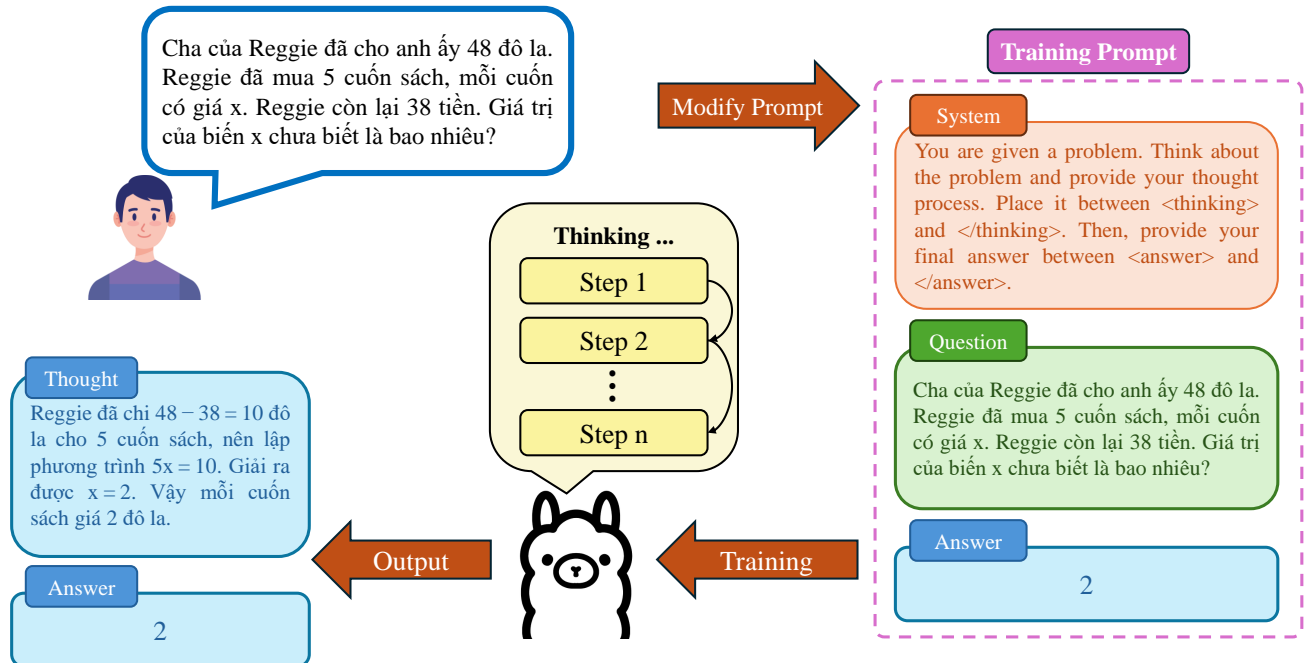


Tutorial: Improving Reasoning Capability of LLMs with Reinforcement Learning

Dinh-Thang Duong, Truong-Binh Duong, Quang-Vinh Dinh

I. Giới thiệu

LLM Reasoning (Tạm dịch: LLM Suy Luận) là khả năng của các mô hình ngôn ngữ lớn trong việc mô phỏng quy trình tư duy logic và giải quyết vấn đề theo từng bước rõ ràng, từ việc trích xuất, liên kết các luận điểm đến diễn giải chi tiết trước khi đưa ra kết quả cuối cùng. Các mô hình như [OpenAI o1](#) hay [DeepSeek-R1](#) đã chứng tỏ thành công ấn tượng khi triển khai khả năng suy luận của mô hình lớn vào nhiều bài toán phức tạp, nâng cao cả độ chính xác lẫn khả năng giải thích của kết quả.



Hình 1: Huấn luyện mô hình ngôn ngữ lớn có khả năng suy luận với dữ liệu giải toán Tiếng Việt.

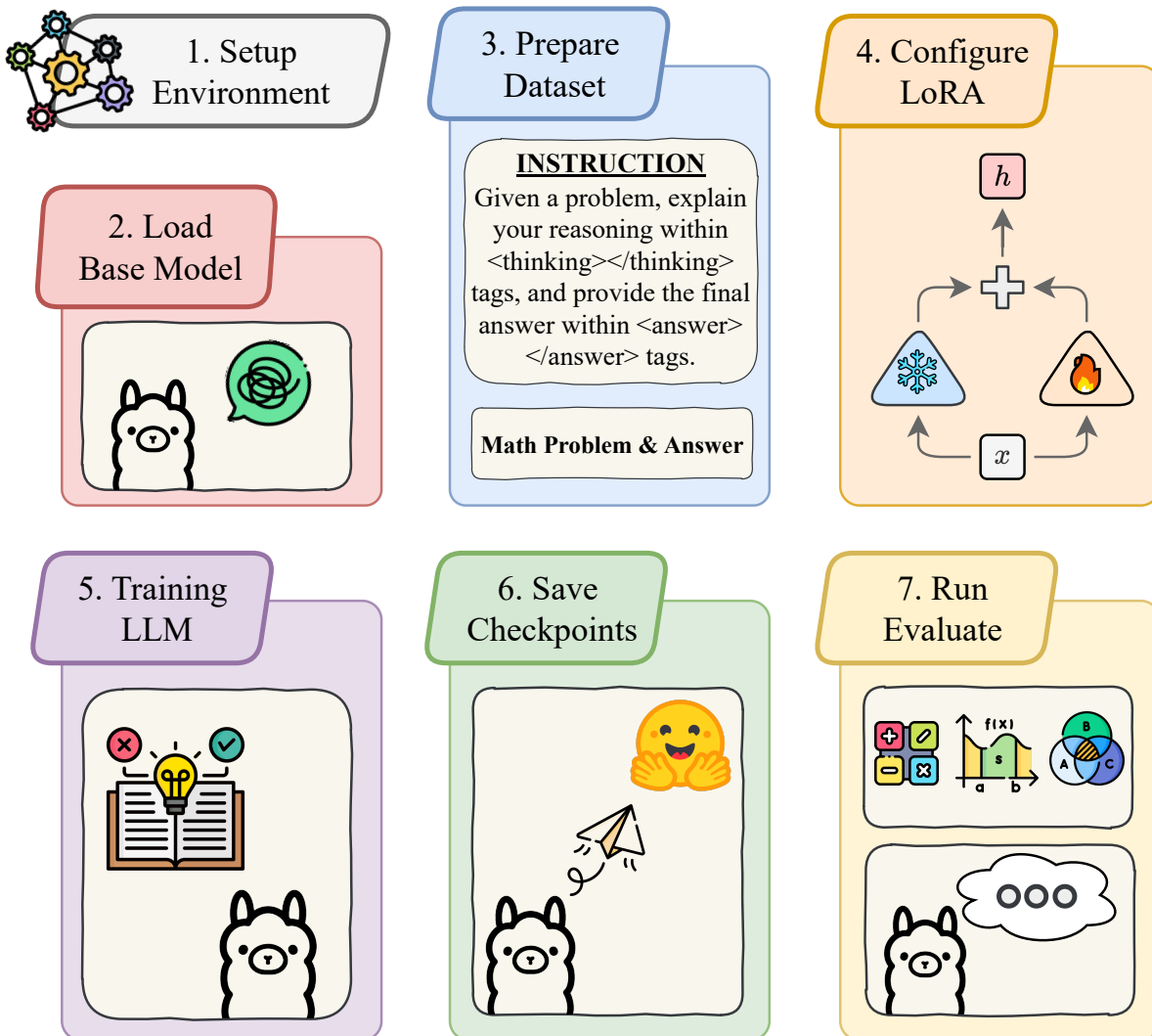
Thay vì chỉ sinh văn bản mạch lạc dựa trên dữ liệu đã học, LLM Reasoning đòi hỏi mô hình “suy nghĩ từng bước” để đưa ra lời giải chính xác và giải thích chi tiết quá trình tư duy, đồng thời cải

thiện tính minh bạch và độ tin cậy trên các nhiệm vụ như giải toán trắc nghiệm, lập luận khoa học hay phân tích ngữ nghĩa.

Trong bài viết này, chúng ta sẽ tìm hiểu cách cài đặt quá trình huấn luyện nhằm tăng cường khả năng suy luận cho một mô hình ngôn ngữ lớn thông qua dữ liệu các câu hỏi toán học. Mục tiêu là xây dựng một mô hình ngôn ngữ lớn chuyên cho việc giải toán trắc nghiệm tiếng Việt, không những có thể đưa ra đáp án đúng mà còn trình bày được chuỗi suy nghĩ logic dẫn đến đáp án đó (dẫn đến việc cải thiện hiệu suất mô hình). Bài toán được định nghĩa với đầu vào và đầu ra như sau:

- **Input:** Nội dung cụ thể của bài toán cần giải quyết.
- **Output:** Lời phản hồi từ mô hình, bao gồm quá trình suy nghĩ và đáp án cho câu hỏi.

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



Hình 2: Huấn luyện khả năng suy luận cho mô hình ngôn ngữ lớn bằng GRPO.

Trong đó:

1. **Setup Environment:** Thiết lập môi trường làm việc, bao gồm cài đặt và import các thư viện phù hợp cho huấn luyện mô hình.
2. **Load Base Model:** Tải mô hình ngôn ngữ lớn (LLM) gốc từ thư viện như Hugging Face làm nền tảng để tinh chỉnh.
3. **Configure LoRA:** Cấu hình LoRA để giảm tài nguyên huấn luyện bằng cách khóa các tầng gốc và chỉ tinh chỉnh các adapter nhẹ.
4. **Prepare Dataset:** Chuẩn bị và tiền xử lý dữ liệu training bao gồm câu hỏi và đáp án đúng theo định dạng đầu vào/đầu ra rõ ràng.
5. **Fine-tune LLM:** Huấn luyện mô hình trên tập dữ liệu toán học đã chuẩn bị bằng GRPO và kỹ thuật LoRA nhằm tối ưu hiệu suất.
6. **Save Checkpoints:** Lưu lại các checkpoint của mô hình trong quá trình huấn luyện để theo dõi tiến độ và phục hồi nếu cần.
7. **Run Evaluate:** Thực hiện suy luận trên các câu hỏi để đánh giá khả năng đọc hiểu của mô hình sau tinh chỉnh.

Thông qua pipeline trên, ta sẽ có một mô hình ngôn ngữ lớn đã được tinh chỉnh với chi phí thấp, có khả năng suy luận và trả lời hiệu quả các câu hỏi toán học.

II. Cài đặt chương trình

II.1. Cài đặt và import các thư viện cần thiết

Đầu tiên, chúng ta cần cài đặt và nạp các thư viện hỗ trợ quá trình huấn luyện với GRPO và LoRA trên nền Llama 3.2 1B. Trong đó:

- unsloth cung cấp lớp FastLanguageModel tối ưu cho inference và fine-tuning.
- vllm hỗ trợ việc tạo sinh output nhanh chóng với cấu hình linh hoạt.
- datasets dùng để tải và xử lý tập dữ liệu MetaMathQA.
- trl bao gồm GRPOConfig và GRPOTrainer cho module học tăng cường.

Chúng ta thực hiện cài đặt unsloth và vllm thông qua lệnh `pip install` dưới đây:

```
1 !pip install unsloth vllm==0.7.3
```

Sau khi cài đặt xong, chúng ta import các thư viện cần thiết để thiết lập mô hình, xử lý dữ liệu và cấu hình GRPO:

```
1 import os
2 import re
3
4 from vllm import SamplingParams
5 from unsloth import FastLanguageModel
6 from datasets import load_dataset, Dataset
7 from trl import GRPOConfig, GRPOTrainer
```

II.2. Load mô hình lớn và setup cài đặt LoRA

Tiếp theo, chúng ta khởi tạo mô hình [Llama 3.2 1B](#) từ Unsloth và kích hoạt LoRA để giảm số tham số cần huấn luyện. Lưu ý rằng, các bạn hoàn toàn có thể thay đổi sang các phiên bản mô hình Llama 3.2 với số lượng tham số lớn hơn để có thể đạt một mô hình với hiệu suất cao hơn nếu phần cứng cho phép.

```
1 max_seq_length = 2048
2 lora_rank = 64
3
4 model, tokenizer = FastLanguageModel.from_pretrained(
5     model_name="meta-llama/Llama-3.2-1B-Instruct",
6     max_seq_length=max_seq_length,
7     load_in_4bit=False,
8     fast_inference=True,
9     max_lora_rank=lora_rank,
10    gpu_memory_utilization=0.8,
11 )
12
13 model = FastLanguageModel.get_peft_model(
14     model,
```

```

15     r=lora_rank,
16     target_modules=[
17         "q_proj", "k_proj", "v_proj", "o_proj",
18         "gate_proj", "up_proj", "down_proj",
19     ],
20     lora_alpha=lora_rank,
21     use_gradient_checkpointing="unsloth",
22     random_state=3407,
23 )

```

- `max_seq_length` định nghĩa độ dài tối đa của input sequence.
- `lora_rank` xác định kích thước của ma trận adapter LoRA.
- `load_in_4bit=False` tắt lượng tử hóa để giữ độ chính xác cho reasoning.
- `fast_inference=True` tối ưu hóa throughput khi sinh text.
- `gpu_memory_utilization=0.8` cho phép sử dụng tối đa 80% bộ nhớ GPU.

II.3. Tải bộ dữ liệu

Trong bài này, ta sẽ tăng cường khả năng suy luận của mô hình LLM trong việc giải toán trắc nghiệm tiếng Việt. Theo đó, bộ dữ liệu [MetaMathQA-40K](#) phiên bản tiếng Việt sẽ được tải về thông qua đoạn code dưới đây:

```

1 dataset = load_dataset("5CD-AI/Vietnamese-meta-math-MetaMathQA-40K-gg-translated", split="
                                train")

```

| Bài toán | Đáp án |
|--|--------|
| Cha của Reggie đã cho anh ấy \$48. Reggie mua 5 cuốn sách, mỗi cuốn giá \$x và còn lại \$38. Hỏi giá trị của \$x là bao nhiêu? | 2 |
| Jean và ba người bạn chơi domino với bộ 28 quân, chia đều cho 4 người. Hỏi mỗi người nhận bao nhiêu quân? | 7 |
| Cally có 10 áo trắng, 5 áo màu, 7 quần đùi và 6 quần dài (tổng 28 bộ). Danny có 6 áo trắng, 8 áo màu, 10 quần đùi và 6 quần dài (tổng 30 bộ). Hỏi họ đã giặt bao nhiêu bộ quần áo? | 58 |

Bảng 1: Một vài mẫu câu hỏi và đáp án trực tiếp (phần tiếng Việt) trong bộ MetaMathQA-40K.

II.4. Cài đặt format prompt cho reasoning

Trước khi huấn luyện, ta cần chuẩn hóa dữ liệu để mô hình học được chuỗi suy nghĩ và đáp án tách biệt. Các bước bao gồm:

1. Dùng regex `answer_pattern` để trích xuất đúng phần “đáp án là...”.

- Đánh dấu bắt đầu/kết thúc chuỗi suy nghĩ bằng `<thinking>...</thinking>` và phần đáp án bằng `<answer>...</answer>`.
- Xây dựng `system_prompt` hướng dẫn model sinh reasoning rồi answer.
- Chuyển danh sách thành `train_dataset` với hai trường `prompt` và `answer`.

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

Cutting Knowledge Date: December 2023

Today Date: 28 Apr 2025

You are given a problem.

Think about the problem and provide your thought process.

Place it between `<thinking>` and `</thinking>`.

Then, provide your final answer between `<SOLUTION></SOLUTION>``<|eot_id|>`

```
<|start_header_id|>user<|end_header_id|>
```

Cha của Reggie đã cho anh ấy 48 đô la. Reggie đã mua 5 cuốn sách, mỗi cuốn có giá x. Reggie còn lại 38 tiền. Giá trị của biến x chưa biết là bao nhiêu?`<|eot_id|>`

```
<|start_header_id|>assistant<|end_header_id|>
```

Hình 3: Nội dung đầy đủ chi tiết của một mẫu prompt khi đưa vào mô hình Llama với một số thông tin về việc reasoning cũng như thông tin về bài toán.

```
1 answer_pattern = re.compile(
2     r"(đáp án là:|đáp án là :|câu trả lời là:|câu trả lời là :)\s*(.*)",
3     re.IGNORECASE
4 )
5
6 formatted_dataset = []
7 for item in dataset:
8     response = item["response_vi"].strip().lower()
9     match = answer_pattern.search(response)
10    if match:
11        answer = match.group(2).strip()
12        formatted_dataset.append({
13            "question": item["query_vi"],
14            "answer": answer
15        })
16
17 reasoning_start = "<thinking>"
18 reasoning_end = "</thinking>"
19 solution_start = "<SOLUTION>"
```

```

20 solution_end = "</SOLUTION>"
21
22 system_prompt = \
23     f"""You are given a problem.
24     Think about the problem and provide your thought process.
25     Place it between {reasoning_start} and {reasoning_end}.
26     Then, provide your final answer between {solution_start}{solution_end}"""
27
28 train_dataset = Dataset.from_list(formatted_dataset[:8000])
29 train_dataset = train_dataset.map(lambda x: {
30     "prompt": [
31         {"role": "system", "content": system_prompt},
32         {"role": "user", "content": x["question"]},
33     ],
34     "answer": x["answer"],
35 })

```

II.5. Định nghĩa các hàm reward cho học tăng cường

Đối với các thuật toán học tăng cường, việc đánh giá hiệu suất mô hình thường được thông qua các hàm điểm thưởng (reward function). Trong trường hợp này, các hàm reward được cân nhắc triển khai đánh giá chất lượng output theo hai tiêu chí: đúng format reasoning và đúng đáp án. Các bạn hoàn toàn có thể tìm hiểu để đưa ra thêm các tiêu chí điểm thưởng phù hợp.

II.5.1. Reward cho output đúng format toán

Đầu tiên, ta tạo ra hai hàm có chức năng tính điểm cho mô hình nếu nó cho ra một format output phù hợp. Ở đây, ta có hàm `match_format_exactly()`, nếu hoàn toàn khớp pattern reasoning + answer, cộng 3.0 điểm. Bên cạnh đó, hàm `match_format_approximately()` sẽ kiểm tra số lượng tag xuất hiện, cộng 0.5 mỗi tag đúng, trừ 1.0 nếu không đúng.

```

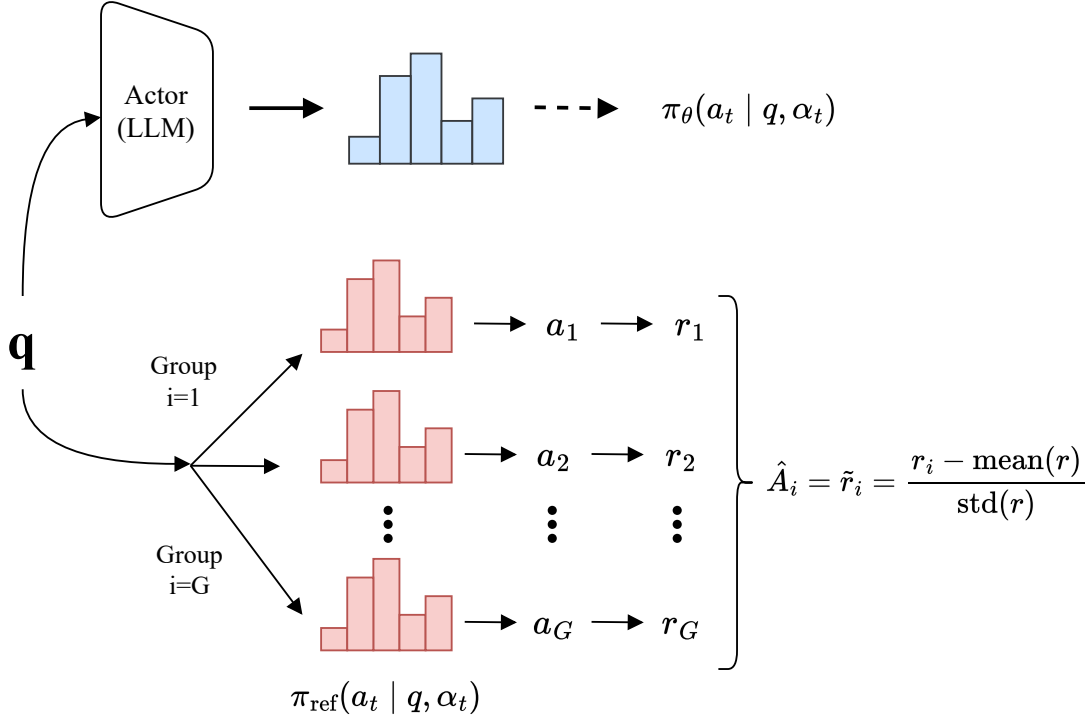
1 match_format = re.compile(
2     rf"^\[s]{0,}"
3     rf"{reasoning_start}.+?{reasoning_end}.+?"
4     rf"{solution_start}(.+){solution_end}"
5     rf"\[s]{0,}\$",
6     flags=re.MULTILINE | re.DOTALL
7 )
8
9 def match_format_exactly(completions, **kwargs):
10     scores = []
11     for completion in completions:
12         score = 0
13         response = completion[0]["content"]
14         if match_format.search(response) is not None:
15             score += 3.0
16         scores.append(score)
17     return scores
18
19 def match_format_approximately(completions, **kwargs):
20     scores = []
21     for completion in completions:

```

```

22     score = 0
23     response = completion[0]["content"]
24     score += 0.5 if response.count(reasoning_start) == 1 else -1.0
25     score += 0.5 if response.count(reasoning_end) == 1 else -1.0
26     score += 0.5 if response.count(solution_start) == 1 else -1.0
27     score += 0.5 if response.count(solution_end) == 1 else -1.0
28     scores.append(score)
29     return scores

```



Hình 4: Minh họa quy trình GRPO, trong đó LLM đóng vai trò actor sinh phân phối xác suất cho hành động, sau đó lấy mẫu thành nhiều nhóm, đánh giá reward cho từng nhóm và chuẩn hóa kết quả để cập nhật chính sách.

II.5.2. Reward cho output đúng đáp án

Tiếp theo, ta định nghĩa hàm `check_answer()` để so khớp chuỗi giữa tag `<answer>`, cộng 3.0 nếu đúng, 1.5 nếu chỉ khác khoảng trắng, trừ 1.5 nếu sai hoàn toàn. Cuối cùng, ta có hàm `check_numbers()` với chức năng trích xuất số, so sánh giá trị float, cộng 1.5 nếu đúng, trừ 0.5 nếu sai.

```

1 match_numbers = re.compile(
2     solution_start + r".*?([\d\.\,]{1,})",
3     flags=re.MULTILINE | re.DOTALL
4 )
5
6 def check_answer(prompts, completions, answer, **kwargs):
7     responses = [completion[0]["content"] for completion in completions]
8

```



```

9     extracted_responses = [
10         guess.group(1)
11         if (guess := match_format.search(r)) is not None else None
12         for r in responses
13     ]
14
15     scores = []
16     for guess, true_answer in zip(extracted_responses, answer):
17         score = 0
18         if guess is None:
19             scores.append(0)
20             continue
21         if guess == true_answer:
22             score += 3.0
23         elif guess.strip() == true_answer.strip():
24             score += 1.5
25         else:
26             score -= 1.5
27         scores.append(score)
28     return scores
29
30 def check_numbers(prompts, completions, answer, **kwargs):
31     question = prompts[0][-1]["content"]
32     responses = [completion[0]["content"] for completion in completions]
33
34     extracted_responses = [
35         guess.group(1)
36         if (guess := match_numbers.search(r)) is not None else None
37         for r in responses
38     ]
39
40     count = getattr(check_numbers, 'counter', 0) + 1
41     check_numbers.counter = count
42     if count % 5 == 0:
43         print('*'*20, f"Question:{question}", f"\nResponse:\n{responses[0]}",
44               f"\nExtracted: {extracted_responses[0]}", f"\nGT Answer: {answer[0]}")
45
46     scores = []
47     for guess, true_answer in zip(extracted_responses, answer):
48         if guess is None:
49             scores.append(0)
50             continue
51         try:
52             true_answer = float(true_answer.strip())
53             # Remove commas like in 123,456
54             guess = float(guess.strip().replace(",", ""))
55             scores.append(1.5 if guess == true_answer else -0.5)
56         except:
57             scores.append(0)
58     return scores

```

II.6. Huấn luyện mô hình

Với tất cả các thông tin trên mà chúng ta đã khai báo, để tiến hành huấn luyện, ta chỉ cần cấu hình và chạy GRPOTrainer để thực hiện training học tăng cường với các hàm reward đã xây dựng:

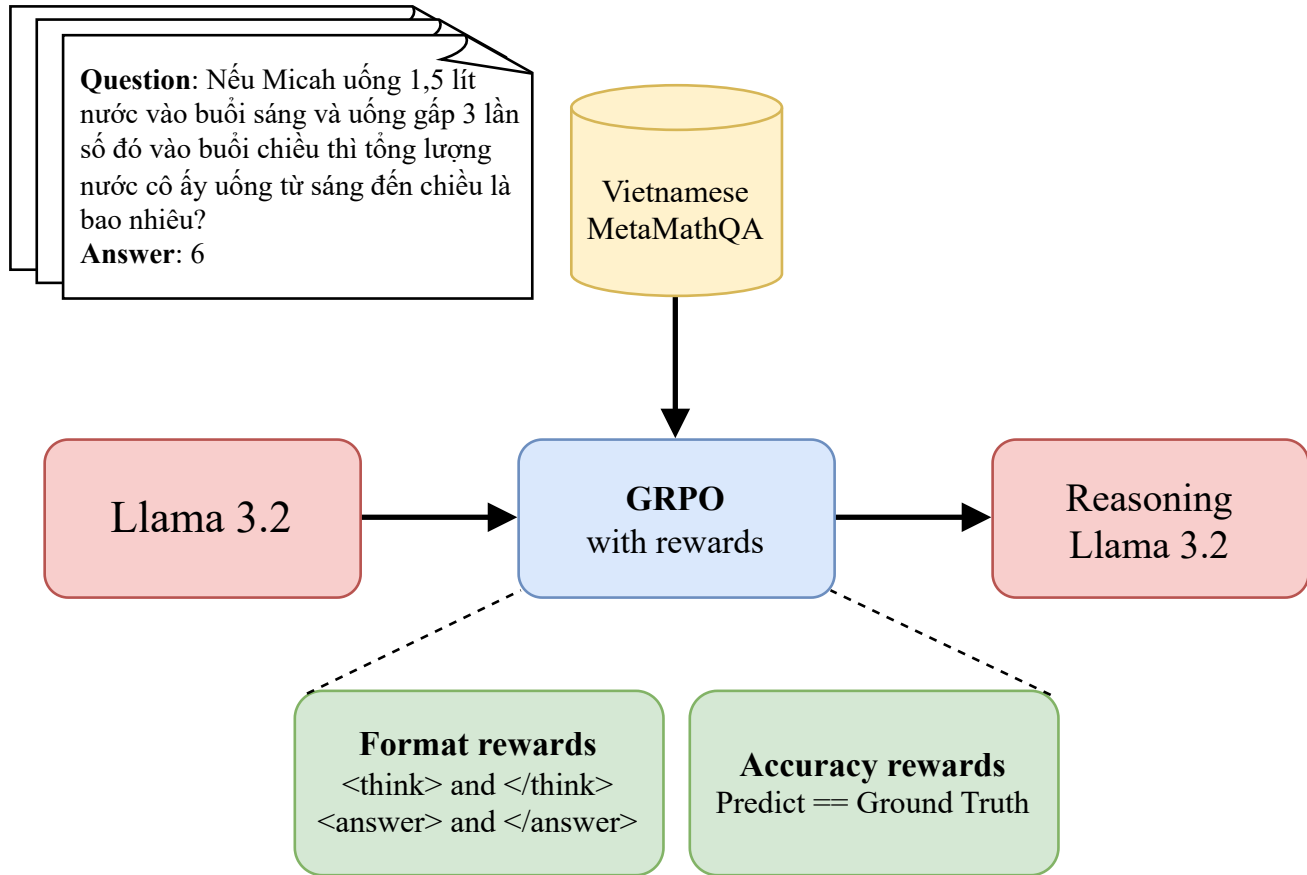
```

1 max_len = max(dataset.map(
2     lambda x: {"tokens": tokenizer.apply_chat_template(
3         x["prompt"], add_generation_prompt=True, tokenize=True)},
4     batched=True,
5 ).map(lambda x: {"length": len(x["tokens"])}))["length"])
6
7 max_prompt_length = max_len + 1
8
9 training_args = GRPOConfig(
10     learning_rate=5e-6,
11     weight_decay=5e-4,
12     warmup_ratio=0.1,
13     lr_scheduler_type="cosine",
14     optim="adamw_torch_fused",
15     logging_steps=1,
16     per_device_train_batch_size=2,
17     gradient_accumulation_steps=64,
18     num_generations=8,
19     max_prompt_length=max_prompt_length,
20     max_completion_length=max_seq_length - max_prompt_length,
21     num_train_epochs=1,
22     max_steps=-1,
23     save_steps=250,
24     max_grad_norm=0.1,
25     report_to="wandb",
26     output_dir="outputs_bz2",
27 )
28
29 trainer = GRPOTrainer(
30     model=model,
31     processing_class=tokenizer,
32     reward_funcs=[
33         match_format_exactly,
34         match_format_approximately,
35         check_answer,
36         check_numbers,
37     ],
38     args=training_args,
39     train_dataset=dataset,
40 )
41 trainer.train()

```

- Tính `max_prompt_length()` dựa trên độ dài token tối đa của prompt đã format.
- Cấu hình `GRPOConfig` với learning rate, weight decay, warmup ratio, scheduler cosine, batch size, số lượng generations, gradient accumulation, và các bước lưu checkpoint.
- Khởi tạo `GRPOTrainer` với danh sách các hàm reward và dataset phù hợp.

- Gọi `trainer.train()` để bắt đầu huấn luyện.



Hình 5: Sử dụng thuật toán học tăng cường GRPO nhằm khuyến khích việc đưa ra các suy luận trước khi tạo sinh phương án cuối cùng của mô hình Llama.

II.7. Lưu trữ và inference

Sau khi huấn luyện xong, ta lưu adapter LoRA đã học và tiến hành thử suy luận trên một ví dụ:

```
1 model.save_lora("grpo_saved_lora")
```

```
1 idx = 0
2 messages = [
3     {"role": "system", "content": system_prompt},
4     {"role": "user", "content": train_dataset[idx]["question"]},
5 ]
6 sampling_params = SamplingParams(
7     temperature = 0.8,
8     top_p = 0.95,
9     max_tokens = 1024,
10 )
11
12 text = tokenizer.apply_chat_template(
13     messages,
```

```

14     add_generation_prompt = True,
15     tokenize = False,
16 )
17
18 path_lora = "grpo_saved_lora"
19 output = model.fast_generate(
20     [text],
21     sampling_params = sampling_params,
22     lora_request = model.load_lora(path_lora),
23 ) [0].outputs[0].text
24
25 print(f"Problem: \n{train_dataset[idx]['question']}")
26 print(f"Response: \n{output}")
27 print("GT Answer:", train_dataset[idx]["answer"])

```

<thinking>

Để tìm giá trị của x , ta cần hiểu rằng Reggie đã chi một khoản tiền $48 - 38 = 10$ đô la để mua các sách. Vì Reggie đã mua 5 cuốn sách với giá x mỗi cuốn nên ta có thể tạo ra một công thức để tính tổng chi phí mua sách: $5x = 10$.

Để giải quyết cho x , ta chia cả hai bên của công thức bằng 5: $x = 10 / 5$.

Solvers lại ta có: $x = 2$.

Vậy giá của mỗi cuốn sách là 2 đô la.

</thinking>

<SOLUTION>2</SOLUTION>

Hình 6: Minh họa một kết quả tạo sinh của mô hình với phần reasoning và đáp án cuối cùng.

Ví dụ minh họa ở Hình 6 cho thấy Llama 3.2 1B sau khi được tinh chỉnh với GRPO và LoRA đã có khả năng sinh chuỗi suy nghĩ chi tiết bằng thẻ <thinking> trước khi đưa ra đáp án cuối cùng trong thẻ <SOLUTION>. Điều này khẳng định tính hiệu quả của việc kết hợp học tăng cường và kỹ thuật LoRA, không chỉ nâng cao độ chính xác mà còn cải thiện tính giải thích được (explainability) cho các nhiệm vụ yêu cầu tư duy logic chặt chẽ. Theo đó, các bạn hoàn toàn có thể mở rộng phương pháp này cho các bài toán phức tạp hơn hoặc kết hợp thêm các tín hiệu đánh giá đa dạng để tiếp tục tối ưu khả năng suy luận của LLMs.

III. Câu hỏi trắc nghiệm

1. Theo bài viết, “LLM Reasoning” là gì?
 - (a) Khả năng mô hình ngôn ngữ lớn sinh văn bản một cách mạch lạc.
 - (b) Khả năng mô hình ngôn ngữ lớn mô phỏng quy trình tư duy logic và giải quyết vấn đề theo từng bước.
 - (c) Khả năng mô hình ngôn ngữ lớn dịch thuật giữa các ngôn ngữ.
 - (d) Khả năng mô hình ngôn ngữ lớn tóm tắt văn bản dài.
2. Mô hình ngôn ngữ lớn (LLM) gốc nào được sử dụng làm nền tảng trong bài viết?
 - (a) OpenAI o1.
 - (b) DeepSeek-R1.
 - (c) Llama-3.2-1B-Instruct.
 - (d) GPT-4.
3. Kỹ thuật nào được sử dụng để giảm số lượng tham số cần huấn luyện, bằng cách khóa các tầng gốc và chỉ tinh chỉnh adapter?
 - (a) Quantization (Lượng tử hóa).
 - (b) Gradient Checkpointing.
 - (c) LoRA (Low-Rank Adaptation).
 - (d) GRPO (Group Relative Policy Optimization).
4. Tập dữ liệu nào được sử dụng để huấn luyện mô hình giải toán tiếng Việt?
 - (a) GSM8K.
 - (b) MATH.
 - (c) Vietnamese-meta-math-MetaMathQA-40K-gg-translated.
 - (d) Alpaca.
5. Cặp thẻ nào được sử dụng để đánh dấu phần chuỗi suy nghĩ của mô hình trong format prompt?
 - (a) `<logic>...</logic>`.
 - (b) `<reasoning>...</reasoning>`.
 - (c) `<stepbystep>...</stepbystep>`.
 - (d) `<thinking>...</thinking>`.
6. Cặp thẻ nào được sử dụng để đánh dấu phần đáp án cuối cùng của mô hình trong format prompt?
 - (a) `<solution>...</solution>`.
 - (b) `<final_answer>...</final_answer>`.

- (c) `<result>...</result>`.
- (d) `<answer>...</answer>`.
7. Trong hàm `reward match_format_exactly()`, mô hình nhận được bao nhiêu điểm nếu output hoàn toàn khớp với pattern reasoning + answer?
- (a) 1.0 điểm.
- (b) 1.5 điểm.
- (c) 3.0 điểm.
- (d) 0.5 điểm.
8. Trong hàm `reward check_answer()`, mô hình nhận được bao nhiêu điểm nếu đáp án đoán ra (guess) hoàn toàn trùng khớp (exact match) với đáp án đúng (true_answer)?
- (a) 1.5 điểm.
- (b) 3.0 điểm.
- (c) -1.5 điểm.
- (d) 1.0 điểm.
9. Thuật toán nào được sử dụng để huấn luyện mô hình với các hàm reward đã định nghĩa?
- (a) PPO (Proximal Policy Optimization).
- (b) DPO (Direct Preference Optimization).
- (c) GRPO (Group Relative Policy Optimization).
- (d) AdamW.
10. Hàm `reward check_numbers()` được thiết kế để đánh giá tiêu chí nào của output?
- (a) Kiểm tra output có tuân thủ đúng định dạng thẻ `<thinking>` và `<answer>` hay không.
- (b) So sánh toàn bộ chuỗi văn bản trong thẻ `<answer>` với đáp án gốc.
- (c) Trích xuất và so sánh giá trị số học (dạng float) trong phần đáp án của output với đáp án gốc.
- (d) Đánh giá độ dài và chi tiết của phần giải thích trong thẻ `<thinking>`.

IV. Phụ lục

1. **Datasets:** Các file dataset được đề cập trong bài có thể được tải tại [đây](#).
2. **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 chủ nhật khi hết deadline phần nội dung này, ad mới copy các tài liệu bài giải vào đường dẫn trên).
3. **Demo:** Web demo và mã nguồn của ứng dụng có thể được truy cập tại [đây](#).
4. **Rubric:**

| Mục | Kiến Thức | Đánh Giá |
|-----|--|---|
| I. | <ul style="list-style-type: none"> - Kiến thức về mô hình lớn (LLMs). - Kiến thức về học tăng cường cơ bản và thuật toán học tăng cường GRPO. - Khái niệm về Reasoning trong mô hình lớn và các cách thức kích hoạt khả năng này cho mô hình lớn. - Kiến thức về thư viện PyTorch và HuggingFace. - Triển khai PEFT (LoRA) và thiết lập các tham số liên quan đến LoRA. | <ul style="list-style-type: none"> - Nắm được lý thuyết về khả năng suy luận của mô hình lớn (LLM Reasoning). - Hiểu được cơ bản về lý thuyết học tăng cường cũng như thuật toán GRPO được ứng dụng. - Nắm được lý thuyết cơ bản về kỹ thuật LoRA. - Có khả năng cài đặt và thực hiện huấn luyện học tăng cường mô hình lớn trên một bộ dữ liệu bất kì nhằm cải thiện khả năng suy luận của một mô hình ngôn ngữ lớn. |

- Hết -