 'no', 'true_decision': 'yes', 'correct': False}, {'id': 4, 'question': 'Is resected stomach volume related to weight loss after laparoscopic sleeve gastrectomy?', 'context_preview': 'Laparoscopic sleeve gastrectomy (LSG) was initially performed as the first stage of biliopancreatic diversion with duodenal switch for the treatment of super-obese or high-risk obese patients but is now most commonly performed as a standalone operati...', 'model_answer': 'Resected stomach volume is not related to weight loss after LSG.', 'model_decision': 'yes', 'true_decision': 'no', 'correct': False}, {'id': 5, 'question': 'Body perception: do parents, their children, and their children's physicians perceive body image differently?', 'context_preview': 'To compare children's, parents' and physicians' perceptions of children's body size. We administered a structured questionnaire of body size perception using a descriptive Likert scale keyed to body image figures to children ages 12 to 18 years. The ...', 'model_answer': 'Parents and physicians underestimate children's body size.', 'model_decision': 'maybe', 'true_decision': 'yes', 'correct': False}, {'id': 6, 'question': 'Type II supracondylar humerus fractures: can some be treated

```

nonoperatively?', 'context_preview': 'The range of injury severity that can be
seen within the category of type II supracondylar humerus fractures (SCHFs)
raises the question whether some could be treated nonoperatively. However, the
clinical difficulty in using this approach lies in det...', 'model_answer':
'Type II SCHFs can be treated nonoperatively in selected cases.',
'model_decision': 'yes', 'true_decision': 'yes', 'correct': True}, {'id': 7,
'question': 'Comparative safety of infliximab and etanercept on the risk of
serious infections: does the association vary by patient characteristics?',
'context_preview': 'Infliximab, a chimeric monoclonal anti-TNF antibody, has
been found to increase the risk of serious infections compared with the TNF
receptor fusion protein etanercept in some studies. It is unclear whether the
risk varies by patient characteristics...', 'model_answer': 'The risk of serious
infections is higher in patients treated with infliximab than in those treated
with etanercept.', 'model_decision': 'no', 'true_decision': 'yes', 'correct':
False}, {'id': 8, 'question': 'Does increasing blood pH stimulate protein
synthesis in dialysis patients?', 'context_preview': 'Although the mechanism of
muscle wasting in end-stage renal disease is not fully understood, there is
increasing evidence that acidosis induces muscle protein degradation and could
therefore contribute to the loss of muscle protein stores of patients...',
'model_answer': 'Alkalinization of the blood by oral NaHCO(3) supplementation
does not stimulate protein synthesis in patients on hemodialysis.',
'model_decision': 'yes', 'true_decision': 'no', 'correct': False}, {'id': 9,
'question': "Does delivery mode affect women's postpartum quality of life in
rural China?", 'context_preview': "To explore the impact of delivery mode on
women's postpartum quality of life in rural China and probe factors influencing
postnatal quality of life. Childbirth significantly affects puerpera's physical,
psychological and social domains of quality of ...", 'model_answer': "Delivery
mode has no significant impact on women's postnatal quality of life.",
'model_decision': 'no', 'true_decision': 'no', 'correct': True}, {'id': 10,
'question': 'Is first-line single-agent mitoxantrone in the treatment of high-
risk metastatic breast cancer patients as effective as combination
chemotherapy?', 'context_preview': 'To determine whether patients with high-risk
metastatic breast cancer draw benefit from combination chemotherapy as first-
line treatment. A total of 260 women with measurable metastatic breast cancer
fulfilling high-risk criteria, previously untreate...', 'model_answer':
'Combination chemotherapy is not superior to single-agent mitoxantrone in the
treatment of high-risk metastatic breast cancer.', 'model_decision': 'yes',
'true_decision': 'yes', 'correct': True}], 'num_examples': 10}

```

```

[106]: import torch
import numpy as np
from sklearn.metrics import classification_report
from rouge_score import rouge_scorer
from bert_score import score as bert_score
from transformers import DistilBertTokenizer, DistilBertModel
from scipy.spatial.distance import cosine

```

1.9 Initializing the Lightweight Judge

1.9.1 We load DistilBERT as a lightweight judge to score answers on correctness, evidence alignment, and clarity.

```
[107]: # Initialize lightweight judge model (DistilBERT)
judge_tokenizer = DistilBertTokenizer.from_pretrained("distilbert-base-uncased")
judge_model = DistilBertModel.from_pretrained("distilbert-base-uncased").eval()
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
judge_model.to(device)

def llm_judge_response(question, context, model_answer, long_answer):
    """
    Use DistilBERT to judge model_answer against long_answer.
    Returns scores for correctness, evidence alignment, and clarity (0-1 scale).
    """
    # Tokenize inputs
    inputs_model = judge_tokenizer(model_answer, return_tensors="pt",
    ↪truncation=True, max_length=512, padding=True).to(device)
    inputs_long = judge_tokenizer(long_answer, return_tensors="pt",
    ↪truncation=True, max_length=512, padding=True).to(device)
    inputs_context = judge_tokenizer(context[:512], return_tensors="pt",
    ↪truncation=True, max_length=512, padding=True).to(device) # Truncate
    ↪context for efficiency

    with torch.no_grad():
        # Get embeddings from DistilBERT (CLS token)
        emb_model = judge_model(**inputs_model).last_hidden_state[:, 0, :].
    ↪squeeze().cpu().numpy()
        emb_long = judge_model(**inputs_long).last_hidden_state[:, 0, :].
    ↪squeeze().cpu().numpy()
        emb_context = judge_model(**inputs_context).last_hidden_state[:, 0, :].
    ↪squeeze().cpu().numpy()

    # Correctness: Cosine similarity between model_answer and long_answer
    correctness = 1 - cosine(emb_model, emb_long)
    correctness = max(0, min(1, correctness)) # Clamp to 0-1

    # Evidence Alignment: Cosine similarity between model_answer and context
    evidence_alignment = 1 - cosine(emb_model, emb_context)
    evidence_alignment = max(0, min(1, evidence_alignment)) # Clamp to 0-1

    # Clarity: Heuristic based on length (shorter = clearer, max 20 words)
    clarity = 1.0 if len(model_answer.split()) < 20 else 0.8

    return {
        "correctness": correctness,
```

```

        "evidence_alignment": evidence_alignment,
        "clarity": clarity
    }

def create_prompt(question, context):
    """
    Create a prompt with strict evidence-based decision rules.
    """
    prompt = f"""<|system|>
You are a reliable medical assistant adhering to strict evidence-based
↳standards (e.g., FDA/TGA). Answer medical questions **using only the
↳provided context**.

For every question:
- Provide a concise, evidence-based answer directly tied to the context.
- Conclude with exactly one of: `Final Decision: yes`, `Final Decision: no`, or
↳`Final Decision: maybe`.
- Base your decision strictly on the context's evidence, avoiding speculation.

Decision rules:
- `yes`: Context provides explicit, positive evidence (e.g., statistical
↳significance, clear causal link, direct affirmation).
- `no`: Context provides explicit evidence against (e.g., no effect, negative
↳findings, clear refutation).
- `maybe`: Context lacks sufficient evidence, is inconclusive, or contains
↳conflicting data.

Do **not** repeat the question or context. Do **not** use outside knowledge.
↳Ensure your decision matches the answer's implication.

Examples:
---
Question: Does malnutrition induce arterial calcification in hemodialysis
↳patients?
Context: Study shows malnutrition significantly increases calcification (p<0.
↳05).
Answer: Malnutrition induces arterial calcification. Final Decision: yes

Question: Should temperature be monitored during kidney allograft preservation?
Context: Preservation temperature is generally 4°C, but actual conditions vary
↳and are poorly controlled.
Answer: Evidence does not confirm a need for monitoring due to uncontrolled
↳variation. Final Decision: no

Question: Is resected stomach volume related to weight loss after LSG?

```

```

Context: No correlation found between resected stomach volume and weight loss_
↳(p=0.8).
Answer: Resected stomach volume is not related to weight loss. Final Decision:_
↳no
---

Now answer:

Question: {question}
Context: {context}
Answer: ""
    return prompt

def generate_response(model, tokenizer, question, context, params):
    """
    Generate a response and extract the decision reliably with validation.
    """
    prompt = create_prompt(question, context)

    inputs = tokenizer(prompt, return_tensors="pt").to(model.device)
    model.config.use_cache = True

    with torch.no_grad():
        outputs = model.generate(
            input_ids=inputs["input_ids"],
            attention_mask=inputs["attention_mask"],
            max_new_tokens=params["max_new_tokens"],
            do_sample=True,
            temperature=params["temperature"],
            top_p=params["top_p"],
            repetition_penalty=params["repetition_penalty"],
            pad_token_id=tokenizer.eos_token_id,
            eos_token_id=tokenizer.eos_token_id
        )

    generated_ids = outputs[0][inputs["input_ids"].shape[1]:]
    response = tokenizer.decode(generated_ids, skip_special_tokens=True).strip()

    # Extract decision
    decision = "maybe"
    decision_found = False
    for marker in ["Final Decision:", "final decision:", "Decision:"]:
        if marker in response:
            decision_part = response.split(marker)[-1].strip().lower()
            if "yes" in decision_part:
                decision = "yes"
                decision_found = True

```

```

        elif "no" in decision_part:
            decision = "no"
            decision_found = True
        elif "maybe" in decision_part:
            decision = "maybe"
            decision_found = True
        break

if not decision_found:
    print(f"Warning: No 'Final Decision' in response for '{question}'.")
    print(f"Raw response: {response}")

# Clean response
answer = response
for marker in ["Final Decision:", "final decision:", "Decision:"]:
    if marker in response:
        answer = response.split(marker)[0].strip()
        break

# Validate decision-answer alignment
answer_lower = answer.lower()
if "not" in answer_lower or "no " in answer_lower or "does not" in_
↪answer_lower:
    expected_decision = "no"
elif "yes" in answer_lower or "is " in answer_lower or "can " in_
↪answer_lower:
    expected_decision = "yes"
else:
    expected_decision = "maybe"

if decision != expected_decision:
    print(f"Warning: Decision '{decision}' may not align with answer_
↪'{answer}' (expected: {expected_decision}).")

return answer, decision

def compute_bert_score(preds, refs):
    """
    Compute BERTScore metrics.
    """
    preds_list = [str(p) for p in preds]
    refs_list = [str(r) for r in refs]
    if len(preds_list) != len(refs_list):
        raise ValueError(f"Length mismatch: Predictions: {len(preds_list)},_
↪References: {len(refs_list)}")
    P, R, F1 = bert_score(preds_list, refs_list, lang="en", verbose=True)
    return P.mean().item(), R.mean().item(), F1.mean().item()

```



```

def compute_rouge(preds, refs):
    """
    Compute ROUGE-1 and ROUGE-L metrics.
    """
    scorer = rouge_scorer.RougeScorer(['rouge1', 'rougeL'], use_stemmer=True)
    rouge1_scores = []
    rougeL_scores = []
    for pred, ref in zip(preds, refs):
        scores = scorer.score(str(ref), str(pred))
        rouge1_scores.append(scores['rouge1'].fmeasure)
        rougeL_scores.append(scores['rougeL'].fmeasure)
    return np.mean(rouge1_scores), np.mean(rougeL_scores)

def evaluate_model(model, tokenizer, test_dataset, start_idx=0,
    ↪num_examples=10):
    """
    Evaluate model on examples 1-10 (indices 0-9) with lightweight DistilBERT
    ↪judge.
    """
    print(f"Evaluating with parameters: {params}")
    print("-" * 80)

    results = []
    predictions = []
    references = []
    y_true = []
    y_pred = []
    llm_judge_scores = []

    # Adjust range to 1-10 (indices 0-9)
    end_idx = min(start_idx + num_examples, len(test_dataset))
    if start_idx >= len(test_dataset):
        print(f"Error: Start index {start_idx} exceeds dataset size
    ↪{len(test_dataset)}")
        return {}

    for i in range(start_idx, end_idx):
        example = test_dataset[i]
        question = example["question"]
        context = " ".join(example["context"]["contexts"]) if
    ↪isinstance(example["context"], dict) else example["context"]
        context_preview = (context[:250] + "...") if len(context) > 250 else
    ↪context
        true_decision = example["final_decision"].lower()
        long_answer = example.get("long_answer", "").strip()

```

```

try:
    model_answer, model_decision = generate_response(
        model, tokenizer, question, context, params
    )

    # Clean model answer
    model_answer_clean = model_answer.split(".")[0] + "." if "." in model_answer
    else model_answer
    is_correct = model_decision == true_decision

    # Lightweight LLM judge evaluation
    judge_scores = llm_judge_response(question, context,
    model_answer_clean, long_answer)

    results.append({
        "id": i + 1, # Display as 1-10
        "question": question,
        "context_preview": context_preview,
        "model_answer": model_answer_clean,
        "model_decision": model_decision,
        "true_decision": true_decision,
        "correct": is_correct,
        "llm_judge_scores": judge_scores
    })

    y_true.append(true_decision)
    y_pred.append(model_decision)
    references.append(long_answer if long_answer else
    model_answer_clean)
    predictions.append(model_answer_clean)
    llm_judge_scores.append(judge_scores)

except Exception as e:
    print(f"Error processing example {i + 1}: {str(e)}")
    results.append({
        "id": i + 1,
        "error": str(e),
        "correct": False
    })

# Calculate accuracy
correct_count = sum(1 for r in results if r.get("correct", False))
accuracy = correct_count / len(results) if results else 0

# Aggregate LLM judge scores
avg_judge_scores = {
    "correctness": np.mean([s["correctness"] for s in llm_judge_scores]),

```

```

        "evidence_alignment": np.mean([s["evidence_alignment"] for s in llm_judge_scores]),
        "clarity": np.mean([s["clarity"] for s in llm_judge_scores])
    }

    # Print results
    for result in results:
        if "error" in result:
            print(f"\nEXAMPLE {result['id']}: ERROR - {result['error']}")
            continue

        print(f"\nEXAMPLE {result['id']}:")
        print(f"Question: {result['question']}")
        print(f"Context: {result['context_preview']}")
        print(f"Model answer: {result['model_answer']}")
        print(f"Model decision: {result['model_decision'].upper()}")
        print(f"True decision: {result['true_decision'].upper()}")
        print(f"Correct: {' ' if result['correct'] else ' '}")
        print(f"LLM Judge Scores: {result['llm_judge_scores']}")
        print("-" * 80)

    # Compute metrics if predictions exist
    if predictions and references:
        bert_p, bert_r, bert_f1 = compute_bert_score(predictions, references)
        bertscore_result = {"precision": bert_p, "recall": bert_r, "f1": bert_f1}

        rouge1, rougeL = compute_rouge(predictions, references)
        rouge_result = {"rouge1": rouge1, "rougeL": rougeL}

        report = classification_report(y_true, y_pred, labels=["yes", "no", "maybe"], zero_division=0, output_dict=True)
        macro_f1 = report["macro avg"]["f1-score"]

        print("\nClassification Report:")
        print(classification_report(y_true, y_pred, labels=["yes", "no", "maybe"], zero_division=0))
        print("\nBERTScore:")
        print(f"Precision: {bert_p:.4f}, Recall: {bert_r:.4f}, F1: {bert_f1:.4f}")

        print("\nROUGE:")
        print(f"ROUGE-1: {rouge1:.4f}, ROUGE-L: {rougeL:.4f}")
        print("\nLLM Judge Average Scores:")
        print(f"Correctness: {avg_judge_scores['correctness']:.4f}")
        print(f"Evidence Alignment: {avg_judge_scores['evidence_alignment']:.4f}")

```

```

        print(f"Clarity: {avg_judge_scores['clarity']:.4f}")
    else:
        bertscore_result = rouge_result = {"precision": 0, "recall": 0, "f1": 0}
        macro_f1 = 0
        avg_judge_scores = {"correctness": 0, "evidence_alignment": 0,
        ↪ "clarity": 0}
        print("\nNo valid predictions for metric computation.")

    return {
        "accuracy": accuracy,
        "macro_f1": macro_f1,
        "bertscore": bertscore_result,
        "rouge": rouge_result,
        "llm_judge_scores": avg_judge_scores,
        "results": results,
        "num_examples": len(results)
    }

# Run the evaluation for examples 1-10
params = {
    "temperature": 0.1,
    "top_p": 0.85,           # Focused output
    "max_new_tokens": 300,   # Avoid truncation
    "repetition_penalty": 1.2 # Prevent repetition
}

# Assuming model, tokenizer, and test_dataset are defined
model = model.eval()
for param in model.parameters():
    param.requires_grad = False

results = evaluate_model(model, tokenizer, test_dataset, start_idx=0,
    ↪ num_examples=10)
print("Evaluation Results:", results)

```

```
tokenizer_config.json:  0%|          | 0.00/48.0 [00:00<?, ?B/s]
```

```
vocab.txt:  0%|          | 0.00/232k [00:00<?, ?B/s]
```

```
tokenizer.json:  0%|          | 0.00/466k [00:00<?, ?B/s]
```

```
config.json:  0%|          | 0.00/483 [00:00<?, ?B/s]
```

Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download. For better performance, install the package with: `pip install huggingface_hub[hf_xet]` or `pip install hf_xet`
 WARNING:huggingface_hub.file_download:Xet Storage is enabled for this repo, but the 'hf_xet' package is not installed. Falling back to regular HTTP download.
 For better performance, install the package with: `pip install`

```

huggingface_hub[hf_xet]` or `pip install hf_xet`
model.safetensors: 0%|          | 0.00/268M [00:00<?, ?B/s]

Evaluating with parameters: {'temperature': 0.1, 'top_p': 0.85,
'max_new_tokens': 300, 'repetition_penalty': 1.2}
-----

/usr/local/lib/python3.11/dist-packages/bitsandbytes/autograd/_functions.py:315:
UserWarning: MatMul8bitLt: inputs will be cast from torch.float32 to float16
during quantization
  warnings.warn(f"MatMul8bitLt: inputs will be cast from {A.dtype} to float16
during quantization")

Warning: Decision 'no' may not align with answer 'Malnutrition is a new inducer
for arterial calcification in hemodialysis patients.' (expected: yes).
Warning: Decision 'maybe' may not align with answer 'The risk of serious
infections is higher in patients treated with infliximab than in those treated
with etanercept. The risk is higher in younger patients and in patients of non-
white race/ethnicity.' (expected: yes).
Warning: Decision 'yes' may not align with answer 'Alkalinization of patients on
hemodialysis does not stimulate protein synthesis and does not improve
nutritional parameters.' (expected: no).
Warning: Decision 'maybe' may not align with answer 'Delivery mode does not
affect women's postpartum quality of life in rural China. Factors influencing
postnatal quality of life include maternal education, husband education, infant
gender, home visit and infant sex.' (expected: no).
Warning: Decision 'yes' may not align with answer 'Combination chemotherapy is
not superior to single-agent mitoxantrone in the treatment of high-risk
metastatic breast cancer patients.' (expected: no).

EXAMPLE 1:
Question: Malnutrition, a new inducer for arterial calcification in hemodialysis
patients?
Context: Arterial calcification is a significant cardiovascular risk factor in
hemodialysis patients. A series of factors are involved in the process of
arterial calcification; however, the relationship between malnutrition and
arterial calcification is still...
Model answer: Malnutrition is a new inducer for arterial calcification in
hemodialysis patients.
Model decision: NO
True decision: YES
Correct:
LLM Judge Scores: {'correctness': np.float32(0.96823794), 'evidence_alignment':
np.float32(0.90713847), 'clarity': 1.0}
-----

EXAMPLE 2:
Question: Should temperature be monitorized during kidney allograft
preservation?

```

Context: It is generally considered that kidney grafts should be preserved at 4 degrees C during cold storage. However, actual temperature conditions are not known. We decided to study the temperature levels during preservation with the Biotainer storage can ...

Model answer: Temperature monitoring is not necessary during preservation of human kidneys.

Model decision: NO

True decision: NO

Correct:

LLM Judge Scores: {'correctness': np.float32(0.9238949), 'evidence_alignment': np.float32(0.90346503), 'clarity': 1.0}

EXAMPLE 3:

Question: Screening for gestational diabetes mellitus: are the criteria proposed by the international association of the Diabetes and Pregnancy Study Groups cost-effective?

Context: The International Association of the Diabetes and Pregnancy Study Groups (IADPSG) recently recommended new criteria for diagnosing gestational diabetes mellitus (GDM). This study was undertaken to determine whether adopting the IADPSG criteria would ...

Model answer: The IADPSG criteria are not cost-effective.

Model decision: NO

True decision: YES

Correct:

LLM Judge Scores: {'correctness': np.float32(0.9117916), 'evidence_alignment': np.float32(0.92423785), 'clarity': 1.0}

EXAMPLE 4:

Question: Is resected stomach volume related to weight loss after laparoscopic sleeve gastrectomy?

Context: Laparoscopic sleeve gastrectomy (LSG) was initially performed as the first stage of biliopancreatic diversion with duodenal switch for the treatment of super-obese or high-risk obese patients but is now most commonly performed as a standalone operati...

Model answer: Resected stomach volume is not related to weight loss.

Model decision: NO

True decision: NO

Correct:

LLM Judge Scores: {'correctness': np.float32(0.9195485), 'evidence_alignment': np.float32(0.8938328), 'clarity': 1.0}

EXAMPLE 5:

Question: Body perception: do parents, their children, and their children's physicians perceive body image differently?

Context: To compare children's, parents' and physicians' perceptions of

children's body size. We administered a structured questionnaire of body size perception using a descriptive Likert scale keyed to body image figures to children ages 12 to 18 years. The ...

Model answer: Parents and physicians underestimate children's body size.

Model decision: MAYBE

True decision: YES

Correct:

LLM Judge Scores: {'correctness': np.float32(0.94213635), 'evidence_alignment': np.float32(0.87143415), 'clarity': 1.0}

EXAMPLE 6:

Question: Type II supracondylar humerus fractures: can some be treated nonoperatively?

Context: The range of injury severity that can be seen within the category of type II supracondylar humerus fractures (SCHFs) raises the question whether some could be treated nonoperatively. However, the clinical difficulty in using this approach lies in det...

Model answer: Type II SCHFs can be successfully treated nonoperatively in children with fractures that have a shaft-condylar angle of >15 degrees, a carrying angle of >90 degrees, and no rotational deformity or coronal angulation.

Model decision: NO

True decision: YES

Correct:

LLM Judge Scores: {'correctness': np.float32(0.9594531), 'evidence_alignment': np.float32(0.9391756), 'clarity': 0.8}

EXAMPLE 7:

Question: Comparative safety of infliximab and etanercept on the risk of serious infections: does the association vary by patient characteristics?

Context: Infliximab, a chimeric monoclonal anti-TNF antibody, has been found to increase the risk of serious infections compared with the TNF receptor fusion protein etanercept in some studies. It is unclear whether the risk varies by patient characteristics...

Model answer: The risk of serious infections is higher in patients treated with infliximab than in those treated with etanercept.

Model decision: MAYBE

True decision: YES

Correct:

LLM Judge Scores: {'correctness': np.float32(0.95569867), 'evidence_alignment': np.float32(0.9165395), 'clarity': 1.0}

EXAMPLE 8:

Question: Does increasing blood pH stimulate protein synthesis in dialysis patients?

Context: Although the mechanism of muscle wasting in end-stage renal disease is

not fully understood, there is increasing evidence that acidosis induces muscle protein degradation and could therefore contribute to the loss of muscle protein stores of patients...

Model answer: Alkalinization of patients on hemodialysis does not stimulate protein synthesis and does not improve nutritional parameters.

Model decision: YES

True decision: NO

Correct:

LLM Judge Scores: {'correctness': np.float32(0.92129904), 'evidence_alignment': np.float32(0.93423694), 'clarity': 1.0}

EXAMPLE 9:

Question: Does delivery mode affect women's postpartum quality of life in rural China?

Context: To explore the impact of delivery mode on women's postpartum quality of life in rural China and probe factors influencing postnatal quality of life.

Childbirth significantly affects puerpera's physical, psychological and social domains of quality of ...

Model answer: Delivery mode does not affect women's postpartum quality of life in rural China.

Model decision: MAYBE

True decision: NO

Correct:

LLM Judge Scores: {'correctness': np.float32(0.9632119), 'evidence_alignment': np.float32(0.9398205), 'clarity': 1.0}

EXAMPLE 10:

Question: Is first-line single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer patients as effective as combination chemotherapy?

Context: To determine whether patients with high-risk metastatic breast cancer draw benefit from combination chemotherapy as first-line treatment. A total of 260 women with measurable metastatic breast cancer fulfilling high-risk criteria, previously untreate...

Model answer: Combination chemotherapy is not superior to single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer patients.

Model decision: YES

True decision: YES

Correct:

LLM Judge Scores: {'correctness': np.float32(0.9527812), 'evidence_alignment': np.float32(0.93100774), 'clarity': 1.0}

Some weights of RobertaModel were not initialized from the model checkpoint at roberta-large and are newly initialized: ['pooler.dense.bias', 'pooler.dense.weight']

You should probably TRAIN this model on a down-stream task to be able to use it

for predictions and inference.

calculating scores...

computing bert embedding.

0%| | 0/1 [00:00<?, ?it/s]

computing greedy matching.

0%| | 0/1 [00:00<?, ?it/s]

done in 0.30 seconds, 32.80 sentences/sec

Classification Report:

	precision	recall	f1-score	support
yes	0.50	0.17	0.25	6
no	0.40	0.50	0.44	4
maybe	0.00	0.00	0.00	0
accuracy			0.30	10
macro avg	0.30	0.22	0.23	10
weighted avg	0.46	0.30	0.33	10

BERTScore:

Precision: 0.9198, Recall: 0.8712, F1: 0.8947

ROUGE:

ROUGE-1: 0.3312, ROUGE-L: 0.2661

LLM Judge Average Scores:

Correctness: 0.9418

Evidence Alignment: 0.9161

Clarity: 0.9800

Evaluation Results: {'accuracy': 0.3, 'macro_f1': 0.23148148148148148, 'bertscore': {'precision': 0.9197982549667358, 'recall': 0.8712231516838074, 'f1': 0.8946866989135742}, 'rouge': {'rouge1': np.float64(0.33117113113439495), 'rougeL': np.float64(0.2660820453056485)}, 'llm_judge_scores': {'correctness': np.float32(0.94180524), 'evidence_alignment': np.float32(0.91608876), 'clarity': np.float64(0.9800000000000001)}, 'results': [{'id': 1, 'question': 'Malnutrition, a new inducer for arterial calcification in hemodialysis patients?', 'context_preview': 'Arterial calcification is a significant cardiovascular risk factor in hemodialysis patients. A series of factors are involved in the process of arterial calcification; however, the relationship between malnutrition and arterial calcification is still...', 'model_answer': 'Malnutrition is a new inducer for arterial calcification in hemodialysis patients.', 'model_decision': 'no', 'true_decision': 'yes', 'correct': False, 'llm_judge_scores': {'correctness': np.float32(0.96823794), 'evidence_alignment': np.float32(0.90713847), 'clarity': 1.0}}, {'id': 2,

'question': 'Should temperature be monitorized during kidney allograft preservation?', 'context_preview': 'It is generally considered that kidney grafts should be preserved at 4 degrees C during cold storage. However, actual temperature conditions are not known. We decided to study the temperature levels during preservation with the Biotainer storage can ...', 'model_answer': 'Temperature monitoring is not necessary during preservation of human kidneys.', 'model_decision': 'no', 'true_decision': 'no', 'correct': True, 'llm_judge_scores': {'correctness': np.float32(0.9238949), 'evidence_alignment': np.float32(0.90346503), 'clarity': 1.0}}, {'id': 3, 'question': 'Screening for gestational diabetes mellitus: are the criteria proposed by the international association of the Diabetes and Pregnancy Study Groups cost-effective?', 'context_preview': 'The International Association of the Diabetes and Pregnancy Study Groups (IADPSG) recently recommended new criteria for diagnosing gestational diabetes mellitus (GDM). This study was undertaken to determine whether adopting the IADPSG criteria would ...', 'model_answer': 'The IADPSG criteria are not cost-effective.', 'model_decision': 'no', 'true_decision': 'yes', 'correct': False, 'llm_judge_scores': {'correctness': np.float32(0.9117916), 'evidence_alignment': np.float32(0.92423785), 'clarity': 1.0}}, {'id': 4, 'question': 'Is resected stomach volume related to weight loss after laparoscopic sleeve gastrectomy?', 'context_preview': 'Laparoscopic sleeve gastrectomy (LSG) was initially performed as the first stage of biliopancreatic diversion with duodenal switch for the treatment of super-obese or high-risk obese patients but is now most commonly performed as a standalone operati...', 'model_answer': 'Resected stomach volume is not related to weight loss.', 'model_decision': 'no', 'true_decision': 'no', 'correct': True, 'llm_judge_scores': {'correctness': np.float32(0.9195485), 'evidence_alignment': np.float32(0.8938328), 'clarity': 1.0}}, {'id': 5, 'question': "Body perception: do parents, their children, and their children's physicians perceive body image differently?", 'context_preview': "To compare children's, parents' and physicians' perceptions of children's body size. We administered a structured questionnaire of body size perception using a descriptive Likert scale keyed to body image figures to children ages 12 to 18 years. The ...", 'model_answer': "Parents and physicians underestimate children's body size.", 'model_decision': 'maybe', 'true_decision': 'yes', 'correct': False, 'llm_judge_scores': {'correctness': np.float32(0.94213635), 'evidence_alignment': np.float32(0.87143415), 'clarity': 1.0}}, {'id': 6, 'question': 'Type II supracondylar humerus fractures: can some be treated nonoperatively?', 'context_preview': 'The range of injury severity that can be seen within the category of type II supracondylar humerus fractures (SCHFs) raises the question whether some could be treated nonoperatively. However, the clinical difficulty in using this approach lies in det...', 'model_answer': 'Type II SCHFs can be successfully treated nonoperatively in children with fractures that have a shaft-condylar angle of >15 degrees, a carrying angle of >90 degrees, and no rotational deformity or coronal angulation.', 'model_decision': 'no', 'true_decision': 'yes', 'correct': False, 'llm_judge_scores': {'correctness': np.float32(0.9594531), 'evidence_alignment': np.float32(0.9391756), 'clarity': 0.8}}, {'id': 7, 'question': 'Comparative safety of infliximab and etanercept on the risk of serious infections: does the association vary by patient

characteristics?', 'context_preview': 'Infliximab, a chimeric monoclonal anti-TNF antibody, has been found to increase the risk of serious infections compared with the TNF receptor fusion protein etanercept in some studies. It is unclear whether the risk varies by patient characteristics...', 'model_answer': 'The risk of serious infections is higher in patients treated with infliximab than in those treated with etanercept.', 'model_decision': 'maybe', 'true_decision': 'yes', 'correct': False, 'llm_judge_scores': {'correctness': np.float32(0.95569867), 'evidence_alignment': np.float32(0.9165395), 'clarity': 1.0}}, {'id': 8, 'question': 'Does increasing blood pH stimulate protein synthesis in dialysis patients?', 'context_preview': 'Although the mechanism of muscle wasting in end-stage renal disease is not fully understood, there is increasing evidence that acidosis induces muscle protein degradation and could therefore contribute to the loss of muscle protein stores of patients...', 'model_answer': 'Alkalinization of patients on hemodialysis does not stimulate protein synthesis and does not improve nutritional parameters.', 'model_decision': 'yes', 'true_decision': 'no', 'correct': False, 'llm_judge_scores': {'correctness': np.float32(0.92129904), 'evidence_alignment': np.float32(0.93423694), 'clarity': 1.0}}, {'id': 9, 'question': '"Does delivery mode affect women's postpartum quality of life in rural China?"', 'context_preview': '"To explore the impact of delivery mode on women's postpartum quality of life in rural China and probe factors influencing postnatal quality of life. Childbirth significantly affects puerpera's physical, psychological and social domains of quality of ..."', 'model_answer': '"Delivery mode does not affect women's postpartum quality of life in rural China."', 'model_decision': 'maybe', 'true_decision': 'no', 'correct': False, 'llm_judge_scores': {'correctness': np.float32(0.9632119), 'evidence_alignment': np.float32(0.9398205), 'clarity': 1.0}}, {'id': 10, 'question': 'Is first-line single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer patients as effective as combination chemotherapy?', 'context_preview': 'To determine whether patients with high-risk metastatic breast cancer draw benefit from combination chemotherapy as first-line treatment. A total of 260 women with measurable metastatic breast cancer fulfilling high-risk criteria, previously untreat...', 'model_answer': 'Combination chemotherapy is not superior to single-agent mitoxantrone in the treatment of high-risk metastatic breast cancer patients.', 'model_decision': 'yes', 'true_decision': 'yes', 'correct': True, 'llm_judge_scores': {'correctness': np.float32(0.9527812), 'evidence_alignment': np.float32(0.93100774), 'clarity': 1.0}}], 'num_examples': 10}