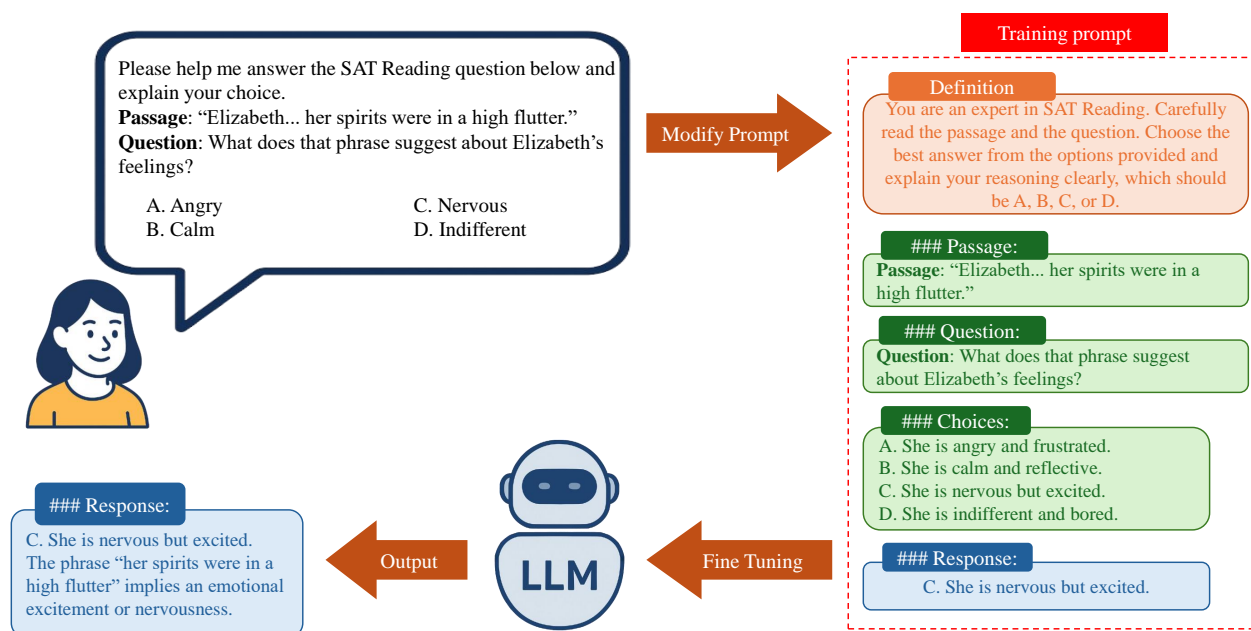


# Tutorial: Boosting LLMs' SAT Reading Proficiency via Instruction Tuning

Dinh-Thang Duong, Yen-Linh Vu, Anh-Khoi Nguyen, Quang-Vinh Dinh

## I. Giới thiệu

**Instruction Tuning (IT)** là một trong những kỹ thuật training mô hình ngôn ngữ lớn (LLMs) rất quan trọng. Trong đó, IT giúp cải thiện khả năng của mô hình cũng như kiểm soát kết quả đầu ra. Là kiểu huấn luyện mô hình có giám sát từ bộ dữ liệu theo cặp (instruction-output), từ đó giúp mô hình thu hẹp khoảng cách giữa từ kế tiếp được sinh ra và sự chỉ dẫn của con người.

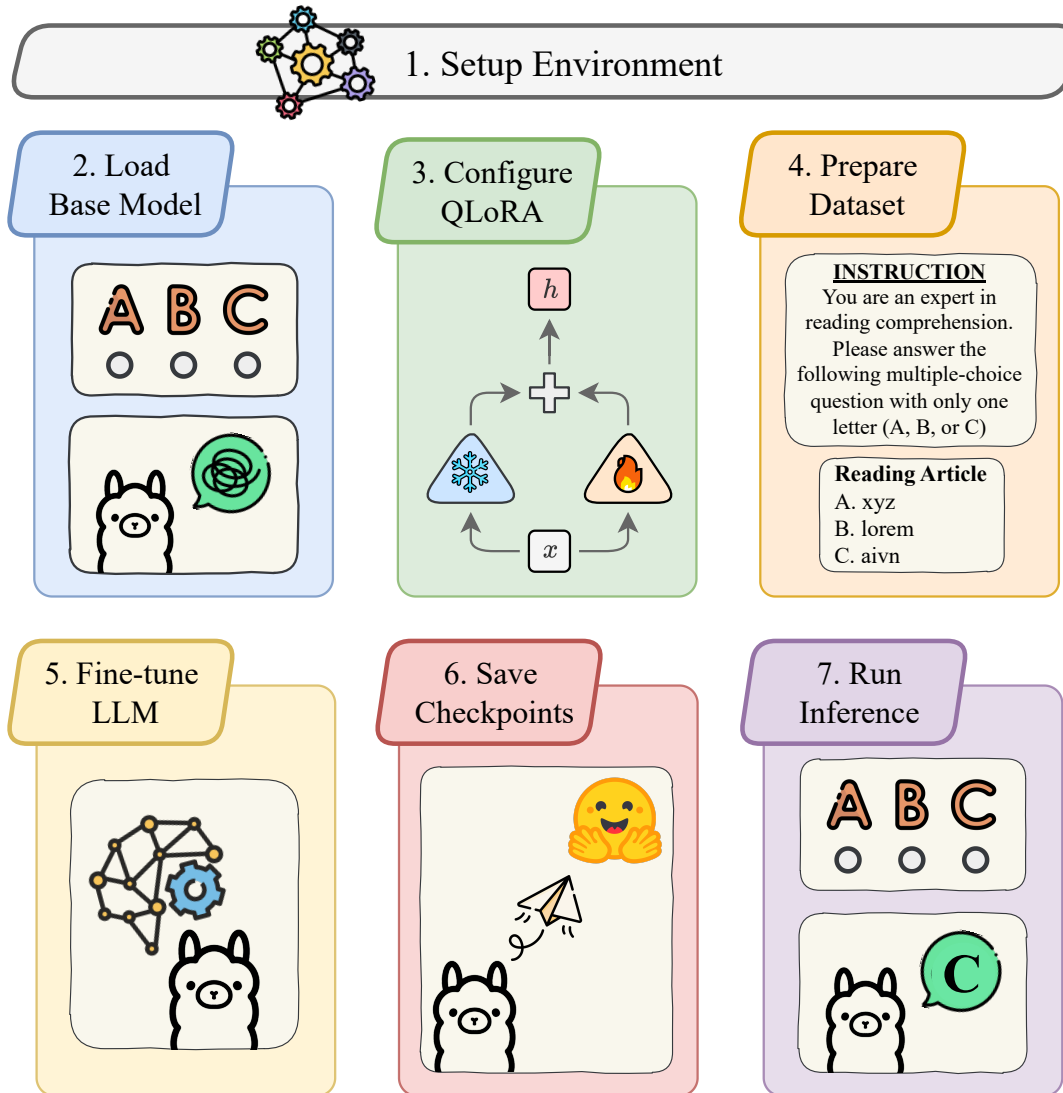


Hình 1: Instruction Fine-tuning mô hình ngôn ngữ lớn với dữ liệu trắc nghiệm.

Bài viết trình bày quá trình huấn luyện một mô hình ngôn ngữ lớn (LLM) với dữ liệu instruction để thực hiện các câu hỏi trắc nghiệm trong bài đọc hiểu thuộc kỳ thi SAT. Mục tiêu là xây dựng một mô hình có khả năng cải thiện độ chính xác trong các bài kiểm tra đọc hiểu SAT, đồng thời hỗ trợ nâng cao kỹ năng đọc hiểu tiếng Anh một cách tổng quát. Bài toán được định nghĩa với đầu vào và đầu ra như sau:

- **Input:** Một câu prompt với lời hướng dẫn (instruction) để LLM thực hiện một bài bài đọc hiểu SAT và lựa chọn phương án đúng.
- **Output:** Lời phản hồi từ mô hình, trong trường hợp này là một trong các phương án trắc nghiệm đã được mô tả trong prompt.

Thông qua nội dung mô tả trên, pipeline các bước xử lý trong bài này được minh họa theo hình dưới đây:



Hình 2: Fine-tuning mô hình ngôn ngữ lớn với dữ liệu instruction giải toán trắc nghiệm.

Trong đó:

### 1. Setup Environment:

Đầu tiên là cài đặt một số thư viện sau để có thể chạy được các mô hình ngôn ngữ lớn từ thư viện HuggingFace:

```
1 !pip install -q -U bitsandbytes
2 !pip install -q -U datasets
3 !pip install -q -U git+https://github.com/huggingface/transformers.git
4 !pip install -q -U git+https://github.com/huggingface/peft.git
5 !pip install -q -U git+https://github.com/huggingface/accelerate.git
6 !pip install -q -U loralib
7 !pip install -q -U einops
```

Sau khi cài đặt hoàn tất, bước tiếp theo là import các thư viện đã tải cũng như một số thư viện khác để phục vụ cho chương trình:

```
1 import json
2 import os
3 import bitsandbytes as bnb
4 import torch
5 import torch.nn as nn
6 import transformers
7
8 from pprint import pprint
9 from tqdm import tqdm
10 from datasets import load_dataset, Dataset
11
12 from peft import (
13     LoraConfig,
14     PeftConfig,
15     PeftModel,
16     get_peft_model,
17     prepare_model_for_kbit_training
18 )
19 from transformers import (
20     AutoConfig,
21     AutoModelForCausalLM,
22     AutoTokenizer,
23     BitsAndBytesConfig
24 )
```

### 2. Load Base Model:

Mô hình ngôn ngữ lớn được sử dụng trong bài này là Llama đến từ Meta (Facebook), mô hình là phiên bản 3.1-Instruct với tổng cộng 8 tỷ tham số. Tuy nhiên, mô hình này cần xin quyền truy cập từ Meta. Để làm điều đó, trước tiên hãy vào [đường link này](#), hãy đăng nhập và xin cấp quyền sau khi điền đủ các thông tin được yêu cầu. Thông thường, Meta sẽ duyệt trong vài phút kể từ lúc gửi. Sau khi được cấp quyền, hãy lấy **User Access Tokens** bằng cách nhấn vào [link này](#) (đã đăng nhập trước), chọn “Create New Token”, chọn “Read”, đặt tên và nhấn “Create token”, và thay Access Tokens này vào đoạn code sau để đăng nhập HuggingFace trên notebook:

```
1 from huggingface_hub import login
2
3 login(token="hf_XXXXXXXXXXXXXXXXX") # your access token
```

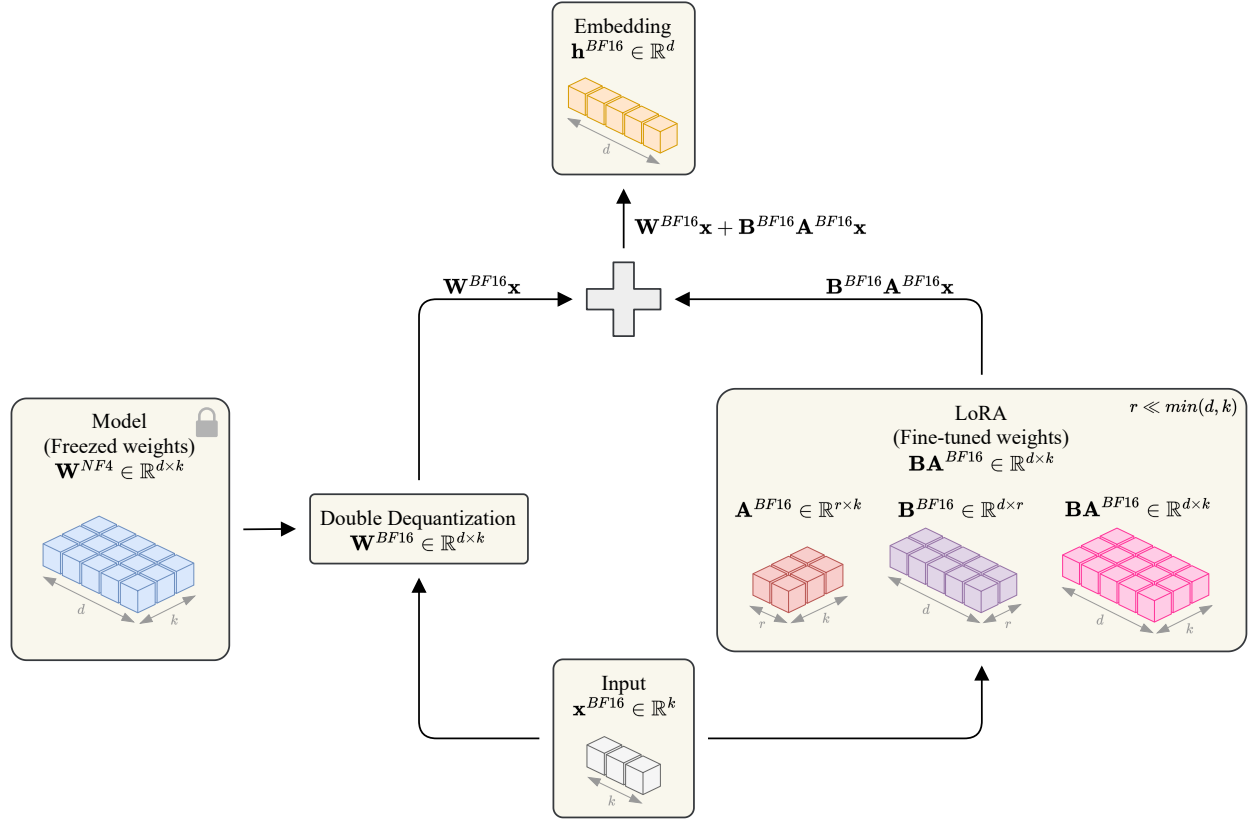
Khi đã đăng nhập thành công, tiến hành khởi tạo mô hình Llama như sau:

```
1 MODEL_NAME = "meta-llama/Llama-3.1-8B-Instruct"
2 # MODEL_NAME = "meta-llama/Llama-3.2-3B-Instruct"
3
4 bnb_config = BitsAndBytesConfig(
5     load_in_4bit=True,
6     bnb_4bit_quant_type="nf4",
7     bnb_4bit_compute_dtype=torch.bfloat16,
8     bnb_4bit_use_double_quant=True
9 )
10
11 model = AutoModelForCausalLM.from_pretrained(
12     MODEL_NAME,
13     quantization_config=bnb_config,
14     device_map="auto",
15     trust_remote_code=True
16 )
17
18 tokenizer = AutoTokenizer.from_pretrained(MODEL_NAME)
19 tokenizer.pad_token = tokenizer.eos_token
20
21 model.gradient_checkpointing_enable()
```

Trong đoạn code trên, mô hình được tải với kỹ thuật lượng tử hóa 4-bit (NF4) thông qua thư viện BitsAndBytes, giúp giảm dung lượng và tối ưu hóa hiệu suất tính toán trên GPU. Tokenizer phù hợp với mô hình cũng được tải về và thiết lập token padding bằng ký tự kết thúc chuỗi (EOS). Tiếp theo, gradient checkpointing được kích hoạt để tiết kiệm bộ nhớ khi huấn luyện. Sau đó, hàm `prepare_model_for_kbit_training()` được áp dụng nhằm điều chỉnh mô hình sẵn sàng cho quá trình fine-tuning.

### 3. Configure QLoRA:

Kể đến, kỹ thuật LoRA được áp dụng lên các module attention (`q_proj` và `v_proj`) để fine-tune mô hình nhanh chóng, nhẹ nhàng, nhưng vẫn hiệu quả.



Hình 3: Trực quan về kỹ thuật QLoRA.

```

1 model = prepare_model_for_kbit_training(model)
2
3 peft_config = LoraConfig(
4     r=8,
5     lora_alpha=16,
6     target_modules=[
7         "q_proj",
8         "v_proj",
9     ],
10    lora_dropout=0.05,
11    bias="none",
12    task_type="CAUSAL_LM"
13 )
14
15 model = get_peft_model(model, peft_config)

```

#### 4. Prepare Dataset:

Để thực hiện nhiệm vụ “đọc bài đọc và trả lời câu hỏi trắc nghiệm”, mô hình LLaMA được fine-tune trên tập dữ liệu `mozilla/sat-reading` từ Hugging Face. Bộ dữ liệu này bao gồm các đoạn văn và câu hỏi trích từ phần đọc hiểu của mười bài thi thử SAT được công bố công khai. Mỗi mục dữ liệu chứa đoạn văn, câu hỏi, bốn lựa chọn trả lời, và đáp án đúng. Một số câu hỏi yêu cầu tham chiếu đến dòng cụ thể trong văn bản, được đánh dấu bằng trường boolean `requires_line`. Các mục liên quan đến biểu đồ và bảng biểu đã bị

loại bỏ nhằm đảm bảo tính nhất quán của tập dữ liệu. Dữ liệu được chia thành ba phần: tập huấn luyện (298 mục), tập kiểm định (39 mục), và tập kiểm tra (38 mục), tổng cộng 375 mục. Việc tải và đọc dữ liệu được thực hiện như sau:

```
1 data = load_dataset("emozilla/sat-reading")
```

Dưới đây là một vài ví dụ điển hình trong tập dữ liệu:

id	text (shortened)	answer	requires line
sat-practice_7-question_3	SAT READING COMPREHENSION TEST This passage is adapted from ... Question 3: Which statement best describes a technique the narrator uses to represent Silas's character before he adopted Eppie? A) ... B) ... C) ... D) ... Answer:	A	false
sat-practice_7-question_10	SAT READING COMPREHENSION TEST This passage is adapted from ... Question 10: As used in line 65, "fine" most nearly means A) acceptable B) delicate C) ornate D) keen Answer:	D	true

Tuy nhiên, để đảm bảo cấu trúc rõ ràng cho các mẫu huấn luyện, hai hàm được định nghĩa nhằm trích xuất thông tin từ văn bản gốc, gồm:

- `extract_sections()` thực hiện phân tách nội dung đầu vào thành các phần rõ ràng như đoạn văn, câu hỏi, các lựa chọn, và đáp án dưới dạng ký tự (A, B, C, D).
- `map_answer()` được dùng để ánh xạ ký tự đáp án sang nội dung đầy đủ tương ứng trong danh sách lựa chọn.

```
1 def extract_sections(text):
2     sections = {
3         "passage": "",
4         "question": "",
5         "choices": [],
6         "answer_letter": ""
7     }
8
9     answer_part = text.split("Answer:")[1].strip()
10    sections["answer_letter"] = answer_part[0] if answer_part else ""
11
12    content = text.split("SAT READING COMPREHENSION TEST")[1].split("Answer:")[0]
```

```

13     blocks = [b.strip() for b in content.split("\n\n") if b.strip()]
14
15     passage_lines = []
16     for line in blocks:
17         if line.startswith("Question"):
18             break
19         passage_lines.append(line)
20     sections["passage"] = "\n".join(passage_lines).strip()
21
22     for block in blocks:
23         if block.startswith("Question"):
24             q_part = block.split(" ", 1) if " " in block else (block, "")
25             sections["question"] = q_part[-1].split("\n")[0].strip()
26             sections["choices"] = [line.strip() for line in block.split("\n")[1:]
27                                   if line.startswith(("A", "B", "C", "D"))]
28     return sections
29
30
31 def map_answer(text, letter):
32     sections = extract_sections(text)
33     for choice in sections["choices"]:
34         if choice.startswith(f"{letter}"):
35             return choice
36     return letter

```

Sau đó, dùng hai hàm vừa định nghĩa để xây dựng một hàm tạo prompt với cấu trúc như mong muốn bằng cách sau:

```

1 LLAMA3_SYSTEM_PROMPT = """You are a helpful AI assistant developed by Meta. Respond
2                             safely and accurately."""
3
4 def generate_prompt(text, answer_letter):
5     sections = extract_sections(text)
6     choices_text = "\n".join(sections['choices'])
7
8     return [
9         {
10             "role": "system",
11             "content": LLAMA3_SYSTEM_PROMPT
12         },
13         {
14             "role": "user",
15             "content": f"""Read the passage and answer the question.
16
17 ### Passage:
18 {sections["passage"]}
19
20 ### Question:
21 {sections["question"]}
22
23 ### Choices:
24 {choices_text}
25
26 Respond with ONLY the letter and full text of the correct answer."""

```

```

27     {
28         "role": "assistant",
29         "content": map_answer(text, answer_letter)
30     }
31 ]

```

Tiếp tục thực hiện hàm tokenize để chuẩn bị đầu vào phù hợp cho mô hình. Cụ thể là:

- Hàm `generate_and_tokenize_prompt()` nhận đầu vào là một cặp (`user_input`, `answer`) và sinh ra prompt hoàn chỉnh thông qua hàm `generate_prompt()`.
- Prompt sau đó được xử lý định dạng bằng `tokenizer.apply_chat_template()` để đảm bảo phù hợp với định dạng hội thoại mà tokenizer yêu cầu.
- Hàm `tokenizer()` được dùng để chuyển prompt thành tensor `input_ids`, cùng với `attention_mask` phục vụ cho quá trình huấn luyện.
- Cuối cùng, nhãn huấn luyện (labels) được gán bằng chính `input_ids` để mô hình học theo dạng tự hồi tiếp (causal language modeling).

```

1 def generate_and_tokenize_prompt(user_input, answer):
2     try:
3         full_prompt = generate_prompt(user_input, answer)
4
5         prompt_str = tokenizer.apply_chat_template(
6             full_prompt,
7             tokenize=False,
8             add_generation_prompt=False
9         )
10
11        tokenized = tokenizer(
12            prompt_str,
13            padding="max_length",
14            truncation=True,
15            max_length=1024,
16            return_tensors="pt"
17        )
18
19        input_ids = tokenized["input_ids"][0]
20        labels = input_ids.clone()
21
22        return {
23            "input_ids": input_ids,
24            "attention_mask": tokenized["attention_mask"][0],
25            "labels": labels
26        }
27
28    except Exception as e:
29        print(f"Error processing sample: {e}")
30        return None

```

Với các hàm đã được chuẩn bị, bước cuối cùng trong phần này là xây dựng bộ dữ liệu huấn luyện theo các bước sau:

- Duyệt qua từng mẫu dữ liệu trong tập huấn luyện gốc.



- Thực hiện tiền xử lý văn bản bằng cách loại bỏ tiêu đề thừa và chuẩn hóa đáp án bằng hàm `map_answer()`.
- Áp dụng hàm `generate_and_tokenize_prompt()` để sinh ra đầu vào tokenized cho từng mẫu.
- Cuối cùng, tập dữ liệu được chia thành hai phần: tập huấn luyện (`train_dataset`) và tập đánh giá (`eval_dataset`) theo tỷ lệ 90/10.

```
1 training_samples = []
2 for sample in tqdm(data["train"]):
3     try:
4         processed_text = sample["text"].replace("SAT READING COMPREHENSION TEST", "").strip()
5         processed_answer = map_answer(sample["text"], sample["answer"].strip())
6
7         tokenized_sample = generate_and_tokenize_prompt(processed_text, processed_answer)
8
9         if tokenized_sample is not None:
10             training_samples.append(tokenized_sample)
11     except Exception as e:
12         print(f"Skipping invalid sample: {e}")
13
14 training_samples = [s for s in training_samples if s is not None]
15
16 train_samples, val_samples = train_test_split(training_samples, test_size=0.1,
17                                               random_state=42)
18
19 train_dataset = Dataset.from_list(train_samples)
20 eval_dataset = Dataset.from_list(val_samples)
```

```

<|begin_of_text|><|start_header_id|>system<|end_header_id|>

Cutting Knowledge Date: December 2023
Today Date: 26 Jul 2024

You are a helpful AI assistant developed by Meta. Respond safely and accurately.
<|eot_id|><|start_header_id|>user<|end_header_id|>

Read the passage and answer the question.

### Passage:
This passage is adapted from Elizabeth Cady Stanton's address to the 1869 Woman
Suffrage Convention in Washington, DC. I urge a sixteenth amendment, because

. . .

and not with her consent would one drop of blood ever be shed, one life sacrificed in vain.

### Question:
Stanton claims that which of the following was a relatively recent historical development?

### Choices:
A) The control of society by men
B) The spread of war and injustice
C) The domination of domestic life by men
D) The acknowledgment of women's true character

Respond with ONLY the letter and full text of the correct answer.<|eot_id|>
<|start_header_id|>assistant<|end_header_id|>

D) The acknowledgment of women's true character<|eot_id|>

```

Hình 4: Trục quan một prompt đã được tổ chức lại.

## 5. Fine-tune LLM:

Sau khi dữ liệu huấn luyện được chuẩn bị, quá trình huấn luyện mô hình được thực hiện theo các bước sau:

- Lớp **LogLossCallback** được định nghĩa nhằm ghi lại giá trị hàm mất mát (loss) sau mỗi bước logging, hỗ trợ theo dõi tiến trình huấn luyện theo thời gian thực.
- Các tham số huấn luyện được thiết lập thông qua **TrainingArguments**, bao gồm:
  - Batch size nhỏ (1), kết hợp với **gradient accumulation** để phù hợp với mô hình lớn.
  - Sử dụng optimizer 8-bit (**paged\_adamw\_8bit**) nhằm tiết kiệm bộ nhớ và tăng hiệu suất.
  - Scheduler dạng cosine với tỉ lệ warmup ban đầu.
  - Đánh giá và lưu mô hình mỗi 50 bước thông qua **eval\_steps** và **save\_steps**.
  - Cấu hình tự động tải lại checkpoint có hiệu suất tốt nhất theo metric loss.

- Hàm `DataCollatorForLanguageModeling()` được sử dụng để xử lý đầu vào theo chuẩn causal language modeling (không áp dụng masking).
- Trước khi huấn luyện, gradient cho đầu vào được kích hoạt và mô hình được biên dịch bằng `torch.compile()` để tối ưu hiệu suất trên phần cứng hiện đại.
- Quá trình huấn luyện được thực hiện thông qua lệnh gọi `Trainer.train()`.

```

1 class LogLossCallback(TrainerCallback):
2     def on_log(self, args, state, control, logs=None, **kwargs):
3         if logs is not None and "loss" in logs:
4             print(f"Step {state.global_step} - Loss: {logs['loss']:.4f}")
5
6 training_args = TrainingArguments(
7     per_device_train_batch_size=1,
8     gradient_accumulation_steps=2,
9     num_train_epochs=2,
10    learning_rate=2e-4,
11    fp16=True,
12    save_total_limit=3,
13    logging_steps=10,
14    output_dir="llama3-8b-sat-reading",
15    optim="paged_adamw_8bit",
16    lr_scheduler_type="cosine",
17    warmup_ratio=0.05,
18    eval_strategy="steps",
19    eval_steps=50,
20    save_strategy="steps",
21    save_steps=50,
22    load_best_model_at_end=True,
23    metric_for_best_model="loss",
24    greater_is_better=False,
25    report_to="none",
26    remove_unused_columns=False
27 )
28
29 data_collator = DataCollatorForLanguageModeling(
30     tokenizer=tokenizer,
31     mlm=False,
32     pad_to_multiple_of=8
33 )
34
35 trainer = Trainer(
36     model=model,
37     train_dataset=train_dataset,
38     eval_dataset=eval_dataset,
39     args=training_args,
40     data_collator=data_collator,
41     callbacks=[LogLossCallback()]
42 )
43
44 model.config.use_cache = False
45 model.enable_input_require_grads()
46 model = torch.compile(model)
47
48 trainer.train()

```

## 6. Save Checkpoints:

Khi quá trình huấn luyện hoàn tất, mô hình được lưu để sử dụng hoặc chia sẻ sau này. Cụ thể:

- Mô hình sau huấn luyện được lưu cục bộ vào thư mục **"trained-model"** bằng hàm `save_pretrained()`.
- Sau đó, mô hình được đẩy lên Hugging Face Hub thông qua hàm `push_to_hub()`, giúp dễ dàng chia sẻ và tái sử dụng trong các dự án khác.
- Để thực hiện bước này, người dùng cần có tài khoản trên Hugging Face và phải cung cấp `use_auth_token=True` để xác thực quyền truy cập.

```
1 model.save_pretrained("trained-model")
2 PEFT_MODEL = "your_huggingface_user_name/instructionTuning-llama-3-1-8B-SAT-reading-
                    solver"
3 model.push_to_hub(PEFT_MODEL, use_auth_token=True)
```

## 7. Run Inference:

Cuối cùng, để đánh giá mô hình sau huấn luyện, quá trình tải mô hình và thực hiện suy luận (inference) trên tập kiểm tra được tiến hành như sau:

- Mô hình được tải từ Hugging Face Hub thông qua ID đã định nghĩa trước trong biến `PEFT_MODEL`
- Tiếp theo là định nghĩa lại các thành phần liên quan đến mô hình tương tự như các đoạn code đã thực hiện bên trên.
- Sau khi tải mô hình và tokenizer, mô hình được gắn adapter PEFT thông qua hàm `PeftModel.from_pretrained()` để khôi phục các tham số đã fine-tuned.

```
1 PEFT_MODEL = "your_huggingface_user_name/instructionTuning-llama-3-1-8B-SAT-reading-
                    solver"
2
3 config = PeftConfig.from_pretrained(PEFT_MODEL)
4
5 bnb_config = BitsAndBytesConfig(
6     load_in_4bit=True,
7     bnb_4bit_use_double_quant=True,
8     bnb_4bit_quant_type="nf4",
9     bnb_4bit_compute_dtype=torch.bfloat16
10 )
11
12 model = AutoModelForCausalLM.from_pretrained(
13     config.base_model_name_or_path,
14     quantization_config=bnb_config,
15     device_map="auto",
16     trust_remote_code=True
17 )
18
19 tokenizer = AutoTokenizer.from_pretrained(config.base_model_name_or_path)
20 tokenizer.pad_token = tokenizer.eos_token
21
22 model = PeftModel.from_pretrained(model, PEFT_MODEL)
```

Khi đã khôi phục mô hình, quá trình định dạng prompt và sinh câu trả lời được thực hiện để đánh giá khả năng suy luận của mô hình trên các mẫu kiểm tra, chi tiết như sau:

- Trong hàm `predict()`, prompt đầu vào được định dạng lại theo chuẩn hội thoại bằng `tokenizer.apply_chat_template`, sau đó thực hiện sinh câu trả lời bằng hàm `model.generate()` với các tham số cấu hình sinh (ví dụ: không sampling, temperature = 0.0, giới hạn số token đầu ra).
- Hàm `extract_answer()` được dùng để trích xuất câu trả lời đầy đủ từ đầu ra của mô hình, và được so sánh với đáp án gốc sau khi ánh xạ bằng `map_answer()`.
- Vòng lặp đánh giá thực hiện suy luận trên 10 mẫu đầu tiên của tập kiểm tra, in ra nội dung câu hỏi, các lựa chọn, dự đoán từ mô hình, đáp án đúng và kết quả đúng/sai.

```

1 generation_config = GenerationConfig(
2     max_new_tokens=64,
3     temperature=0.0,
4     top_p=1.0,
5     do_sample=False,
6     repetition_penalty=1.0,
7     eos_token_id=tokenizer.eos_token_id,
8     pad_token_id=tokenizer.eos_token_id
9 )
10
11 def predict(text):
12     messages = format_test_prompt(text)
13
14     prompt_text = tokenizer.apply_chat_template(
15         messages,
16         add_generation_prompt=True,
17         tokenize=False
18     )
19
20     inputs = tokenizer(prompt_text, return_tensors="pt").to(model.device)
21
22     with torch.no_grad():
23         outputs = model.generate(
24             input_ids=inputs["input_ids"],
25             attention_mask=inputs["attention_mask"],
26             generation_config=generation_config
27         )
28
29     output_text = tokenizer.decode(outputs[0], skip_special_tokens=True)
30     return extract_answer(output_text)
31
32 for i in range(10):
33     print("="*100)
34     sample = data["test"][i]
35     input_text = sample["text"]
36     true_answer = sample["answer"].strip()
37
38     predicted_answer = predict(input_text)
39
40     true_answer_full = map_answer(input_text, true_answer)
41

```

```
42     pred_choice = extract_choice_letter(predicted_answer)
43     true_choice = extract_choice_letter(true_answer_full)
44
45     print(f"### Sample {i+1}")
46     print(f"[Question]\n{extract_sections(input_text)[\"question\"]}")
47     print(f"[Choices]\n{extract_sections(input_text)[\"choices\"]}")
48     print(f"\n[Model Prediction]\n{predicted_answer}")
49     print(f"\n[Ground Truth]\n{true_answer_full}")
50     print(f"""\nResult: {"CORRECT" if pred_choice == true_choice else "INCORRECT"}""")
51     print(f"="*100 + "\n")
```

## II. Câu hỏi trắc nghiệm

1. Trong ngữ cảnh về mô hình ngôn ngữ lớn (LLMs), Instruction Tuning được hiểu như thế nào?
  - (a) Huấn luyện mô hình trên task mới mà không cần mẫu dữ liệu nào.
  - (b) Điều chỉnh kết quả của mô hình trong quá trình deploy.
  - (c) Huấn luyện mô hình để tuân theo các yêu cầu cụ thể (instruction).
  - (d) Giảm kích thước của mô hình để tăng độ hiệu quả.
2. Instruction Tuning là một dạng của kiểu học gì?
  - (a) Supervised Learning.
  - (b) Self-supervised Learning.
  - (c) Unsupervised Learning.
  - (d) Reinforcement Learning.
3. Khi thực hiện Instruction Tuning, hàm loss nào sau đây có thể được sử dụng?
  - (a) Mean Squared Error (MSE).
  - (b) Hinge Loss.
  - (c) Kullback-Leibler Divergence.
  - (d) Cross-Entropy Loss.
4. Mệnh đề sau đúng hay sai: “Ta luôn luôn nên áp dụng Instruction Tuning để giảm thiểu chi phí tính toán trong quá trình training”?
  - (a) Đúng.
  - (b) Sai.
5. Trong LLMs, khái niệm prompt được hiểu như thế nào?
  - (a) Một phần mềm hỗ trợ giúp tối ưu hóa hiệu suất của mô hình.
  - (b) Một kỹ thuật mã hóa thông tin riêng tư trong quá trình đào tạo mô hình.
  - (c) Một câu hỏi hoặc yêu cầu mà người dùng đưa ra cho mô hình.
  - (d) Một thuật toán đặc biệt để phân loại dữ liệu đầu vào.
6. Trong LLMs, ta nên áp dụng kỹ thuật nào sau đây để cải thiện khả năng thực hiện một task cụ thể nào đó của mô hình mà không cần training?
  - (a) Sử dụng kỹ thuật Transfer Learning.
  - (b) Ứng dụng khả năng Zero-shot Learning của mô hình.
  - (c) Lập trình thủ công tại bước hậu xử lý cho mỗi task.
  - (d) Mở rộng bộ dữ liệu training.

7. Kỹ thuật prompting nào dưới đây cung cấp cho mô hình chỉ một ví dụ về task cần làm?
- (a) One-shot Learning.
  - (b) Few-shot Learning.
  - (c) Continuous Learning.
  - (d) Transfer Learning.
8. Câu nào sau đây mô tả đúng về kỹ thuật Parameter Efficient Fine-tuning (PEFT)?
- (a) Một kỹ thuật dùng để huấn luyện mô hình trên bộ dữ liệu cực lớn.
  - (b) Một kỹ thuật liên quan đến việc cập nhật một phần nhỏ tham số của mô hình khi huấn luyện.
  - (c) Một kỹ thuật huấn luyện dành riêng cho các mô hình có kích thước nhỏ (dưới 1 tỷ tham số).
  - (d) Một kỹ thuật để tăng chi phí tính toán của mô hình.
9. Mệnh đề sau đúng hay sai: “Low-Rank Adaptation (LoRA) là một kỹ thuật về PEFT.”?
- (a) Đúng.
  - (b) Sai.
10. So với LoRA, QLoRA có điểm gì khác biệt gì trong việc huấn luyện LLMs?
- (a) QLoRA lượng tử hóa (quantize) tham số mô hình; LoRA thì không.
  - (b) QLoRA sử dụng nhiều tham số mô hình; LoRA ít hơn.
  - (c) QLoRA tối ưu khả năng tổng quát của mô hình; LoRA tối ưu trên một task cụ thể.
  - (d) QLoRA giảm kích thước mô hình; LoRA tăng lên.



### III. Phụ lục

1. **Datasets:** Các file dataset được đề cập trong bài có thể được tải tại [đây](#).
2. **Hint:** Các file code gợi ý có thể được tải tại [đây](#).
3. **Solution:** Các file code cài đặt hoàn chỉnh và phần trả lời nội dung trắc nghiệm có thể được tải tại [đây](#) (**Lưu ý:** Sáng thứ 6 khi hết deadline phần nội dung này, ad mới copy các tài liệu bài giải nêu trên vào đường dẫn).
4. **Demo:** Web demo và mã nguồn của ứng dụng có thể được truy cập tại [đây](#).
5. **Rubric:**

Mục	Kiến Thức	Đánh Giá
I.	<ul style="list-style-type: none"> <li>- Kiến thức về mô hình lớn (LLMs).</li> <li>- Kiến thức về kỹ thuật Instruction Tuning trong LLMs.</li> <li>- Kiến thức về thư viện PyTorch và HuggingFace.</li> <li>- Triển khai PEFT (LoRA) và thiết lập các tham số liên quan đến LoRA.</li> </ul>	<ul style="list-style-type: none"> <li>- Nắm được lý thuyết về mô hình lớn (LLMs) và kỹ thuật Instruction Tuning trong LLMs.</li> <li>- Nắm được lý thuyết cơ bản về kỹ thuật LoRA.</li> <li>- Có khả năng cài đặt và thực hiện Instruction Tuning mô hình lớn trên bộ dữ liệu bất kì, trong trường hợp này về bộ dữ liệu IELTS Reading.</li> </ul>

- Hết -