

CS336 Assignment 5 (alignment): Alignment and Reasoning RL

Version 1.0.2

CS336 Staff

Spring 2025

1 Assignment Overview

In this assignment, you will gain some hands-on experience with training language models to reason when solving math problems.

What you will implement.

1. Zero-shot prompting baseline for the MATH dataset of competition math problems Hendrycks et al. [2021].
2. Supervised finetuning, given reasoning traces from a stronger reasoning model (DeepSeek R1, DeepSeek-AI et al. 2025).
3. Expert Iteration for improving reasoning performance with verified rewards.
4. Group-Relative Policy Optimization (GRPO) for improving reasoning performance with verified rewards.

For those interested, we will have an **entirely optional** part of the assignment on aligning language models to human preferences, which will be released in the next few days.

What you will run.

1. Measure Qwen 2.5 Math 1.5B zero-shot prompting performance (our baseline).
2. Run SFT on Qwen 2.5 Math 1.5B with reasoning traces from R1.
3. Run Expert Iteration on Qwen 2.5 Math 1.5B with verified rewards.
4. Run GRPO on Qwen 2.5 Math 1.5B with verified rewards.

What the code looks like. All the assignment code as well as this writeup are available on GitHub at:

`github.com/stanford-cs336/assignment5-alignment`

Please `git clone` the repository. If there are any updates, we will notify you and you can `git pull` to get the latest.

1. `cs336_alignment/`: This is where you'll write your code for assignment 5. Note that there's no code in here (aside from a little starter code), so you should be able to do whatever you want from scratch.

2. `cs336_alignment/prompts/*`: For your convenience, we’ve provided text files with prompts to minimize possible errors caused by copying-and-pasting prompts from the PDF to your code.
3. `tests/*.py`: This contains all the tests that you must pass. **You are only expected to pass the tests in `tests/test_sft.py` and `tests/test_grpo.py`—the rest of the tests are for the non-mandatory parts of the assignment.** These tests invoke the hooks defined in `tests/adapters.py`. You’ll implement the adapters to connect your code to the tests. Writing more tests and/or modifying the test code can be helpful for debugging your code, but your implementation is expected to pass the original provided test suite.
4. `README.md`: This file contains some basic instructions on setting up your environment.

What you can use. We expect you to build most of the RL related components from scratch. You may use tools like vLLM to generate text from language models (§3.1). In addition, you may use HuggingFace Transformers to load the Qwen 2.5 Math 1.5B model and tokenizer and run forward passes (§4.1), but you may not use any of the training utilities (e.g., the `Trainer` class).

How to submit. You will submit the following files to Gradescope:

- `writeup.pdf`: Answer all the written questions. Please typeset your responses.
- `code.zip`: Contains all the code you’ve written.

2 Reasoning with Language Models

2.1 Motivation

One of the remarkable use cases of language models is in building generalist systems that can handle a wide range of natural language processing tasks. In this assignment, we will focus on a developing use case for language models: mathematical reasoning. It will serve as a testbed for us to set up evaluations, perform supervised finetuning, and experiment with teaching LMs to reason using reinforcement learning (RL).

There are going to be two differences from the way we’ve done our past assignments.

- First, we are not going to be using our language model codebase and models from earlier. We would ideally like to use base language models trained from previous assignments, but finetuning those models will not give us a satisfying result—these models are far too weak to display non-trivial mathematical reasoning capabilities. Because of this, we are going to switch to a modern, high-performance language model that we can access (Qwen 2.5 Math 1.5B Base) and do most of our work on top of that model.
- Second, we are going to introduce a new benchmark with which to evaluate our language models. Up until this point, we have embraced the view that cross-entropy is a good surrogate for many downstream tasks. However, the point of this assignment will be to bridge the gap between base models and downstream tasks and so we will have to use evaluations that are separate from cross-entropy. We will use the MATH 12K dataset from Hendrycks et al. [2021], which consists of challenging high-school competition mathematics problems. We will evaluate language model outputs by comparing them against a reference answer.

2.2 Chain-of-Thought Reasoning and Reasoning RL

An exciting recent trend in language models is the use of *chain-of-thought* reasoning to improve performance across a variety of tasks. Chain-of-thought refers to the process of reasoning through a problem step-by-step, generating intermediate reasoning steps before arriving at a final answer.

Chain-of-thought reasoning with LLMs. Early chain-of-thought approaches finetuned language models to solve simple mathematical tasks like arithmetic by using a “scratchpad” to break the problem into intermediate steps [Nye et al., 2021]. Other work prompts a strong model to “think step by step” before answering, finding that this significantly improves performance on mathematical reasoning tasks like grade-school math questions [Wei et al., 2023].

Learning to reason with expert iteration. The Self-Taught Reasoner (STaR) [Zelikman et al., 2022] frames reasoning as a bootstrapping loop: a pretrained model first samples diverse chains-of-thought (CoTs), keeps only those that lead to correct answers, and then finetunes on these “expert” traces. Iterating this cycle can improve the LM’s reasoning capabilities and solve rate. STaR demonstrated that this version of expert iteration [Anthony et al., 2017] using automatic, string match-based verification of generated answers can bootstrap reasoning skills without human-written reasoning traces.

Reasoning RL with verified rewards, o1, and R1. Recent work has explored using more powerful reinforcement learning algorithms with verified rewards to improve reasoning performance. OpenAI’s o1 (and subsequent o3/o4) [OpenAI et al., 2024], DeepSeek’s R1 [DeepSeek-AI et al., 2025], and Moonshot’s kimi k1.5 [Team et al., 2025] use policy gradient methods [Sutton et al., 1999] to train on math and code tasks where string matching or unit tests verify correctness, demonstrating remarkable improvements in competition math and coding performance. Later works such as Open-R1 [Face, 2025], SimpleRL-Zoo [Zeng et al., 2025], and TinyZero [Pan et al., 2025] confirm that pure reinforcement learning with verified rewards—even on models as small as 1.5B parameters—can improve reasoning performance.

Our setup: model and dataset. In the following sections, we will consider progressively more complex approaches to train a base language model to reason step-by-step in order to solve math problems. For this assignment, we will be using the Qwen 2.5 Math 1.5B Base model, which was continually pretrained from the Qwen 2.5 1.5B model on high-quality synthetic math pretraining data [Yang et al., 2024]. The MATH dataset is available on the Together cluster at `/data/a5-alignment/MATH`.

Tip for Open-Source Auditors: Alternative Datasets

Unfortunately, the MATH dataset is not publicly available due to a copyright claim. If you are following along at home, you can use one of the following open-source mathematical reasoning datasets:

- Countdown [Pan et al., 2025], available here: a simple synthetic task based on the British TV show Countdown that has served as a popular testbed for small-scale reasoning RL.
- GSM8K [Cobbe et al., 2021a], available here: grade-school math problems, which are easier than MATH but should allow you to debug correctness and get familiar with the reasoning RL pipeline.
- Tulu 3 SFT Math [Lambert et al., 2025], available here: synthetic math problems generated using GPT-4o and Claude 3.5 Sonnet. Because these are synthetic, some answers (or even the questions) may not be entirely correct.
- Some other math SFT dataset linked here.

To obtain short ground-truth labels (e.g., 1/2) if they are not provided directly, you can process the ground-truth column with a math answer parser such as Math-Verify.

3 Measuring Zero-Shot MATH Performance

We’ll start by measuring the performance of our base language model on the 5K example test set of MATH. Establishing this baseline is useful for understanding how each of the later approaches affects model behavior.

Unless otherwise specified, for experiments on MATH we will use the following prompt from the DeepSeek R1-Zero model [DeepSeek-AI et al., 2025]. We will refer to this as the `r1_zero` prompt:

```
A conversation between User and Assistant. The User asks a question, and the Assistant
→ solves it. The Assistant first thinks about the reasoning process in the mind and
→ then provides the User with the answer. The reasoning process is enclosed within
→ <think> </think> and answer is enclosed within <answer> </answer> tags, respectively,
→ i.e., <think> reasoning process here </think> <answer> answer here </answer>.
User: {question}
Assistant: <think>
```

The `r1_zero` prompt is located in the text file `cs336_alignment/prompts/r1_zero.prompt`.

In the prompt, `question` refers to some question that we insert (e.g., Natalia sold clips to 48 of her friends in April, and then she sold half as many clips in May. How many clips did Natalia sell altogether in April and May?). The expectation is that the model plays the role of the assistant, and starts generating the thinking process (since we have already included a left think tag `<think>`), closes the thinking process with `</think>` and then generates a final symbolic answer within the answer tags, like `<answer> 4x + 10 </answer>`. The purpose of having the model generate tags like `<answer> </answer>` is so that we can easily parse the model's output and compare it against a ground truth answer, and so that we can stop response generation when we see the right answer tag `</answer>`.

Note on prompt choice. It turns out that the `r1_zero` prompt is not the best choice for maximizing downstream performance after RL, because of a mismatch between the prompt and how the Qwen 2.5 Math 1.5B model was pretrained. Liu et al. [2025] finds that simply prompting the model with the question (and nothing else) starts with a very high accuracy, e.g., matching the `r1_zero` prompt after 100+ steps of RL. Their findings suggest that Qwen 2.5 Math 1.5B was already pretrained on such question-answer pairs.

Nonetheless, we choose the `r1_zero` prompt for this assignment because RL with it shows clear accuracy improvements in a short number of steps, allowing us to walk through the mechanics of RL and sanity check correctness quickly, even if we don't manage the best final performance. As a reality check, you will compare directly to the `question_only` prompt later in the assignment.

3.1 Using vLLM for offline language model inference

To evaluate our language models, we're going to have to generate continuations (responses) for a variety of prompts. While one could certainly implement their own functions for generation (e.g., as you did in assignment 1), efficient implementation of RL requires high-performance inference techniques, and implementing these inference techniques are beyond the scope of this assignment. Therefore, in this assignment we will recommend using vLLM for offline batched inference. vLLM is a high-throughput and memory-efficient inference engine for language models that incorporates a variety of useful efficiency techniques (e.g., optimized CUDA kernels, PagedAttention for efficient attention KV caching [Kwon et al., 2023], etc.). To use vLLM to generate continuations for a list of prompts:¹

```
from vllm import LLM, SamplingParams

# Sample prompts.
prompts = [
    "Hello, my name is",
    "The president of the United States is",
    "The capital of France is",
    "The future of AI is",
]
```

¹Example taken from https://github.com/vllm-project/vllm/blob/main/examples/offline_inference.py.

```

# Create a sampling params object, stopping generation on newline.
sampling_params = SamplingParams(
    temperature=1.0, top_p=1.0, max_tokens=1024, stop=["\n"]
)

# Create an LLM.
llm = LLM(model=<path to model>)

# Generate texts from the prompts. The output is a list of RequestOutput objects
# that contain the prompt, generated text, and other information.
outputs = llm.generate(prompts, sampling_params)

# Print the outputs.
for output in outputs:
    prompt = output.prompt
    generated_text = output.outputs[0].text
    print(f"Prompt: {prompt!r}, Generated text: {generated_text!r}")

```

In the example above, the LLM can be initialized with the name of a HuggingFace model (which will be automatically downloaded and cached if it isn't found locally), or a path to a HuggingFace model. Since downloads can take a long time (especially for larger models, e.g., 70B parameters) and to conserve cluster disk space (so everyone doesn't have their own independent copy of the pre-trained models), we have downloaded the following pre-trained models at the following the paths on the Together cluster. **Please do not re-download these models on the Together cluster:**

- Qwen 2.5 Math 1.5B Base (for reasoning experiments):
/data/a5-alignment/models/Qwen2.5-Math-1.5B
- Llama 3.1 8B Base (for optional instruction tuning experiments):
/data/a5-alignment/models/Llama-3.1-8B
- Llama 3.3 70B Instruct (for optional instruction tuning experiments):
/data/a5-alignment/models/Llama-3.3-70B-Instruct

3.2 Zero-shot MATH Baseline

Prompting setup. To evaluate zero-shot performance on the MATH test set, we'll simply load the examples and prompt the language model to answer the question using the `r1_zero` prompt from above.

Evaluation metric. When we evaluate a multiple-choice or binary response task, the evaluation metric is clear—we test whether the model outputs exactly the correct answer.

In math problems we assume that there is a known ground truth (e.g. 0.5) but we cannot simply test whether the model outputs exactly 0.5—it can also answer `<answer> 1/2 </answer>`. Because of this, we must address the tricky problem of matching for semantically equivalent responses from the LM when we evaluate MATH.

To this end, we want to come up with some answer parsing function that takes as input the model's output and a known ground-truth, and returns a boolean indicating whether the model's output is correct. For example, a reward function could receive the model's string output ending in `<answer> She sold 15 clips. </answer>` and the gold answer 72, and return `True` if the model's output is correct and `False` otherwise (in this case, it should return `False`).

For our MATH experiments, we will use a fast and fairly accurate answer parser used in recent work on reasoning RL [Liu et al., 2025]. This reward function is implemented at `cs336_alignment.drgrepo_grader.r1_zero_reward_fn`, and you should use it to evaluate performance on MATH unless otherwise specified.

Generation hyperparameters. When generating responses, we'll sample with temperature 1.0, top-p 1.0, max generation length 1024. The prompt asks the model to end its answer with the string `</answer>`, and therefore we can direct vLLM to stop when the model outputs this string:

```
# Based on Dr. GRPO: stop when the model completes its answer
# https://github.com/sail-sg/understand-r1-zero/blob/
# c18804602b85da9e88b4aeeb6c43e2f08c594fbc/train_zero_math.py#L167
sampling_params.stop = ["</answer>"]
sampling_params.include_stop_str_in_output = True
```

Problem (math_baseline): 4 points

- (a) Write a script to evaluate Qwen 2.5 Math 1.5B zero-shot performance on MATH. This script should (1) load the MATH validation examples from `/data/a5-alignment/MATH/validation.jsonl`, (2) format them as string prompts to the language model using the `r1_zero` prompt, and (3) generate outputs for each example. This script should also (4) calculate evaluation metrics and (5) serialize the examples, model generations, and corresponding evaluation scores to disk for analysis in subsequent problems.

It might be helpful for your implementation to include a method `evaluate_vllm` with arguments similar to the following, as you will be able to reuse it later:

```
def evaluate_vllm(
    vllm_model: LLM,
    reward_fn: Callable[[str, str], dict[str, float]],
    prompts: List[str],
    eval_sampling_params: SamplingParams
) -> None:
    """
    Evaluate a language model on a list of prompts,
    compute evaluation metrics, and serialize results to disk.
    """
```

```
python scripts/math_baseline.py \
    --model Qwen/Qwen2.5-Math-1.5B \
    --prompt_path
/home/zks/Disk/2025FirstSemester/CS336LargeLanguage
Model/LAB/assignment5-
alignment/cs336_alignment/prompts/r1_zero.prompt \
    --out_dir ./eval_gsm8k_r1zero
```

Deliverable: A script to evaluate baseline zero-shot MATH performance.

```
{
  "prompt_path": "/home/zks/Disk/2025FirstSemester/CS336LargeLanguageModel/LAB/a
ssignment5-alignment/cs336_alignment/prompts/r1_zero.prompt",
  "dataset": "openai/gsm8k:main",
  "split": "test",
  "max_examples": -1,
  "sampling": {
    "temperature": 0.0,
    "top_p": 1.0,
    "max_tokens": 512,
    "seed": 42
  },
  "metrics": {
    "n": 1319,
    "format_rate": 0.5041698256254739,
    "answer_accuracy": 0.17134192570128887,
    "reward_mean": 0.17134192570128887,
    "counts": {
      "format=1 answer=1": 226,
      "format=1 answer=0": 439,
      "format=0 answer=0": 654,
      "format=0 answer=1": 0
    }
  }
}
```

- (b) Run your evaluation script on Qwen 2.5 Math 1.5B. How many model generations fall into each of the following categories: (1) correct with both format and answer reward 1, (2) format reward 1 and answer reward 0, (3) format reward 0 and answer reward 0? Observing at least 10 cases where format reward is 0, do you think the issue is with the base model's output, or the parser? Why? What about in (at least 10) cases where format reward is 1 but answer reward is 0?

Deliverable: Commentary on the model and reward function performance, including examples of each category.

3) format reward 0 并不只在 base model , 也很大概率出在 parser / scorer 上。而且证据相当明显。许多是输出的格式不符合parse的要求直接format就失败了, 还有部分直接进行了复读。

format reward is 1 but answer reward is 0:这些样本 格式全部是合格的 (format_reward=1) , 但 答案错误 / 不符合 gold_final / 不是期望的“纯数值”形式

- (c) How well does the Qwen 2.5 Math 1.5B zero-shot baseline perform on MATH?

Deliverable: 1-2 sentences with evaluation metrics.

4 Supervised Finetuning for MATH

Algorithm 1 Supervised Finetuning (SFT)

Input initial policy model $\pi_{\theta_{\text{init}}}$; SFT dataset \mathcal{D}

- 1: policy model $\pi_{\theta} \leftarrow \pi_{\theta_{\text{init}}}$
- 2: **for** step = 1, ..., **n_sft_steps** **do**
- 3: Sample a batch of question-response pairs \mathcal{D}_b from \mathcal{D}
- 4: Compute the cross-entropy loss of the responses given the questions using the model π_{θ}
- 5: Update the model parameters θ by taking a gradient step with respect to the cross-entropy loss
- 6: **end for**

Output π_{θ}

Supervised finetuning for reasoning. In this section, we will finetune our base model on the MATH dataset (Algorithm 1). As our goal is to improve the model’s reasoning ability, rather than finetune it to directly predict correct answers, we will finetune it to first generate a chain-of-thought reasoning trace followed by an answer. To this end, we have made available a dataset of such reasoning traces, obtained from DeepSeek R1 DeepSeek-AI et al. [2025], in

`/data/a5-alignment/MATH/sft.jsonl`

When training a reasoning model in practice, SFT is often used as a warm-start for a second RL finetuning step. There are two main reasons for this. First, SFT requires high-quality annotated data (i.e., with pre-existing reasoning traces), whereas RL requires only the correct answer for feedback. Second, even in settings where annotated data is plentiful, RL can still unlock performance gains by finding better policies than the SFT data. Unfortunately, the models we use are not big enough to show effects when composing SFT and RL, so for this assignment we will treat these two phases separately.

4.1 Using HuggingFace Models

Loading a HuggingFace model and tokenizer. To load a HuggingFace model and tokenizer from a local dir (in bfloat16 and with FlashAttention-2 to save memory), you can use the following starter code:

```
from transformers import AutoModelForCausalLM, AutoTokenizer

model = AutoModelForCausalLM.from_pretrained(
    "/data/a5-alignment/models/Qwen2.5-Math-1.5B",
    torch_dtype=torch.bfloat16,
    attn_implementation="flash_attention_2",
)
tokenizer = AutoTokenizer.from_pretrained("/data/a5-alignment/models/Qwen2.5-Math-1.5B")
```

Forward pass. After we’ve loaded the model, we can run a forward pass on a batch of input IDs and get the logits (with the `.logits`) attribute of the output. Then, we can compute the loss between the model’s predicted logits and the actual labels:

```
input_ids = train_batch["input_ids"].to(device)
labels = train_batch["labels"].to(device)

logits = model(input_ids).logits
loss = F.cross_entropy(..., ...)
```


Saving a trained model. To save the model to a directory after training is finished, you can use the `.save_pretrained()` function, passing in the path to the desired output directory. Make sure to save under `/data/yourusername` since they can be quite large. We recommend also saving the tokenizer as well (even if you didn't modify it), just so the model and tokenizer are self-contained and loadable from a single directory.

```
# Save the model weights
model.save_pretrained(save_directory=output_dir)
tokenizer.save_pretrained(save_directory=output_dir)
```

Gradient accumulation. Despite loading the model in `bfloat16` and using FlashAttention-2, even an 80GB GPU does not have enough memory to support reasonable batch sizes. To use larger batch sizes, we can use a technique called *gradient accumulation*. The basic idea behind gradient accumulation is that rather than updating our model weights (i.e., taking an optimizer step) after every batch, we'll *accumulate* the gradients over several batches before taking a gradient step. Intuitively, if we had a larger GPU, we should get the same results from computing the gradient on a batch of 32 examples all at once, vs. splitting them up into 16 batches of 2 examples each and then averaging at the end.

Gradient accumulation is straightforward to implement in PyTorch. Recall that each weight tensor has an attribute `.grad` that stores its gradient. Before we call `loss.backward()`, the `.grad` attribute is `None`. After we call `loss.backward()`, the `.grad` attribute contains the gradient. Normally, we'd take an optimizer step, and then zero the gradients with `optimizer.zero_grad()`, which resets the `.grad` field of the weight tensors:

```
for inputs, labels in data_loader:
    # Forward pass.
    logits = model(inputs)
    loss = loss_fn(logits, labels)

    # Backward pass.
    loss.backward()

    # Update weights.
    optimizer.step()
    # Zero gradients in preparation for next iteration.
    optimizer.zero_grad()
```

To implement gradient accumulation, we'll just call the `optimizer.step()` and `optimizer.zero_grad()` every `k` steps, where `k` is the number of gradient accumulation steps. We divide the loss by `gradient_accumulation_steps` before calling `loss.backward()` so that the gradients are averaged across the gradient accumulation steps.

```
gradient_accumulation_steps = 4
for idx, (inputs, labels) in enumerate(data_loader):
    # Forward pass.
    logits = model(inputs)
    loss = loss_fn(logits, labels) / gradient_accumulation_steps

    # Backward pass.
    loss.backward()

    if (idx + 1) % gradient_accumulation_steps == 0:
        # Update weights every `gradient_accumulation_steps` batches.
        optimizer.step()
```

```
# Zero gradients every `gradient_accumulation_steps` batches.
optimizer.zero_grad()
```

As a result, our effective batch size when training is multiplied by k , the number of gradient accumulation steps.

4.2 SFT Helper Methods

Next, we will implement some helper methods that you will use during SFT and in the later RL experiments. As a quick note on nomenclature: in the following sections, we will interchangeably refer to a model's completion given a prompt as an “output”, “completion”, or “response”.

Tokenizing prompts and outputs. For each pair of question and target output (q, o) , we will tokenize the question and output separately and concatenate them. Then, we can score the log-probabilities of the output with our SFT model (or in later sections, our RL policy). Moreover, we will need to construct a `response_mask`: a boolean mask that is `True` for all tokens in the response, and `False` for all question and padding tokens. We will use this mask in the training loop to ensure that we only compute the loss on the response tokens.

Problem (tokenize_prompt_and_output): Prompt and output tokenization (2 points)

Deliverable: Implement a method `tokenize_prompt_and_output` that tokenizes the question and output separately, concatenates them together, and constructs a `response_mask`. The following interface is recommended:

```
def tokenize_prompt_and_output(prompt_strs, output_strs, tokenizer):
    """Tokenize the prompt and output strings, and construct a mask that is 1 for the response tokens and 0 for other tokens (prompt or padding)."""
```

Args:

`prompt_strs`: `list[str]` List of prompt strings.

`output_strs`: `list[str]` List of output strings.

`tokenizer`: `PreTrainedTokenizer` Tokenizer to use for tokenization.

Returns:

`dict[str, torch.Tensor]`. Let `prompt_and_output_lens` be a list containing the lengths of the tokenized prompt and output strings. Then the returned dictionary should have the following keys:

`input_ids` `torch.Tensor` of shape `(batch_size, max(prompt_and_output_lens) - 1)`: the tokenized prompt and output strings, with the final token sliced off.

`labels` `torch.Tensor` of shape `(batch_size, max(prompt_and_output_lens) - 1)`: shifted input ids, i.e., the input ids without the first token.

`response_mask` `torch.Tensor` of shape `(batch_size, max(prompt_and_output_lens) - 1)`: a mask on the response tokens in the labels.

To test your code, implement `[adapters.run_tokenize_prompt_and_output]`. Then, run the test with `uv run pytest -k test_tokenize_prompt_and_output` and make sure your implementation passes it.

Logging per-token entropies. When doing RL, it is often useful to keep track of per-token entropies to see if the predictive distribution of the model is becoming (over)confident. We will implement this now and compare how each of our finetuning approaches affects the model’s predictive entropy.

The entropy of a discrete distribution $p(x)$ with support \mathcal{X} is defined as

$$H(p) = - \sum_{x \in \mathcal{X}} p(x) \log p(x). \quad (1)$$

Given our SFT or RL model’s logits, we will compute the per-token entropy, i.e., the entropy of each next-token prediction.

Problem (compute_entropy): Per-token entropy (1 point)

Deliverable: Implement a method `compute_entropy` that computes the per-token entropy of next-token predictions.

The following interface is recommended:

```
def compute_entropy(logits: torch.Tensor) -> torch.Tensor:
```

Get the entropy of the next-token predictions (i.e., entropy over the vocabulary dimension).

Args:

logits: `torch.Tensor` Tensor of shape `(batch_size, sequence_length, vocab_size)` containing unnormalized logits.

Returns:

torch.Tensor Shape `(batch_size, sequence_length)`. The entropy for each next-token prediction.

Note: you should use a numerically stable method (e.g., using `logsumexp`) to avoid overflow.

To test your code, implement `[adapters.run_compute_entropy]`. Then run `uv run pytest -k test_compute_entropy` and ensure your implementation passes.

Getting log-probabilities from a model. Obtaining log-probabilities from a model is a primitive that we will need in both SFT and RL.

For a prefix x , an LM producing next-token logits $f_\theta(x) \in \mathbb{R}^{|\mathcal{V}|}$, and a label $y \in \mathcal{V}$, the log-probability of y is

$$\log p_\theta(y | x) = \log [\text{softmax}(f_\theta(x))]_y, \quad (2)$$

where the notation $[x]_y$ denotes the y -th element of the vector x .

You will want to use a numerically stable method to compute this, and are free to use methods from `torch.nn.functional`. We also suggest including an argument to optionally compute and return token entropies.

Problem (get_response_log_probs): Response log-probs (and entropy) (2 points)

Deliverable: Implement a method `get_response_log_probs` that gets per-token conditional log-probabilities (given the previous tokens) from a causal language model, and optionally the entropy of the model’s next-token distribution.

The following interface is recommended:

```
def get_response_log_probs(
    model: PreTrainedModel,
    input_ids: torch.Tensor,
    labels: torch.Tensor,
    return_token_entropy: bool = False,
) -> dict[str, torch.Tensor]:
```

Args:

model: PreTrainedModel HuggingFace model used for scoring (placed on the correct device and in inference mode if gradients should not be computed).

input_ids: torch.Tensor shape (batch_size, sequence_length), concatenated prompt + response tokens as produced by your tokenization method.

labels: torch.Tensor shape (batch_size, sequence_length), labels as produced by your tokenization method.

return_token_entropy: bool If **True**, also return per-token entropy by calling `compute_entropy`.

Returns:

dict[str, torch.Tensor].

"log_probs" shape (batch_size, sequence_length), conditional log-probabilities $\log p_{\theta}(x_t | x_{<t})$.

"token_entropy" optional, shape (batch_size, sequence_length), per-token entropy for each position (present only if `return_token_entropy=True`).

Implementation tips:

- Obtain logits with `model(input_ids).logits`.

To test your code, implement `[adapters.run_get_response_log_probs]`. Then run `uv run pytest -k test_get_response_log_probs` and ensure the test passes.

SFT microbatch train step. The loss we minimize in SFT is the negative log-likelihood of the target output given the prompt. To compute this loss, we need to compute the log-probabilities of the target output given the prompt and sum over all tokens in the output, masking the tokens in the prompt and padding tokens.

We will implement a helper function for this, that we will also make use of later during RL.

Problem (masked_normalize): Masked normalize (1 point)

Deliverable: Implement a method `masked_normalize` that sums over tensor elements and normalizes by a constant while respecting a boolean mask.

The following interface is recommended:

```
def masked_normalize(
    tensor: torch.Tensor,
    mask: torch.Tensor,
    normalize_constant: float,
    dim: int | None = None,
) -> torch.Tensor:
```

Sum over a dimension and normalize by a constant, considering only those elements where `mask == 1`.

Args:

tensor: `torch.Tensor` The tensor to sum and normalize.

mask: `torch.Tensor` Same shape as **tensor**; positions with 1 are included in the sum.

normalize_constant: `float` the constant to divide by for normalization.

dim: `int` | `None` the dimension to sum along before normalization. If `None`, sum over all dimensions.

Returns:

`torch.Tensor` the normalized sum, where masked elements (`mask == 0`) don't contribute to the sum.

To test your code, implement `[adapters.run_masked_normalize]`. Then run `uv run pytest -k test_masked_normalize` and ensure it passes.

SFT microbatch train step. We are now ready to implement a single microbatch train step for SFT (recall that for a train minibatch, we iterate over many microbatches if `gradient_accumulation_steps > 1`).

Problem (`sft_microbatch_train_step`): Microbatch train step (3 points)

Deliverable: Implement a single micro-batch update for SFT, including cross-entropy loss, summing with a mask, and gradient scaling.

The following interface is recommended:

```
def sft_microbatch_train_step(
    policy_log_probs: torch.Tensor,
    response_mask: torch.Tensor,
    gradient_accumulation_steps: int,
    normalize_constant: float = 1.0,
) -> tuple[torch.Tensor, dict[str, torch.Tensor]]:
```

Execute a forward-and-backward pass on a microbatch.

Args:

policy_log_probs (`batch_size`, `sequence_length`), per-token log-probabilities from the SFT policy being trained.

response_mask (`batch_size`, `sequence_length`), 1 for response tokens, 0 for prompt/padding.

gradient_accumulation_steps Number of microbatches per optimizer step.

normalize_constant The constant by which to divide the sum. It is fine to leave this as 1.0.

Returns:

`tuple[torch.Tensor, dict[str, torch.Tensor]]`.

loss scalar tensor. The microbatch loss, adjusted for gradient accumulation. We return this so we can log it.

metadata Dict with metadata from the underlying loss call, and any other statistics you might want to log.

Implementation tips:

- You should call `loss.backward()` in this function. Make sure to adjust for gradient accumulation.

To test your code, implement `[adapters.run_sft_microbatch_train_step]`. Then run `uv run pytest -k test_sft_microbatch_train_step` and confirm it passes.

Logging generations in-the-loop. It's always good practice to do some in-the-loop logging that involves generation from your model, and reasoning SFT/RL is no exception. Write a function `log_generations` that will prompt your model to generate responses for some given prompts (e.g., sampled from the validation set). It's a good idea to log at least the following for each example:

1. The input prompt.
2. The response generated by the SFT/RL model.
3. The ground-truth answer.
4. The reward information, including format, answer, and total reward.
5. The average token entropy of the response.
6. The average response length, average response length for correct responses, and average response length for incorrect responses.

Problem (`log_generations`): Logging generations (1 point)

Deliverable: Implement a function `log_generations` that can be used to log generations from your model.

4.3 SFT Experiment

Using the pieces above, you will now implement the full SFT procedure (Algorithm 1) to finetune the Qwen 2.5 Math 1.5B Base model on the MATH dataset. Each example in `/data/a5-alignment/MATH/sft.jsonl` consists of a formatted prompt and a target response, where the target response includes a chain-of-thought reasoning trace and the final answer. In particular, each example is a JSON element of type `{"prompt": str, "response": str}`.

In order to track the progress of your model over the course of training, you should periodically evaluate it on the MATH validation set. You should run your script with 2 GPUs, using one GPU for the policy model and the other for the vLLM instance to evaluate the policy. To get this to work, here is some starter code to initialize vLLM and to load the policy weights into the vLLM instance before every rollout phase:

```
from vllm.model_executor import set_random_seed as vllm_set_random_seed

def init_vllm(model_id: str, device: str, seed: int, gpu_memory_utilization: float = 0.85):
    """
    Start the inference process, here we use vLLM to hold a model on
    a GPU separate from the policy.
```

```

"""
vllm_set_random_seed(seed)

# Monkeypatch from TRL:
# https://github.com/huggingface/trl/blob/
# 22759c820867c8659d00082ba8cf004e963873c1/trl/trainer/grpo_trainer.py
# Patch vLLM to make sure we can
# (1) place the vLLM model on the desired device (world_size_patch) and
# (2) avoid a test that is not designed for our setting (profiling_patch).
world_size_patch = patch("torch.distributed.get_world_size", return_value=1)
profiling_patch = patch(
    "vllm.worker.worker.Worker._assert_memory_footprint_increased_during_profiling",
    return_value=None
)
with world_size_patch, profiling_patch:
    return LLM(
        model=model_id,
        device=device,
        dtype=torch.bfloat16,
        enable_prefix_caching=True,
        gpu_memory_utilization=gpu_memory_utilization,
    )

def load_policy_into_vllm_instance(policy: PreTrainedModel, llm: LLM):
    """
    Copied from https://github.com/huggingface/trl/blob/
    22759c820867c8659d00082ba8cf004e963873c1/trl/trainer/grpo_trainer.py#L670.
    """
    state_dict = policy.state_dict()
    llm_model = llm.llm_engine.model_executor.driver_worker.model_runner.model
    llm_model.load_weights(state_dict.items())

```

You may find it helpful to log metrics with respect to both the train and validation steps (this will also be useful in later RL experiments). To do this in wandb, you can use the following code:

```

# Setup wandb metrics
wandb.define_metric("train_step") # the x-axis for training
wandb.define_metric("eval_step") # the x-axis for evaluation

# everything that starts with train/ is tied to train_step
wandb.define_metric("train/*", step_metric="train_step")

# everything that starts with eval/ is tied to eval_step
wandb.define_metric("eval/*", step_metric="eval_step")

```

Lastly, we suggest that you use gradient clipping with clip value 1.0.

Problem (sft_experiment): Run SFT on the MATH dataset (2 points) (2 H100 hrs)

1. Run SFT on the reasoning SFT examples (provided in /data/a5-alignment/MATH/sft.jsonl) using the Qwen 2.5 Math 1.5B base model, varying the number of unique examples for SFT in

sft训练

```
bash scripts/run_sft_sweep.sh
```

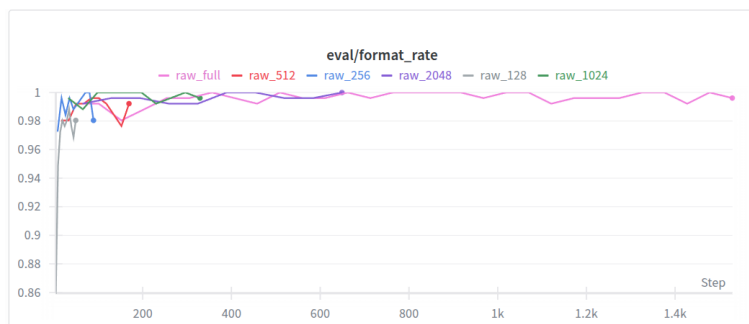


the range {128, 256, 512, 1024}, along with using the full dataset. Tune the learning rate and batch size to achieve at least 15% validation accuracy when using the full dataset.

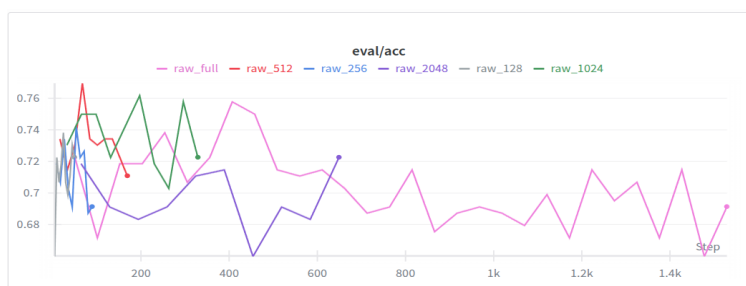
Deliverable: Validation accuracy curves associated with different dataset sizes.

2. Filter the reasoning SFT examples to only include examples that produce the correct answer. Run SFT on the (full) filtered dataset and report the size of the filtered dataset and the validation accuracy you achieve.

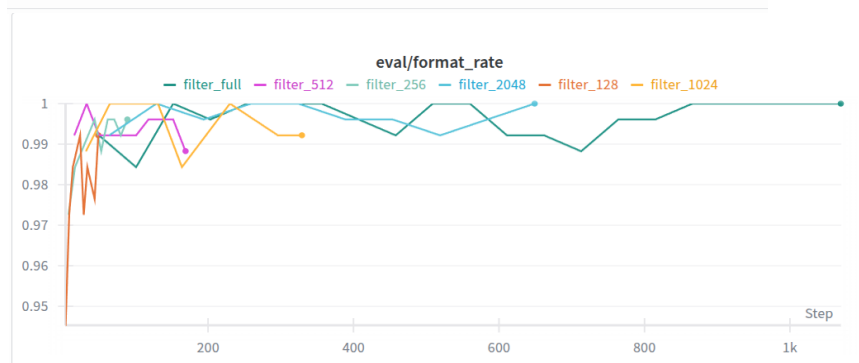
Deliverable: Report the size of the dataset and the validation accuracy curve you achieve. Compare your findings to the previous SFT experiment.



raw直接训练，大概在2048之后基本达到formate的最高点，使用full反而有些许下降



在2048和full都明显下降，说明推理能力下滑，之前具体差别不大



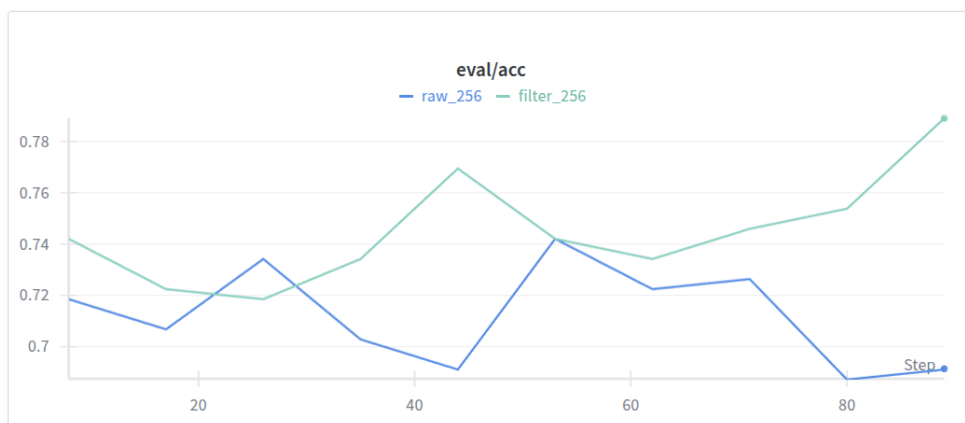
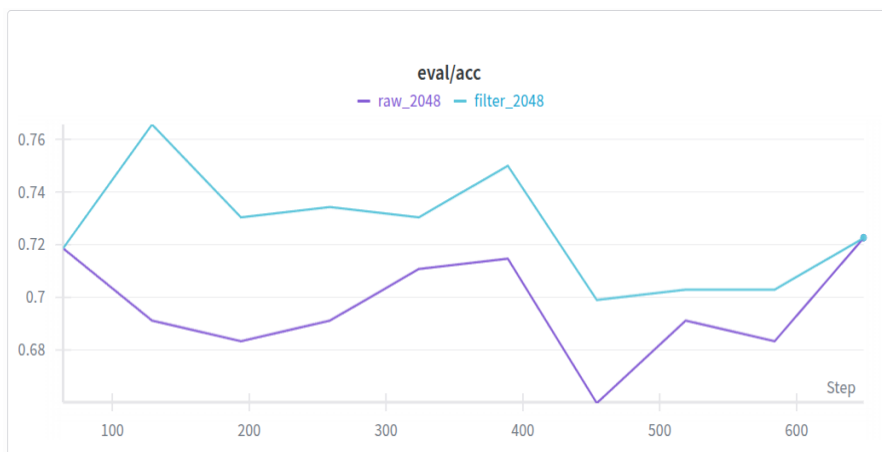
大概在1024左右效果最佳





format的效果差别不大，但是使用filter后，正确的数据，对于整个的正确率提升明显

并且在各个尺度上都差距明显



5 Expert Iteration for MATH

In the previous section, we observed that we can improve the performance of our SFT model by filtering out bad examples from the SFT data. In this section, we will go one step further: we will apply this filtering procedure to reasoning traces we generate from our base model itself. This process is known in the literature as *expert iteration* [Anthony et al., 2017], and in the context of language models has been explored in Cobbe et al. [2021b], Zelikman et al. [2022], Dohan et al. [2022], Gulcehre et al. [2023].

Algorithm 2 Expert iteration (EI)

Input initial policy model $\pi_{\theta_{\text{init}}}$; reward function R ; task questions \mathcal{D}

```
1: policy model  $\pi_{\theta} \leftarrow \pi_{\theta_{\text{init}}}$ 
2: for step = 1, ..., n_ei_steps do
3:   Sample a batch of questions  $\mathcal{D}_b$  from  $\mathcal{D}$ 
4:   Set the old policy model  $\pi_{\theta_{\text{old}}} \leftarrow \pi_{\theta}$ 
5:   Sample  $G$  outputs  $\{o^{(i)}\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | q)$  for each question  $q \in \mathcal{D}_b$ 
6:   Compute rewards  $\{r^{(i)}\}_{i=1}^G$  for each sampled output  $o^{(i)}$  by running reward function  $R(q, o^{(i)})$ 
7:   Filter out wrong outputs (i.e.,  $o^{(i)}$  with  $r^{(i)} = 0$ ) to obtain a dataset  $\mathcal{D}_{\text{sft}}$  of correct question-response pairs
8:    $\pi_{\theta} \leftarrow \text{SFT}(\pi_{\theta}, \mathcal{D}_{\text{sft}})$  (Algorithm 1)
9: end for
```

Output π_{θ}

Next, we will run expert iteration on the MATH dataset.

As a tip, you should pass a `min_tokens` value to your vLLM `SamplingParams`, which will ensure that you do not generate an empty string (which could then cause a NaN downstream depending on your implementation). This can be done with

```
sampling_min_tokens = 4
sampling_params = SamplingParams(
    temperature=sampling_temperature,
    max_tokens=sampling_max_tokens,
    min_tokens=sampling_min_tokens,
    n=G,
    seed=seed,
)
```

As in SFT, you should use gradient clipping with clip value 1.0.

Problem (expert_iteration_experiment): Run expert iteration on the MATH dataset (2 points) (6 H100 hrs)

Run expert iteration on the MATH dataset (provided at `/data/a5-alignment/MATH/train.jsonl`)

```
bash scripts/run_sft_sweep_ei.sh
```



using the Qwen 2.5 Math 1.5B Base model, varying the number of rollouts G per question and the number of epochs used in the SFT step, and using `n_ei_steps` = 5. Vary the batch size for each expert iteration step (i.e., the size of \mathcal{D}_b) in {512, 1024, 2048}. (You do not need to try all possible combinations of these hyperparameters. Just enough to draw conclusions about each is fine.) Log the entropy of the model's reponses over training. Make sure to have vLLM terminate generations at the second answer tag `</answer>`, as done in the SFT section.

Deliverable: Validation accuracy curves associated with different rollout configurations. Try at least 2 different rollout counts and epoch counts.

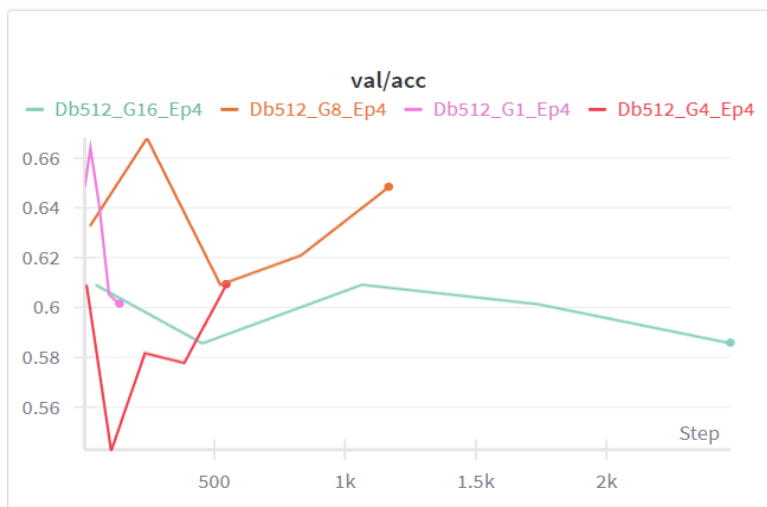
Deliverable: A model that achieves validation accuracy of at least 15% on MATH.

Deliverable: A brief 2 sentence discussion comparing to your SFT performance, as well as performance across EI steps.

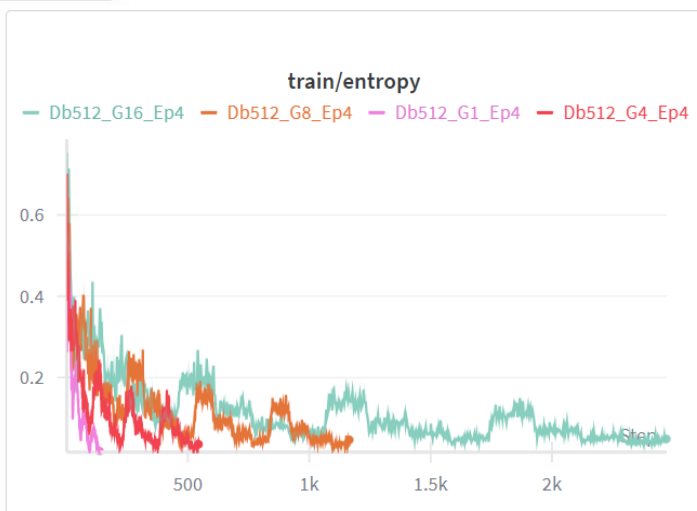
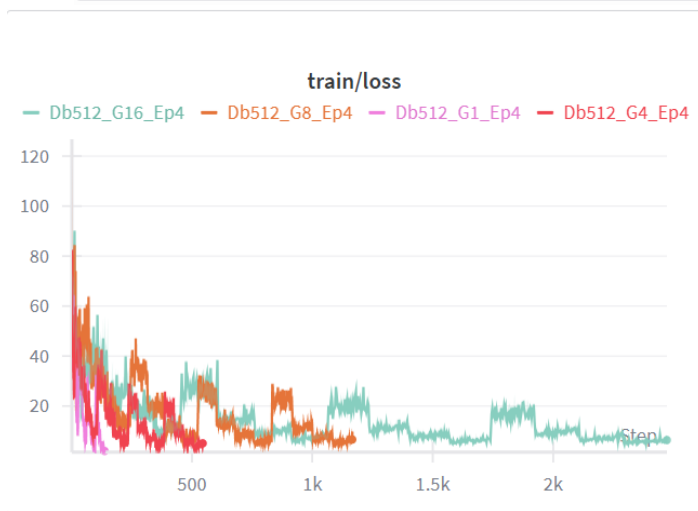
Deliverable: A plot of the entropy of the model's responses over training.

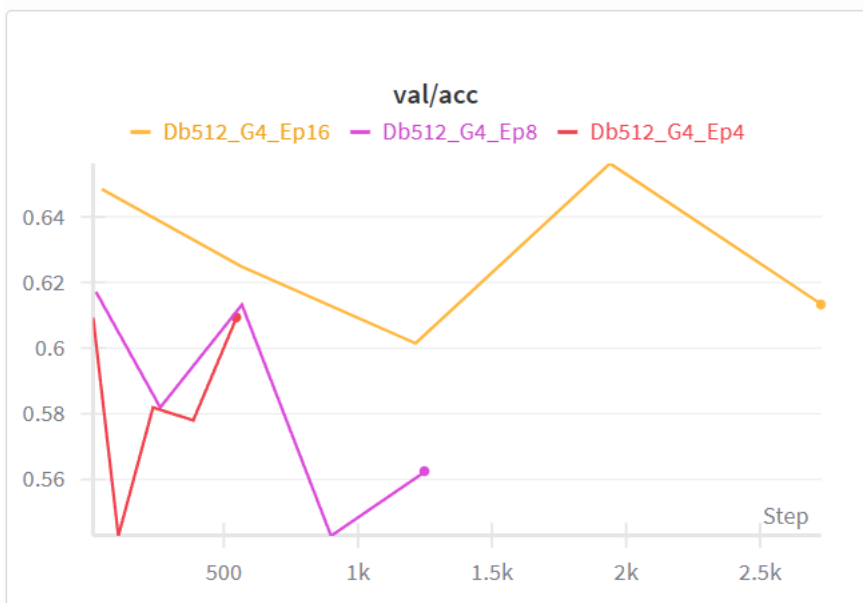


batchsize
和上一个类似



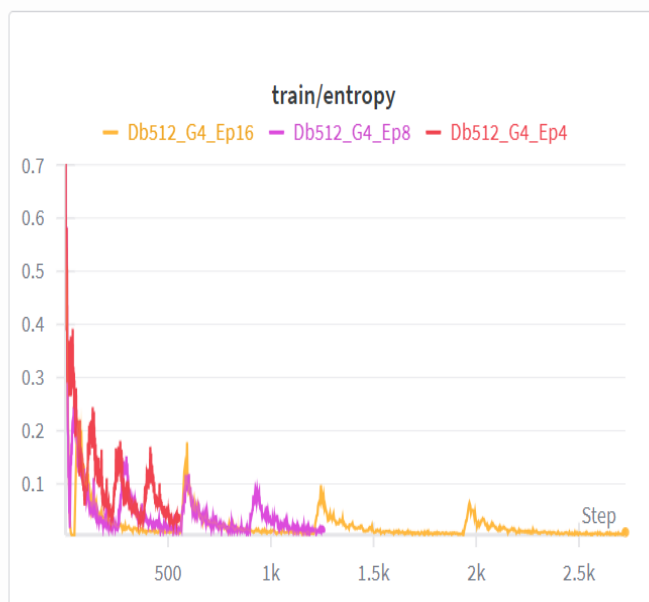
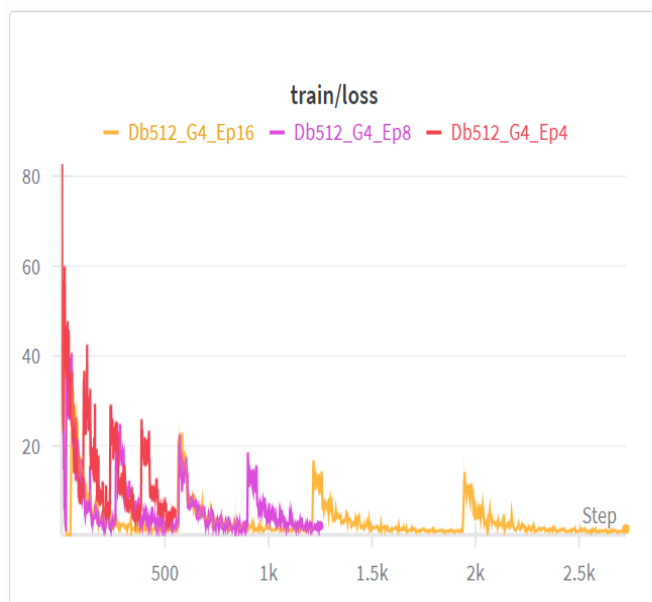
rollout
基本上1-8上升，结果rollout取到16的时候下降
EI 本质是用“当前策略采样出的专家”做监督学习。G 越大，最好的样本往往越“离群”（极端高分、极端长、极端罕见）





epoch
明显越大越好

⋮ ▾ train 2



6 Primer on Policy Gradients

An exciting new finding in language model research is that performing RL against verified rewards with strong base models can lead to significant improvements in their reasoning capabilities and performance [OpenAI et al., 2024, DeepSeek-AI et al., 2025]. The strongest such open reasoning models, such as DeepSeek R1 and Kimi k1.5 [Team et al., 2025], were trained using policy gradients, a powerful reinforcement learning algorithm that can optimize arbitrary reward functions.

We provide a brief introduction to policy gradients for RL on language models below. Our presentation is based closely on a couple great resources which walk through these concepts in more depth: OpenAI’s Spinning Up in Deep RL [Achiam, 2018a] and Nathan Lambert’s Reinforcement Learning from Human Feedback (RLHF) Book [Lambert, 2024].

6.1 Language Models as Policies

A causal language model (LM) with parameters θ defines a probability distribution over the next token $a_t \in \mathcal{V}$ given the current text prefix s_t (the state/observation). In the context of RL, we think of the next token a_t as an *action* and the current text prefix s_t as the *state*. Hence, the LM is a *categorical stochastic policy*

$$a_t \sim \pi_\theta(\cdot \mid s_t), \quad \pi_\theta(a_t \mid s_t) = [\text{softmax}(f_\theta(s_t))]_{a_t}. \quad (3)$$

Two primitive operations will be needed in optimizing the policy with policy gradients:

1. *Sampling from the policy*: drawing an action a_t from the categorical distribution above;
2. *Scoring the log-likelihood of an action*: evaluating $\log \pi_\theta(a_t \mid s_t)$.

Generally, when doing RL with LLMs, s_t is the partial completion/solution produced so far, and each a_t is the next token of the solution; the episode ends when an end-of-text token is emitted, like `<|end_of_text|>`, or `</answer>` in the case of our `r1_zero` prompt.

6.2 Trajectories

A (finite-horizon) trajectory is the interleaved sequence of states and actions experienced by an agent:

$$\tau = (s_0, a_0, s_1, a_1, \dots, s_T, a_T), \quad (4)$$

where T is the length of the trajectory, i.e., a_T is an end-of-text token or we have reached a maximum generation budget in tokens.

The initial state is drawn from the start distribution, $s_0 \sim \rho_0(s_0)$; in the case of RL with LLMs, $\rho_0(s_0)$ is a distribution over formatted prompts. In general settings, state transitions follow some environment

dynamics $s_{t+1} \sim P(\cdot \mid s_t, a_t)$. In RL with LLMs, the environment is deterministic: the next state is the old prefix concatenated with the emitted token, $s_{t+1} = s_t \parallel a_t$. Trajectories are also called *episodes* or *rollouts*; we will use these terms interchangeably.

6.3 Rewards and Return

A scalar reward $r_t = R(s_t, a_t)$ judges the immediate quality of the action taken at state s_t . For RL on verified domains, it is standard to assign zero reward to intermediate steps and a *verified reward* to the terminal action

$$r_T = R(s_T, a_T) := \begin{cases} 1 & \text{if the trajectory } s_T \parallel a_T \text{ matches the ground-truth according to our reward function} \\ 0 & \text{otherwise.} \end{cases}$$

The *return* $R(\tau)$ aggregates rewards along the trajectory. Two common choices are *finite-horizon undiscounted* returns

$$R(\tau) := \sum_{t=0}^T r_t, \quad (5)$$

and *infinite-horizon discounted* returns

$$R(\tau) := \sum_{t=0}^{\infty} \gamma^t r_t, \quad 0 < \gamma < 1. \quad (6)$$

In our case, we will use the undiscounted formulation since episodes have a natural termination point (end-of-text or max generation length).

The objective of the agent is to maximize the expected return

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)], \quad (7)$$

leading to the optimization problem

$$\theta^* = \arg \max_{\theta} J(\theta). \quad (8)$$

6.4 Vanilla Policy Gradient

Next, let us attempt to learn policy parameters θ with *gradient ascent* on the expected return:

$$\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} J(\theta_k). \quad (9)$$

The core identity that we will use to do this is the REINFORCE policy gradient, shown below.

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t) R(\tau) \right]. \quad (10)$$

Deriving the policy gradient. How did we get this equation? For completeness, we will give a derivation of this identity below. We will make use of a few identities.

1. The probability of a trajectory is given by

$$P(\tau \mid \theta) = \rho_0(s_0) \prod_{t=0}^T P(s_{t+1} \mid s_t, a_t) \pi_{\theta}(a_t \mid s_t). \quad (11)$$

Therefore, the log-probability of a trajectory is:

$$\log P(\tau \mid \theta) = \log \rho_0(s_0) + \sum_{t=0}^T [\log P(s_{t+1} \mid s_t, a_t) + \log \pi_{\theta}(a_t \mid s_t)]. \quad (12)$$

2. The log-derivative trick:

$$\nabla_{\theta} P = P \nabla_{\theta} \log P. \quad (13)$$

3. The environment terms are constant in θ . ρ_0 , $P(\cdot | \cdot)$ and $R(\tau)$ do not depend on the policy parameters, so

$$\nabla_{\theta} \rho_0 = \nabla_{\theta} P = \nabla_{\theta} R(\tau) = 0. \quad (14)$$

Applying the facts above:

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)] \quad (15)$$

$$= \nabla_{\theta} \sum_{\tau} P(\tau | \theta) R(\tau) \quad (16)$$

$$= \sum_{\tau} \nabla_{\theta} P(\tau | \theta) R(\tau) \quad (17)$$

$$= \sum_{\tau} P(\tau | \theta) \nabla_{\theta} \log P(\tau | \theta) R(\tau) \quad (\text{Log-derivative trick}) \quad (18)$$

$$= \mathbb{E}_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log P(\tau | \theta) R(\tau)], \quad (19)$$

and therefore, plugging in the log-probability of a trajectory and using the fact that the environment terms are constant in θ , we get the *vanilla* or REINFORCE policy gradient:

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) R(\tau) \right]. \quad (20)$$

Intuitively, this gradient will increase the log probability of every action in a trajectory that has high return, and decrease them otherwise.

Sample estimate of the gradient. Given a batch of N rollouts $\mathcal{D} = \{\tau^{(i)}\}_{i=1}^N$ collected by sampling a starting state $s_0^{(i)} \sim \rho_0(s_0)$ and then running the policy π_{θ} in the environment, we form an unbiased estimator of the gradient as

$$\hat{g} = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) R(\tau^{(i)}). \quad (21)$$

This vector is used in the gradient-ascent update $\theta \leftarrow \theta + \alpha \hat{g}$.

6.5 Policy Gradient Baselines

The main issue with vanilla policy gradient is the high variance of the gradient estimate. A common technique to mitigate this is to subtract from the reward a *baseline* function b that depends only on the state. This is a type of *control variate* [Ross, 2022]: the idea is to decrease the variance of the estimator by subtracting a term that is correlated with it, without introducing bias.

Let us define the baselined policy gradient as:

$$B = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (R(\tau) - b(s_t)) \right]. \quad (22)$$

As an example, a reasonable baseline is the on-policy value function $V^{\pi}(s) = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau) | s_t = s]$, i.e., the expected return if we start at $s_t = s$ and follow the policy π_{θ} from there. Then, the quantity $(R(\tau) - V^{\pi}(s_t))$ is, intuitively, how much better the realized trajectory is than expected.

As long as the baseline depends only on the state, the baselined policy gradient is unbiased. We can see this by rewriting the baselined policy gradient as

$$B = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) R(\tau) \right] - \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t) \right]. \quad (23)$$

Focusing on the baseline term, we see that

$$\mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t) \right] = \sum_{t=0}^T \mathbb{E}_{s_t} \left[b(s_t) \mathbb{E}_{a_t \sim \pi_\theta(\cdot | s_t)} [\nabla_\theta \log \pi_\theta(a_t | s_t)] \right]. \quad (24)$$

In general, the expectation of the score function is zero: $\mathbb{E}_{x \sim P_\theta} [\nabla_\theta \log P_\theta(x)] = 0$. Therefore, the expression in Eq. 24 is zero and

$$B = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) R(\tau) \right] - 0 = \nabla_\theta J(\pi_\theta), \quad (25)$$

so we conclude that the baselined policy gradient is unbiased. We will later run an experiment to see whether baselining improves downstream performance.

A note on policy gradient “losses.” When we implement policy gradient methods in a framework like PyTorch, we will define a so-called policy gradient loss `pg_loss` such that calling `pg_loss.backward()` will populate the gradient buffers of our model parameters with our approximate policy gradient \hat{g} . In math, it can be stated as

$$\text{pg_loss} = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \log \pi_\theta(a_t^{(i)} | s_t^{(i)}) (R(\tau^{(i)}) - b(s_t^{(i)})). \quad (26)$$

`pg_loss` is not a loss in the canonical sense—it’s not meaningful to report `pg_loss` on the train or validation set as an evaluation metric, and a good validation `pg_loss` doesn’t indicate that our model is generalizing well. The `pg_loss` is really just some scalar such that when we call `pg_loss.backward()`, the gradients we obtain through backprop are the approximate policy gradient \hat{g} .

When doing RL, you should always **log and report train and validation rewards**. These are the “meaningful” evaluation metrics and what we are attempting to optimize with policy gradient methods.

6.6 Off-Policy Policy Gradient

REINFORCE is an *on-policy* algorithm: the training data is collected by the same policy that we are optimizing. To see this, let us write out the REINFORCE algorithm:

1. Sample a batch of rollouts $\{\tau^{(i)}\}_{i=1}^N$ from the current policy π_θ .
2. Approximate the policy gradient as $\nabla_\theta J(\pi_\theta) \approx \hat{g} = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t^{(i)} | s_t^{(i)}) R(\tau^{(i)})$.
3. Update the policy parameters using the computed gradient: $\theta \leftarrow \theta + \alpha \hat{g}$.

We need to do a lot of inference to sample a new batch of rollouts, only to take just one gradient step. The behavior of an LM generally cannot change significantly in a single step, so this on-policy approach is highly inefficient.

Off-policy policy gradient. In off-policy learning, we instead have rollouts sampled from some policy other than the one we are optimizing. Off-policy variants of popular policy gradient algorithms like PPO and GRPO use rollouts from a previous version of the policy $\pi_{\theta_{\text{old}}}$ to optimize the current policy π_{θ} . The off-policy policy gradient estimate is

$$\hat{g}_{\text{off-policy}} = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \frac{\pi_{\theta}(a_t^{(i)} | s_t^{(i)})}{\pi_{\theta_{\text{old}}}(a_t^{(i)} | s_t^{(i)})} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) R(\tau^{(i)}). \quad (27)$$

This looks like an importance sampled version of the vanilla policy gradient, with reweighting terms $\frac{\pi_{\theta}(a_t^{(i)} | s_t^{(i)})}{\pi_{\theta_{\text{old}}}(a_t^{(i)} | s_t^{(i)})}$. Indeed, Eq. 27 can be derived by importance sampling and applying an approximation that is reasonable as long as π_{θ} and $\pi_{\theta_{\text{old}}}$ are not too different: see Degris et al. [2013] for more on this.

7 Group Relative Policy Optimization

Next, we will describe Group Relative Policy Optimization (GRPO), the variant of policy gradient that you will implement and experiment with for solving math problems.

7.1 GRPO Algorithm

Advantage estimation. The core idea of GRPO is to sample many outputs for each question from the policy π_{θ} and use them to compute a baseline. This is convenient because we avoid the need to learn a neural value function $V_{\phi}(s)$, which can be hard to train and is cumbersome from the systems perspective. For a question q and group outputs $\{o^{(i)}\}_{i=1}^G \sim \pi_{\theta}(\cdot | q)$, let $r^{(i)} = R(q, o^{(i)})$ be the reward for the i -th output. DeepSeekMath [Shao et al., 2024] and DeepSeek R1 [DeepSeek-AI et al., 2025] compute the group-normalized reward for the i -th output as

$$A^{(i)} = \frac{r^{(i)} - \text{mean}(r^{(1)}, r^{(2)}, \dots, r^{(G)})}{\text{std}(r^{(1)}, r^{(2)}, \dots, r^{(G)}) + \text{advantage_eps}}, \quad (28)$$

where `advantage_eps` is a small constant to prevent division by zero. Note that this *advantage* $A^{(i)}$ is the same for each token in the response, i.e., $A_t^{(i)} = A^{(i)}, \forall t \in 1, \dots, |o^{(i)}|$, so we drop the t subscript in the following.

High-level algorithm. Before we dive into the GRPO objective, let us first get an idea of the train loop by writing out the algorithm from Shao et al. [2024] in Algorithm 3.²

GRPO objective. The GRPO objective combines three ideas:

1. Off-policy policy gradient, as in Eq. 27.
2. Computing advantages $A^{(i)}$ with group normalization, as in Eq. 28.
3. A clipping mechanism, as in Proximal Policy Optimization (PPO, Schulman et al. [2017]).

The purpose of clipping is to maintain stability when taking many gradient steps on a single batch of rollouts. It works by keeping the policy π_{θ} from straying too far from the old policy.

²This is a special case of DeepSeekMath’s GRPO with a verified reward function, no KL term, and no iterative update of the reference and reward model.

Algorithm 3 Group Relative Policy Optimization (GRPO)

Input initial policy model $\pi_{\theta_{\text{init}}}$; reward function R ; task questions \mathcal{D}

```
1: policy model  $\pi_{\theta} \leftarrow \pi_{\theta_{\text{init}}}$ 
2: for step = 1, ..., n_grpo_steps do
3:   Sample a batch of questions  $\mathcal{D}_b$  from  $\mathcal{D}$ 
4:   Set the old policy model  $\pi_{\theta_{\text{old}}} \leftarrow \pi_{\theta}$ 
5:   Sample  $G$  outputs  $\{o^{(i)}\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot | q)$  for each question  $q \in \mathcal{D}_b$ 
6:   Compute rewards  $\{r^{(i)}\}_{i=1}^G$  for each sampled output  $o^{(i)}$  by running reward function  $R(q, o^{(i)})$ 
7:   Compute  $A^{(i)}$  with group normalization (Eq. 28)
8:   for train step = 1, ..., n_train_steps_per_rollout_batch do
9:     Update the policy model  $\pi_{\theta}$  by maximizing the GRPO-Clip objective (to be discussed, Eq. 29)
10:  end for
11: end for
Output  $\pi_{\theta}$ 
```

Let us first write out the full GRPO-Clip objective, and then we can build some intuition on what the clipping does:

$$J_{\text{GRPO-Clip}}(\theta) = \mathbb{E}_{q \sim \mathcal{D}, \{o^{(i)}\}_{i=1}^G \sim \pi_{\theta}(\cdot | q)} \left[\frac{1}{G} \sum_{i=1}^G \frac{1}{|o^{(i)}|} \sum_{t=1}^{|o^{(i)}|} \underbrace{\min \left(\frac{\pi_{\theta}(o_t^{(i)} | q, o_{<t}^{(i)})}{\pi_{\theta_{\text{old}}}(o_t^{(i)} | q, o_{<t}^{(i)})} A^{(i)}, \text{clip} \left(\frac{\pi_{\theta}(o_t^{(i)} | q, o_{<t}^{(i)})}{\pi_{\theta_{\text{old}}}(o_t^{(i)} | q, o_{<t}^{(i)})}, 1 - \epsilon, 1 + \epsilon \right) A^{(i)} \right)}_{\text{per-token objective}} \right]. \quad (29)$$

The hyperparameter $\epsilon > 0$ controls how much the policy can change. To see this, we can rewrite the per-token objective in a more intuitive way following Achiam [2018a,b]. Define the function

$$g(\epsilon, A^{(i)}) = \begin{cases} (1 + \epsilon)A^{(i)} & \text{if } A^{(i)} \geq 0 \\ (1 - \epsilon)A^{(i)} & \text{if } A^{(i)} < 0. \end{cases} \quad (30)$$

We can rewrite the per-token objective as

$$\text{per-token objective} = \min \left(\frac{\pi_{\theta}(o_t^{(i)} | q, o_{<t}^{(i)})}{\pi_{\theta_{\text{old}}}(o_t^{(i)} | q, o_{<t}^{(i)})} A^{(i)}, g(\epsilon, A^{(i)}) \right)$$

We can now reason by cases. When the advantage $A^{(i)}$ is positive, the per-token objective simplifies to

$$\text{per-token objective} = \min \left(\frac{\pi_{\theta}(o_t^{(i)} | q, o_{<t}^{(i)})}{\pi_{\theta_{\text{old}}}(o_t^{(i)} | q, o_{<t}^{(i)})}, 1 + \epsilon \right) A^{(i)}.$$

Since $A^{(i)} > 0$, the objective goes up if the action $o_t^{(i)}$ becomes more likely under π_{θ} , i.e., if $\pi_{\theta}(o_t^{(i)} | q, o_{<t}^{(i)})$ increases. The clipping with \min limits how much the objective can increase: once $\pi_{\theta}(o_t^{(i)} | q, o_{<t}^{(i)}) > (1 + \epsilon)\pi_{\theta_{\text{old}}}(o_t^{(i)} | q, o_{<t}^{(i)})$, this per-token objective hits its maximum value of $(1 + \epsilon)A^{(i)}$. So, the policy π_{θ} is not incentivized to go very far from the old policy $\pi_{\theta_{\text{old}}}$.

Analogously, when the advantage $A^{(i)}$ is negative, the model tries to drive down $\pi_{\theta}(o_t^{(i)} | q, o_{<t}^{(i)})$, but is not incentivized to decrease it below $(1 - \epsilon)\pi_{\theta_{\text{old}}}(o_t^{(i)} | q, o_{<t}^{(i)})$ (refer to Achiam [2018b] for the full argument).

7.2 Implementation

Now that we have a high-level understanding of the GRPO training loop and objective, we will start implementing pieces of it. Many of the pieces implemented in the SFT and EI sections will also be reused for GRPO.

Computing advantages (group-normalized rewards). First, we will implement the logic to compute advantages for each example in a rollout batch, i.e., the group-normalized rewards. We will consider two possible ways to obtain group-normalized rewards: the approach presented above in Eq. 28, and a recent simplified approach.

Dr. GRPO [Liu et al., 2025] highlights that normalizing by $\text{std}(r^{(1)}, r^{(2)}, \dots, r^{(G)})$ rewards questions in a batch with low variation in answer correctness, which may not be desirable. They propose simply removing the normalization step, computing

$$A^{(i)} = r^{(i)} - \text{mean}(r^{(1)}, r^{(2)}, \dots, r^{(G)}). \quad (31)$$

We will implement both variants and compare their performance later in the assignment.

Problem (compute_group_normalized_rewards): Group normalization (2 points)

Deliverable: Implement a method `compute_group_normalized_rewards` that calculates raw rewards for each rollout response, normalizes them within their groups, and returns both the normalized and raw rewards along with any metadata you think is useful.

The following interface is recommended:

```
def compute_group_normalized_rewards(
    reward_fn,
    rollout_responses,
    repeated_ground_truths,
    group_size,
    advantage_eps,
    normalize_by_std,
):
```

Compute rewards for each group of rollout responses, normalized by the group size.

Args:

reward_fn: Callable[[str, str], dict[str, float]] Scores the rollout responses against the ground truths, producing a dict with keys "reward", "format_reward", and "answer_reward".

rollout_responses: list[str] Rollouts from the policy. The length of this list is rollout_batch_size = n_prompts_per_rollout_batch * group_size.

repeated_ground_truths: list[str] The ground truths for the examples. The length of this list is rollout_batch_size, because the ground truth for each example is repeated group_size times.

group_size: int Number of responses per question (group).

advantage_eps: float Small constant to avoid division by zero in normalization.

normalize_by_std: bool If True, divide by the per-group standard deviation; otherwise subtract only the group mean.

Returns:

tuple[torch.Tensor, torch.Tensor, dict[str, float]].

advantages shape (rollout_batch_size,). Group-normalized rewards for each rollout response.

raw_rewards shape (rollout_batch_size,). Unnormalized rewards for each rollout response.

metadata your choice of other statistics to log (e.g. mean, std, max/min of rewards).

To test your code, implement `[adapters.run_compute_group_normalized_rewards]`. Then, run the test with `uv run pytest -k test_compute_group_normalized_rewards` and make sure your implementation passes it.

Naive policy gradient loss. Next, we will implement some methods for computing “losses”.

As a **reminder/disclaimer**, these are not really losses in the canonical sense and should not be reported as evaluation metrics. When it comes to RL, you should instead track the train and validation returns, among other metrics (cf. Section 6.5 for discussion).

We will start with the naive policy gradient loss, which simply multiplies the advantage by the log-probability of actions (and negates). With question q , response o , and response token o_t , the naive per-token policy gradient loss is

$$-A_t \cdot \log p_\theta(o_t|q, o_{<t}). \quad (32)$$

Problem (compute_naive_policy_gradient_loss): Naive policy gradient (1 point)

Deliverable: Implement a method `compute_naive_policy_gradient_loss` that computes the per-token policy-gradient loss using raw rewards or pre-computed advantages.

The following interface is recommended:

```
def compute_naive_policy_gradient_loss(
    raw_rewards_or_advantages: torch.Tensor,
    policy_log_probs: torch.Tensor,
) -> torch.Tensor:
```

Compute the policy-gradient loss at every token, where `raw_rewards_or_advantages` is either the raw reward or an already-normalized advantage.

Args:

raw_rewards_or_advantages: `torch.Tensor` Shape (batch_size, 1), scalar reward/advantage for each rollout response.

policy_log_probs: `torch.Tensor` Shape (batch_size, sequence_length), logprobs for each token.

Returns:

torch.Tensor Shape (batch_size, sequence_length), the per-token policy-gradient loss (to be aggregated across the batch and sequence dimensions in the training loop).

Implementation tips:

- Broadcast the `raw_rewards_or_advantages` over the `sequence_length` dimension.

To test your code, implement `[adapters.run_compute_naive_policy_gradient_loss]`. Then run `uv run pytest -k test_compute_naive_policy_gradient_loss` and ensure the test passes.

GRPO-Clip loss. Next, we will implement the more interesting GRPO-Clip loss.

The per-token GRPO-Clip loss is

$$-\min \left(\frac{\pi_{\theta}(o_t|q, o_{<t})}{\pi_{\theta_{\text{old}}}(o_t|q, o_{<t})} A_t, \text{clip} \left(\frac{\pi_{\theta}(o_t|q, o_{<t})}{\pi_{\theta_{\text{old}}}(o_t|q, o_{<t})}, 1 - \epsilon, 1 + \epsilon \right) A_t \right). \quad (33)$$

Problem (compute_grpo_clip_loss): GRPO-Clip loss (2 points)

Deliverable: Implement a method `compute_grpo_clip_loss` that computes the per-token GRPO-Clip loss.

The following interface is recommended:

```
def compute_grpo_clip_loss(
    advantages: torch.Tensor,
    policy_log_probs: torch.Tensor,
    old_log_probs: torch.Tensor,
    cliprange: float,
) -> tuple[torch.Tensor, dict[str, torch.Tensor]]:
```

Args:

advantages: `torch.Tensor` Shape (batch_size, 1), per-example advantages A .

policy_log_probs: `torch.Tensor` Shape (batch_size, sequence_length), per-token log probs from the policy being trained.

old_log_probs: `torch.Tensor` Shape (batch_size, sequence_length), per-token log probs from the old policy.

cliprange: `float` Clip parameter ϵ (e.g. 0.2).

Returns:

```
tuple[torch.Tensor, dict[str, torch.Tensor]].
```

loss `torch.Tensor` of shape (batch_size, sequence_length), the per-token clipped loss.

metadata dict containing whatever you want to log. We suggest logging whether each token was clipped or not, i.e., whether the clipped policy gradient loss on the RHS of the min was lower than the LHS.

Implementation tips:

- Broadcast `advantages` over `sequence_length`.

To test your code, implement `[adapters.run_compute_grpo_clip_loss]`. Then run `uv run pytest -k test_compute_grpo_clip_loss` and ensure the test passes.

Policy gradient loss wrapper. We will be running ablations comparing three different versions of policy gradient:

- no_baseline:** Naive policy gradient loss without a baseline, i.e., advantage is just the raw rewards $A = R(q, o)$.
- reinforce_with_baseline:** Naive policy gradient loss but using our group-normalized rewards as the advantage. If \bar{r} are the group-normalized rewards from `compute_group_normalized_rewards` (which may or may not be normalized by the group standard deviation), then $A = \bar{r}$.

(c) `grpo_clip`: GRPO-Clip loss.

For convenience, we will implement a wrapper that lets us easily swap between these three policy gradient losses.

Problem (`compute_policy_gradient_loss`): Policy-gradient wrapper (1 point)

Deliverable: Implement `compute_policy_gradient_loss`, a convenience wrapper that dispatches to the correct loss routine (`no_baseline`, `reinforce_with_baseline`, or `grpo_clip`) and returns both the per-token loss and any auxiliary statistics.

The following interface is recommended:

```
def compute_policy_gradient_loss(
    policy_log_probs: torch.Tensor,
    loss_type: Literal["no_baseline", "reinforce_with_baseline", "grpo_clip"],
    raw_rewards: torch.Tensor | None = None,
    advantages: torch.Tensor | None = None,
    old_log_probs: torch.Tensor | None = None,
    cliprange: float | None = None,
) -> tuple[torch.Tensor, dict[str, torch.Tensor]]:
```

Select and compute the desired policy-gradient loss.

Args:

policy_log_probs (`batch_size`, `sequence_length`), per-token log-probabilities from the policy being trained.

loss_type One of `"no_baseline"`, `"reinforce_with_baseline"`, or `"grpo_clip"`.

raw_rewards Required if `loss_type == "no_baseline"`; shape (`batch_size`, 1).

advantages Required for `"reinforce_with_baseline"` and `"grpo_clip"`; shape (`batch_size`, 1).

old_log_probs Required for `"grpo_clip"`; shape (`batch_size`, `sequence_length`).

cliprange Required for `"grpo_clip"`; scalar ϵ used for clipping.

Returns:

`tuple[torch.Tensor, dict[str, torch.Tensor]]`.

loss (`batch_size`, `sequence_length`), per-token loss.

metadata dict, statistics from the underlying routine (e.g., clip fraction for GRPO-Clip).

Implementation tips:

- Delegate to `compute_naive_policy_gradient_loss` or `compute_grpo_clip_loss`.
- Perform argument checks (see assertion pattern above).
- Aggregate any returned metadata into a single dict.

To test your code, implement `[adapters.run_compute_policy_gradient_loss]`. Then run `uv run pytest -k test_compute_policy_gradient_loss` and verify it passes.

Masked mean. Up to this point, we have the logic needed to compute advantages, log probabilities, per-token losses, and helpful statistics like per-token entropies and clip fractions. To reduce our per-token loss tensors of shape (`batch_size`, `sequence_length`) to a vector of losses (one scalar for each example), we

will compute the mean of the loss over the sequence dimension, but only over the indices corresponding to the response (i.e., the token positions for which `mask[i, j]==1`).

Normalizing by the sequence length has been canonical in most codebases for doing RL with LLMs, but it is not obvious that this is the right thing to do—you may notice, looking at our statement of the policy gradient estimate in (21), that there is no normalization factor $\frac{1}{T(i)}$. We will start with this standard primitive, often referred to as a `masked_mean`, but will later test out using the `masked_normalize` method that we implemented during SFT.

We will allow specification of the dimension over which we compute the mean, and if `dim` is `None`, we will compute the mean over all masked elements. This may be useful to obtain average per-token entropies on the response tokens, clip fractions, etc.

Problem (`masked_mean`): Masked mean (1 point)

Deliverable: Implement a method `masked_mean` that averages tensor elements while respecting a boolean mask.

The following interface is recommended:

```
def masked_mean(
    tensor: torch.Tensor,
    mask: torch.Tensor,
    dim: int | None = None,
) -> torch.Tensor:
```

Compute the mean of `tensor` along a given dimension, considering only those elements where `mask == 1`.

Args:

tensor: `torch.Tensor` The data to be averaged.

mask: `torch.Tensor` Same shape as `tensor`; positions with 1 are included in the mean.

dim: `int` | `None` Dimension over which to average. If `None`, compute the mean over all masked elements.

Returns:

torch.Tensor The masked mean; shape matches `tensor.mean(dim)` semantics.

To test your code, implement `[adapters.run_masked_mean]`. Then run `uv run pytest -k test_masked_mean` and ensure it passes.

GRPO microbatch train step. Now we are ready to implement a single microbatch train step for GRPO (recall that for a train minibatch, we iterate over many microbatches if `gradient_accumulation_steps > 1`).

Specifically, given the raw rewards or advantages and log probs, we will compute the per-token loss, use `masked_mean` to aggregate to a scalar loss per example, average over the batch dimension, adjust for gradient accumulation, and backpropagate.

Problem (`grpo_microbatch_train_step`): Microbatch train step (3 points)

Deliverable: Implement a single micro-batch update for GRPO, including policy-gradient loss, averaging with a mask, and gradient scaling.

The following interface is recommended:

```
def grpo_microbatch_train_step(
    policy_log_probs: torch.Tensor,
    response_mask: torch.Tensor,
    gradient_accumulation_steps: int,
    loss_type: Literal["no_baseline", "reinforce_with_baseline", "grpo_clip"],
    raw_rewards: torch.Tensor | None = None,
    advantages: torch.Tensor | None = None,
    old_log_probs: torch.Tensor | None = None,
    cliprange: float | None = None,
) -> tuple[torch.Tensor, dict[str, torch.Tensor]]:
```

Execute a forward-and-backward pass on a microbatch.

Args:

policy_log_probs (batch_size, sequence_length), per-token log-probabilities from the policy being trained.

response_mask (batch_size, sequence_length), 1 for response tokens, 0 for prompt/padding.

gradient_accumulation_steps Number of microbatches per optimizer step.

loss_type One of "no_baseline", "reinforce_with_baseline", "grpo_clip".

raw_rewards Needed when loss_type == "no_baseline"; shape (batch_size, 1).

advantages Needed when loss_type != "no_baseline"; shape (batch_size, 1).

old_log_probs Required for GRPO-Clip; shape (batch_size, sequence_length).

cliprange Clip parameter ϵ for GRPO-Clip.

Returns:

tuple[torch.Tensor, dict[str, torch.Tensor]].

loss scalar tensor. The microbatch loss, adjusted for gradient accumulation. We return this so we can log it.

metadata Dict with metadata from the underlying loss call, and any other statistics you might want to log.

Implementation tips:

- You should call `loss.backward()` in this function. Make sure to adjust for gradient accumulation.

To test your code, implement `[adapters.run_grpo_microbatch_train_step]`. Then run `uv run pytest -k test_grpo_microbatch_train_step` and confirm it passes.

Putting it all together: GRPO train loop. Now we will put together a complete train loop for GRPO. You should refer to the algorithm in Section 7.1 for the overall structure, using the methods we've implemented where appropriate.

Below we provide some starter hyperparameters. If you have a correct implementation, you should see reasonable results with these.

```
n_grpo_steps: int = 200
learning_rate: float = 1e-5
```



```

advantage_eps: float = 1e-6
rollout_batch_size: int = 256
group_size: int = 8
sampling_temperature: float = 1.0
sampling_min_tokens: int = 4 # As in Expiter, disallow empty string responses
sampling_max_tokens: int = 1024
epochs_per_rollout_batch: int = 1 # On-policy
train_batch_size: int = 256 # On-policy
gradient_accumulation_steps: int = 128 # microbatch size is 2, will fit on H100
gpu_memory_utilization: float = 0.85
loss_type: Literal[
    "no_baseline",
    "reinforce_with_baseline",
    "grpo_clip",
] = "reinforce_with_baseline"
use_std_normalization: bool = True
optimizer = torch.optim.AdamW(
    policy.parameters(),
    lr=learning_rate,
    weight_decay=0.0,
    betas=(0.9, 0.95),
)

```

These default hyperparameters will start you in the on-policy setting—for each rollout batch, we take a single gradient step. In terms of hyperparameters, this means that `train_batch_size` is equal to `rollout_batch_size`, and `epochs_per_rollout_batch` is equal to 1.

Here are some sanity check asserts and constants that should remove some edge cases and point you in the right direction:

```

assert train_batch_size % gradient_accumulation_steps == 0, (
    "train_batch_size must be divisible by gradient_accumulation_steps"
)
micro_train_batch_size = train_batch_size // gradient_accumulation_steps
assert rollout_batch_size % group_size == 0, (
    "rollout_batch_size must be divisible by group_size"
)
n_prompts_per_rollout_batch = rollout_batch_size // group_size
assert train_batch_size >= group_size, (
    "train_batch_size must be greater than or equal to group_size"
)
n_microbatches_per_rollout_batch = rollout_batch_size // micro_train_batch_size

```

And here are a few additional tips:

- Remember to use the `r1_zero` prompt, and direct vLLM to stop generation at the second answer tag `</answer>`, as in the previous experiments.
- We suggest using `typer` for argument parsing.
- Use gradient clipping with clip value 1.0.
- You should routinely log validation rewards (e.g., every 5 or 10 steps). You should evaluate on at least 1024 validation examples to compare hyperparameters, as CoT/RL evaluations can be noisy.
- With our implementation of the losses, GRPO-Clip should only be used when off-policy (since it requires the old log-probabilities).
- In the off-policy setting with multiple epochs of gradient updates per rollout batch, it would be wasteful to recompute the old log-probabilities for each epoch. Instead, we can compute the old log-probabilities

- once and reuse them for each epoch.
- You should not differentiate with respect to the old log-probabilities.
- You should log some or all of the following for each optimizer update:
 - The loss.
 - Gradient norm.
 - Token entropy.
 - Clip fraction, if off-policy.
 - Train rewards (total, format, and answer).
 - Anything else you think could be useful for debugging.

Problem (grpо_train_loop): GRPO train loop (5 points)

Deliverable: Implement a complete train loop for GRPO. Begin training a policy on MATH and confirm that you see validation rewards improving, along with sensible rollouts over time. Provide a plot with the validation rewards with respect to steps, and a few example rollouts over time.

8 GRPO Experiments

Now we can start experimenting with our GRPO train loop, trying out different hyperparameters and algorithm tweaks. Each experiment will take 2 GPUs, one for the vLLM instance and one for the policy.

Note on stopping runs early. if you see significant differences between hyperparameters before 200 GRPO steps (e.g., a config diverges or is clearly suboptimal), you should of course feel free to stop the experiment early, saving time and compute for later runs. The GPU hours mentioned below are a rough estimate.

Problem (grpо_learning_rate): Tune the learning rate (2 points) (6 H100 hrs)

Starting with the suggested hyperparameters above, perform a sweep over the learning rates and report the final validation answer rewards (or note divergence if the optimizer diverges).

Deliverable: Validation reward curves associated with multiple learning rates.

Deliverable: A model that achieves validation accuracy of at least 25% on MATH.

Deliverable: A brief 2 sentence discussion on any other trends you notice on other logged metrics.

For the rest of the experiments, you can use the learning rate that performed best in your sweep above.

Effect of baselines. Continuing on with the hyperparameters above (except with your tuned learning rate), we will now investigate the effect of baselining. We are in the on-policy setting, so we will compare the loss types:

- `no_baseline`
- `reinforce_with_baseline`

Note that `use_std_normalization` is `True` in the default hyperparameters.

Problem (grpо_baselines): Effect of baselining (2 points) (2 H100 hrs)

Train a policy with `reinforce_with_baseline` and with `no_baseline`.

Deliverable: Validation reward curves associated with each loss type.

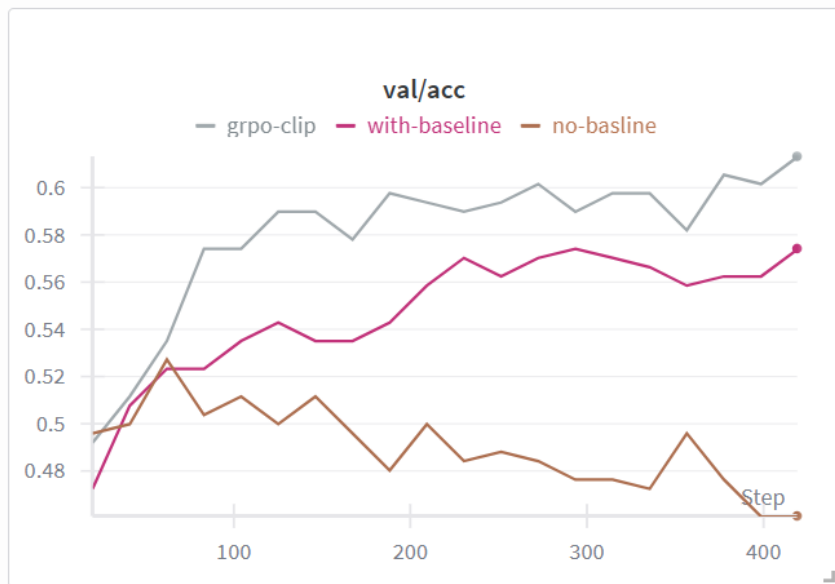
Deliverable: A brief 2 sentence discussion on any other trends you notice on other logged metrics.

For the next few experiments, you should use the best loss type found in the above experiment.

测试grpо baseline的作用

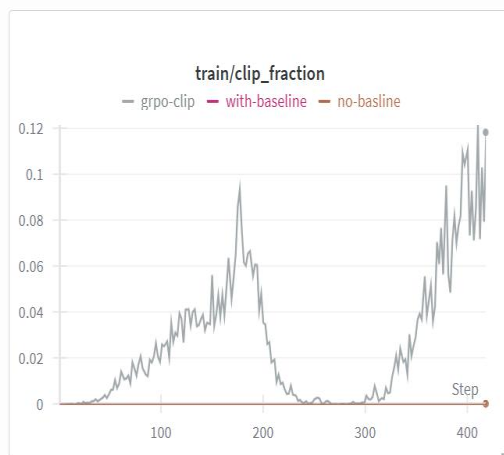
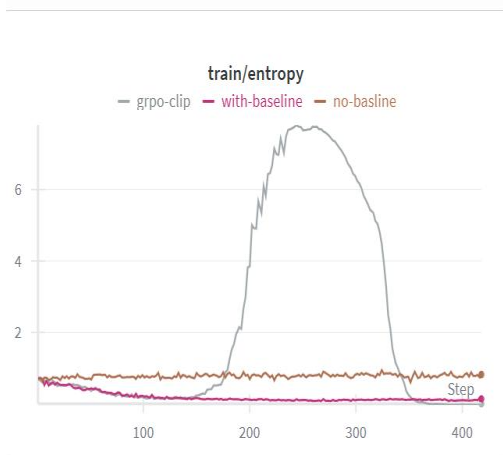
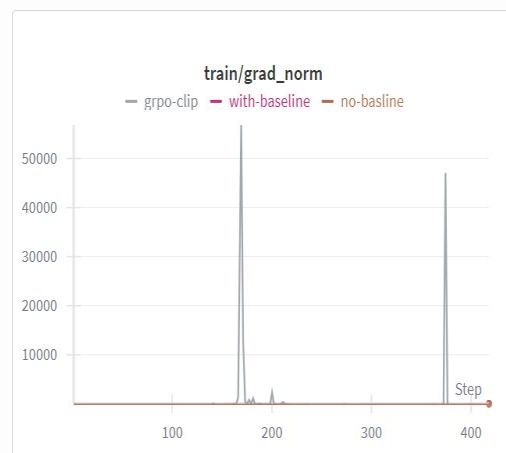
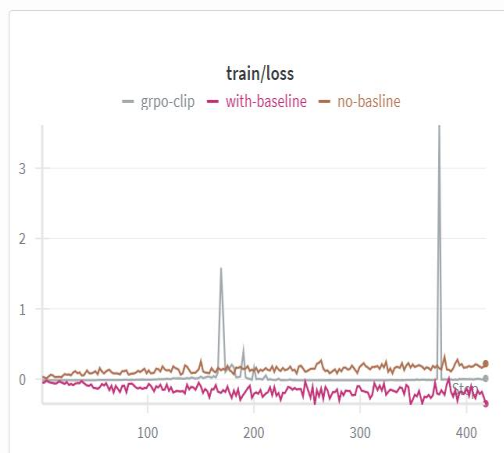
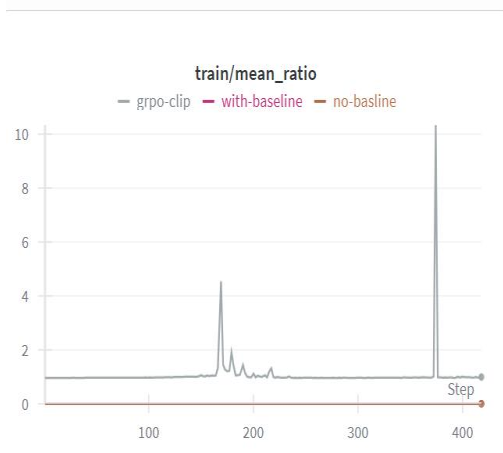
```
bash scripts/run_grpo_loss_sweep.sh
```

val 1



效果很明显，GRPO>baseline>no baseline

与你的 no-baseline 和 with-baseline 目标相比，即使在 on-policy 设置下，GRPO-Clip 依然表现出一种独特的熵变化轨迹：在训练早期，策略熵会短暂上升，随后逐渐下降并收敛。这一现象并非源于 off-policy 更新，而是来自 clipped ratio 目标本身的梯度结构。在训练初期，clipping 会限制高优势 token 概率的快速提升，使得模型更容易通过“摊平”概率分布来降低损失；随着训练推进、优势信号变得更加一致，策略开始将概率质量集中到高奖励 token 上，从而导致熵下降并实现收敛。



Length normalization. As hinted at when we were implementing `masked_mean`, it is not necessary or even correct to average losses over the sequence length. The choice of how to sum over the loss is an important hyperparameter which results in different types of credit attribution to policy actions.

Let us walk through an example from Lambert [2024] to illustrate this. Inspecting the GRPO train step, we start out by obtaining per-token policy gradient losses (ignoring clipping for a moment):

```
advantages # (batch_size, 1)
per_token_probability_ratios # (batch_size, sequence_length)
per_token_loss = -advantages * per_token_probability_ratios
```

where we have broadcasted the advantages over the sequence length. Let's compare two approaches to aggregating these per-token losses:

- The `masked_mean` we implemented, which averages over the unmasked tokens in each sequence.
- Summing over the unmasked tokens in each sequence, and dividing by a constant scalar (which our `masked_normalize` method supports with `constant_normalizer != 1.0`) [Liu et al., 2025, Yu et al., 2025].

We will consider an example where we have a batch size of 2, the first response has 4 tokens, and the second response has 7 tokens. Then, we can see how these normalization approaches affect the gradient.

```
from your_utils import masked_mean, masked_normalize

ratio = torch.tensor([
    [1, 1, 1, 1, 1, 1, 1,],
    [1, 1, 1, 1, 1, 1, 1,],
], requires_grad=True)

advs = torch.tensor([
    [2, 2, 2, 2, 2, 2, 2,],
    [2, 2, 2, 2, 2, 2, 2,],
])

masks = torch.tensor([
    # generation 1: 4 tokens
    [1, 1, 1, 1, 0, 0, 0,],
    # generation 2: 7 tokens
    [1, 1, 1, 1, 1, 1, 1,],
])

# Normalize with each approach
max_gen_len = 7
masked_mean_result = masked_mean(ratio * advs, masks, dim=1)
masked_normalize_result = masked_normalize(
    ratio * advs, masks, dim=1, constant_normalizer=max_gen_len)

print("masked_mean", masked_mean_result)
print("masked_normalize", masked_normalize_result)

# masked_mean tensor([2., 2.], grad_fn=<DivBackward0>)
# masked_normalize tensor([1.1429, 2.0000], grad_fn=<DivBackward0>)

masked_mean_result.mean().backward()
print("ratio.grad", ratio.grad)
# ratio.grad:
```

```
# tensor([[0.2500, 0.2500, 0.2500, 0.2500, 0.0000, 0.0000, 0.0000],
#         [0.1429, 0.1429, 0.1429, 0.1429, 0.1429, 0.1429, 0.1429]])
ratio.grad.zero_()

masked_normalize_result.mean().backward()
print("ratio.grad", ratio.grad)
# ratio.grad:
# tensor([[0.1429, 0.1429, 0.1429, 0.1429, 0.0000, 0.0000, 0.0000],
#         [0.1429, 0.1429, 0.1429, 0.1429, 0.1429, 0.1429, 0.1429]])
```

Problem (think_about_length_normalization): Think about length normalization (1 point)

Deliverable: Compare the two approaches (without running experiments yet). What are the pros and cons of each approach? Are there any specific settings or examples where one approach seems better?

✓ masked_normalize (除以固定常数)

优点

- 每个 token 接收相同强度的 credit
- 长序列自然带来更多总梯度，更符合：
 - token-level decision view
 - language modeling / per-token policy optimization
- 在以下场景更合理：
 - Advantage 可以被理解为“每一步都适用”
 - 希望模型为每个 action 都负责
 - GRPO / PPO-like 方法中常见 (Liu et al., Yu et al.)

缺点

- 长序列主导训练 (gradient domination)
- 如果 reward 本质是 sequence-level 的，会过度惩罚/奖励长输出
- 对长度分布非常敏感 (length bias)

优缺点对比

✓ masked_mean (按长度平均)

优点

- 每个样本 (trajectory / response) 贡献相同的总梯度
- 不会因为生成得更长就被“奖励更多更新”
- 在以下场景更合理：
 - Advantage 是 sequence-level reward (如 RLHF)
 - 奖励并不随着 token 数线性增长
 - 希望“生成多 ≠ 学得更多”

缺点

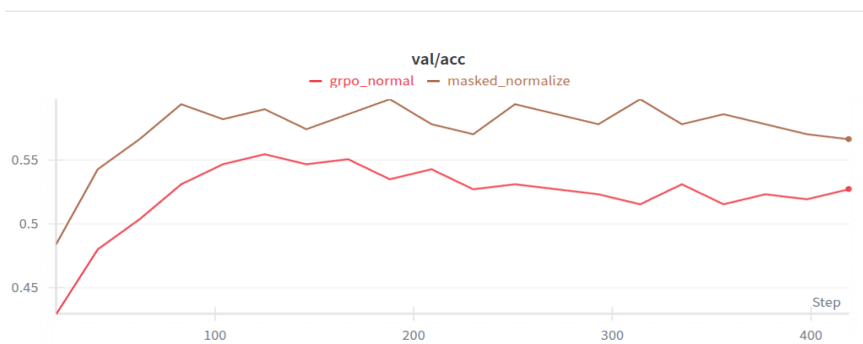
- 对长序列的单个 token 学习信号较弱
- 如果 reward 实际上与 token 数相关，会低估长输出的影响
- 会隐式鼓励模型生成更长序列 (因为梯度被平均稀释)

Now, let's compare masked_mean with masked_normalize empirically.

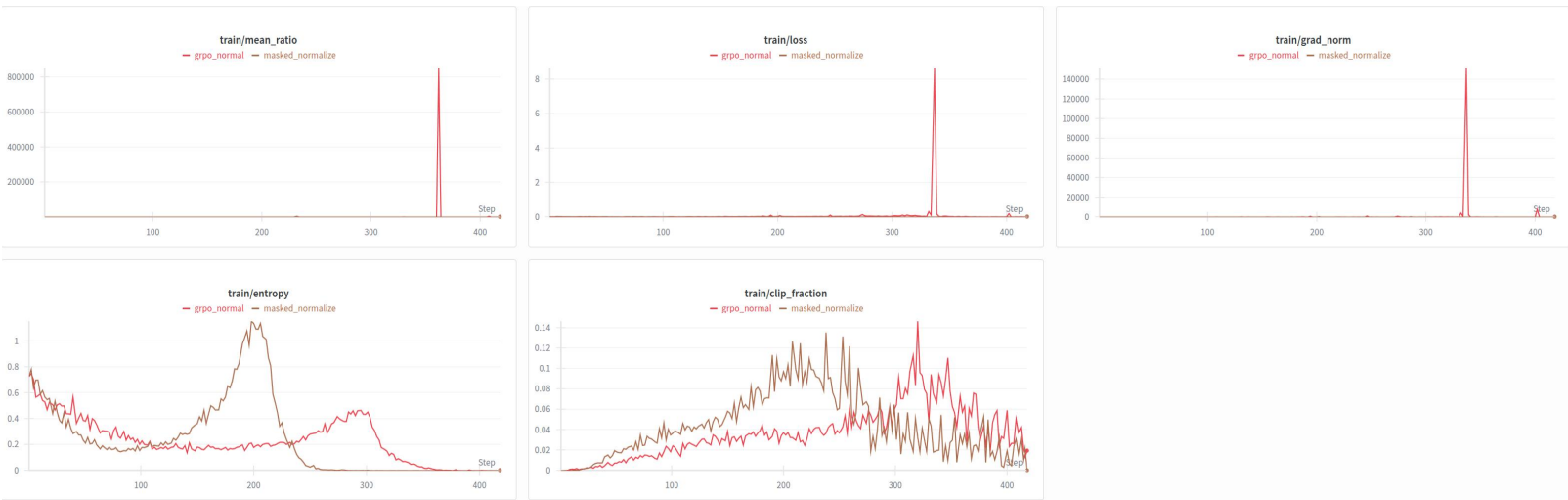
Problem (grpo_length_normalization): Effect of length normalization (2 points) (2 H100 hrs)

Deliverable: Compare normalization with masked_mean and masked_normalize with an end-to-end GRPO training run. Report the validation answer reward curves. Comment on the findings, including any other metrics that have a noticeable trend.

Hint: consider metrics related to stability, such as the gradient norm.



mask normalize明显更好一些，在GRPO中，因为是action负责而不是对于token负责。



Fix to the better performing length normalization approach for the following experiments.

Normalization with group standard deviation. Recall our standard implementation of `compute_group_normalized_rewards` (based on Shao et al. [2024], DeepSeek-AI et al. [2025]), where we normalized by the group standard deviation. Liu et al. [2025] notes that dividing by the group standard deviation could introduce unwanted biases to the training procedure: questions with lower standard deviations (e.g., too easy or too hard questions with all rewards almost all 1 or all 0) would receive higher weights during training.

Liu et al. [2025] propose removing the normalization by the standard deviation, which we have already implemented in `compute_group_normalized_rewards` and will now test.

Problem (`grpo_group_standard_deviation`): Effect of standard deviation normalization (2 points) (2 H100 hrs)

Deliverable: Compare the performance of `use_std_normalization == True` and `use_std_normalization == False`. Report the validation answer reward curves. Comment on the findings, including any other metrics that have a noticeable trend. Hint: consider metrics related to stability, such as the gradient norm.

Fix to the better performing group normalization approach for the following experiments.



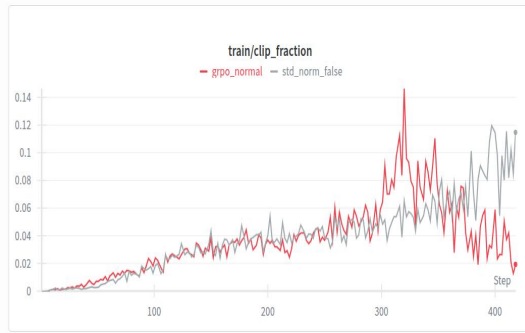
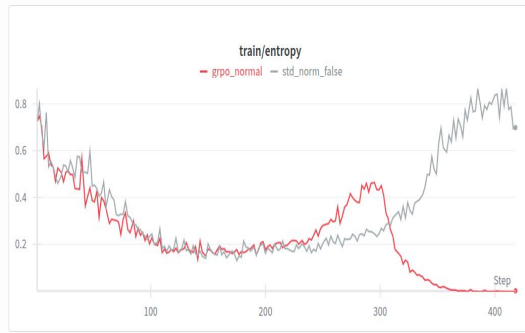
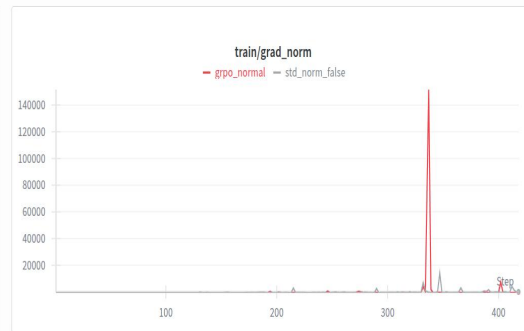
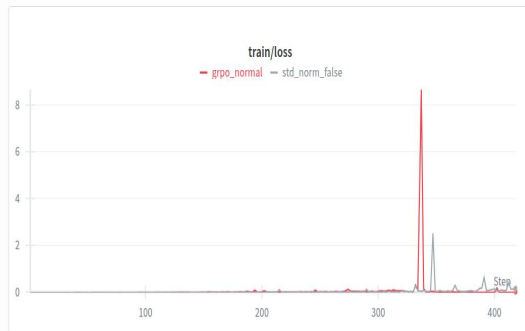
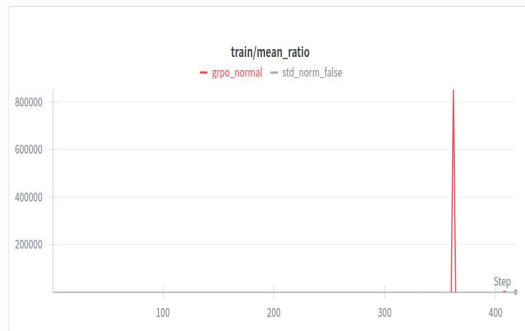
标准差归一化 (group standard deviation normalization) 的影响

我们比较了在 GRPO 训练中 是否对组内奖励按标准差进行归一化 (即 `use_std_normalization=True` 与 `False`) 对训练效果的影响。

实验结果表明, 当 启用标准差归一化 时, 训练过程整体更加不稳定。具体来说, 验证集 answer reward 曲线波动更大, 在部分训练阶段会出现明显抖动; 同时, 梯度范数 (gradient norm) 更容易出现尖峰。这是因为当某些问题的组内奖励方差较小 (例如问题非常简单或非常困难, 所有回答奖励几乎全为 1 或 0) 时, 除以接近零的标准差会放大这些样本的优势值, 从而导致过大的梯度更新。

相比之下, 关闭标准差归一化 (`use_std_normalization=False`) 后, 训练过程明显更加稳定。验证集 answer reward 提升更加平滑, 梯度范数整体更小且变化更一致, 熵的变化也更加温和。这与 Liu et al. [2025] 的分析一致, 即标准差归一化可能会无意中的对低方差问题赋予过高权重, 从而引入不必要的训练偏置。

基于上述观察, 我们在后续所有实验中 固定使用 `use_std_normalization=False`, 因为该设置在保持甚至提升验证性能的同时, 显著改善了训练稳定性。



为什么“除以 std”可能有问题？

当你用：

$$A = \frac{r - \bar{r}}{\sigma_r}$$

如果某一组问题：

- 非常简单（reward 几乎全是 1）
- 或非常难（reward 几乎全是 0）

那么：

- $\sigma_r \approx 0$
- 即使加了 `eps`，优势会被 人为放大

结果：

- 这些“信息量很低”的问题
 - 👉 在训练中反而获得更大的权重
- 导致 梯度 spikes、不稳定更新

Off-policy versus on-policy. The hyperparameters we have experimented with so far are all on-policy: we take only a single gradient step per rollout batch, and therefore we are almost exactly using the “principled” approximation \hat{g} to the policy gradient (besides the length and advantage normalization choices mentioned above).

While this approach is theoretically justified and stable, it is inefficient. Rollouts require slow generation from the policy and therefore are the dominating cost of GRPO; it seems wasteful to only take a single gradient step per rollout batch, which may be insufficient to meaningfully change the policy’s behavior.

We will now experiment with off-policy training, where we take multiple gradient steps (and even multiple epochs) per rollout batch.

Problem (grpo_off_policy): Implement off-policy GRPO

Deliverable: Implement off-policy GRPO training.

Depending on your implementation of the full GRPO train loop above, you may already have the infrastructure to do this. If not, you need to implement the following:

- You should be able to take multiple epochs of gradient steps per rollout batch, where the number of epochs and optimizer updates per rollout batch are controlled by `rollout_batch_size`, `epochs_per_rollout_batch`, and `train_batch_size`.
- Edit your main training loop to get response logprobs from the policy after each rollout batch generation phase and before the inner loop of gradient steps—these will be the `old_log_probs`. We suggest using `torch.inference_mode()`.
- You should use the “GRPO-Clip” loss type.

Now we can use the number of epochs and optimizer updates per rollout batch to control the extent to which we are off-policy.

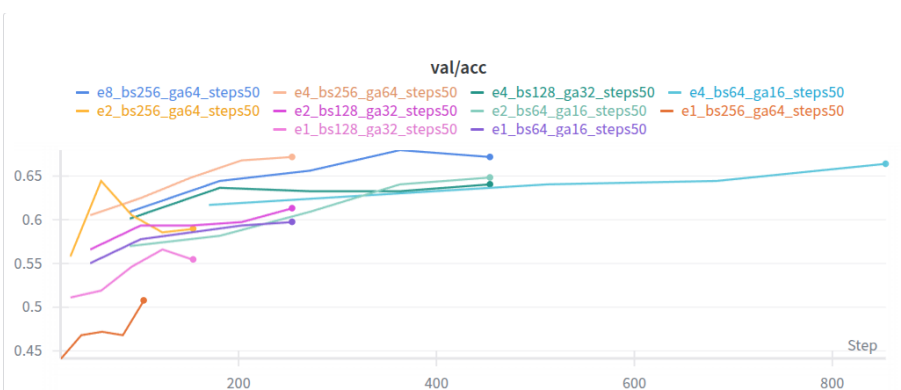
Problem (grpo_off_policy_sweep): Off-policy GRPO hyperparameter sweep (4 points) (12 H100 hrs)

Deliverable: Fixing `rollout_batch_size = 256`, choose a range over `epochs_per_rollout_batch` and `train_batch_size` to sweep over. First do a broad sweep for a limited number of GRPO steps (<50) to get a sense of the performance landscape, and then a more focused sweep for a larger number of GRPO steps (200). Provide a brief experiment log explaining the ranges you chose.

Compare to your on-policy run with `epochs_per_rollout_batch = 1` and `train_batch_size = 256`, reporting plots with respect to number of validation steps as well as with respect to wall-clock time.

Report the validation answer reward curves. Comment on the findings, including any other metrics that have a noticeable trend such as entropy and response length. Compare the entropy of the model’s responses over training to what you observed in the EI experiment.

Hint: you will need to change `gradient_accumulation_steps` to keep memory usage constant.



测试off-policy

scripts/run_grpo_offpolicy_coarse_sweep.sh



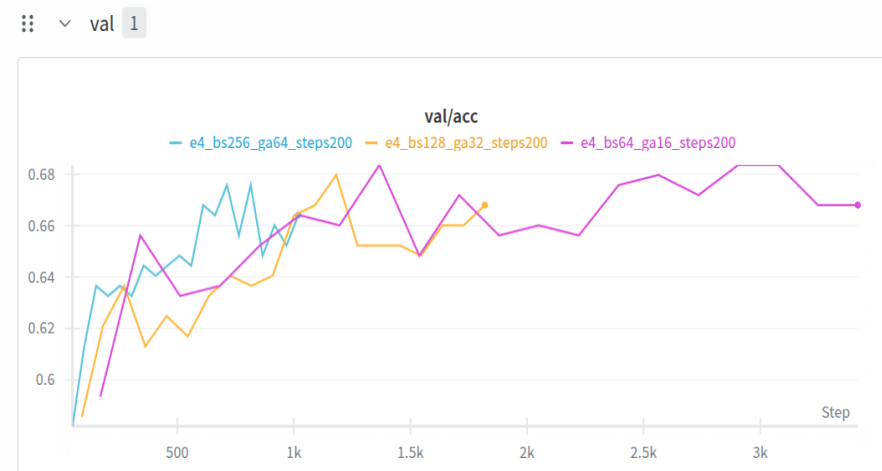
从前50看来，效果最好的是e4_bs256_ga64，：采用 Epochs=4 的 Off-Policy 配置取得了最好的权衡。虽然它在每一步（per step）的方差略高于 On-Policy 基线，但在墙钟时间（wall-clock time）上达到目标 25% 准确率的速度大约快了 2–3 倍。

指标（Metrics）：

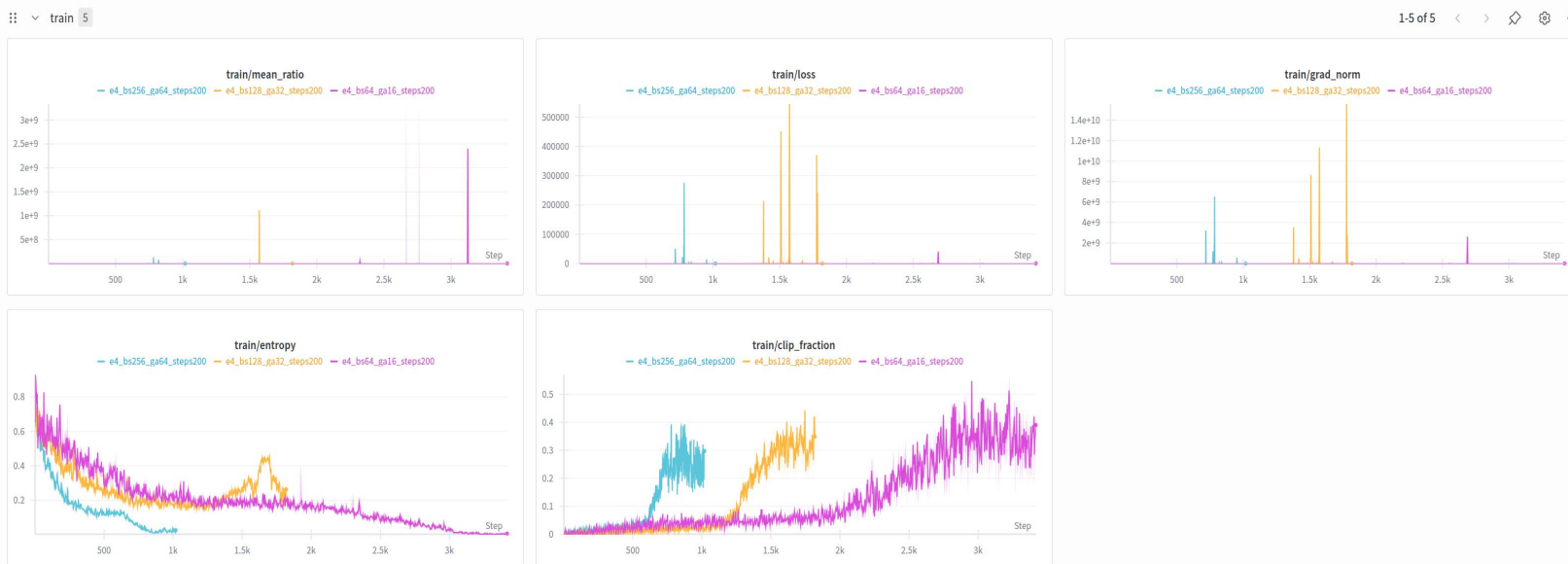
模型的回复长度稳步增长，说明它学会了进行更深入的推理，以解决更复杂的问题。

熵（Entropy）：

与 Expert Iteration 不同（其熵会很快塌缩，表现为记忆化/死记硬背），GRPO 能在更长时间内维持较高的熵，表明模型在优化策略的同时仍保留了一定的探索能力。不过，当 Epochs > 4 时，熵会过早塌缩，并且与验证性能下降相关（可能对 rollout batch 发生了过拟合）。



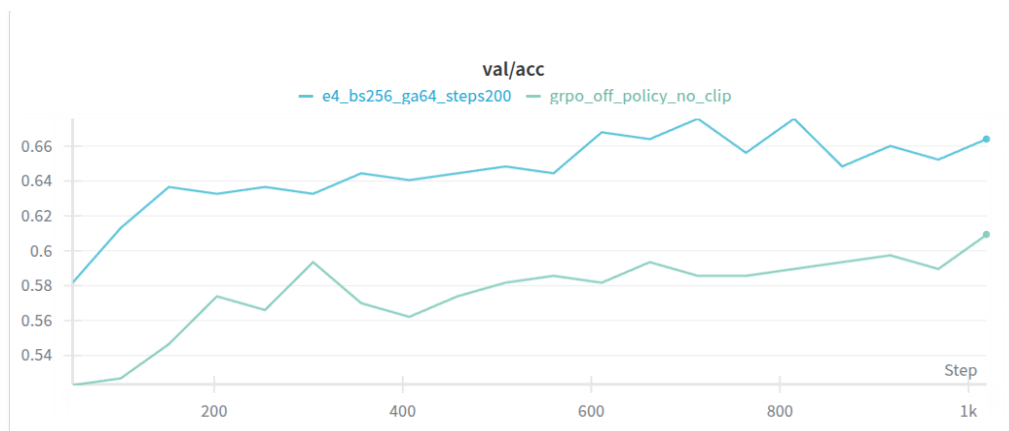
最终成绩最好的还是e4_bs64_ga16
相对于其他的他的entropy明显小一些
相比于EI 一直高的entropy，RL方法的entropy有明显的逐渐下降



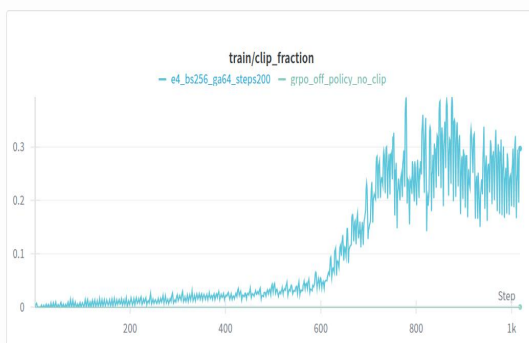
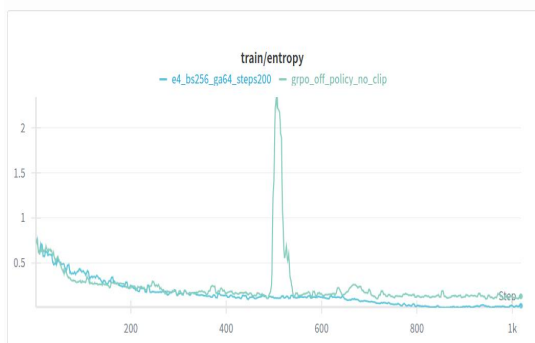
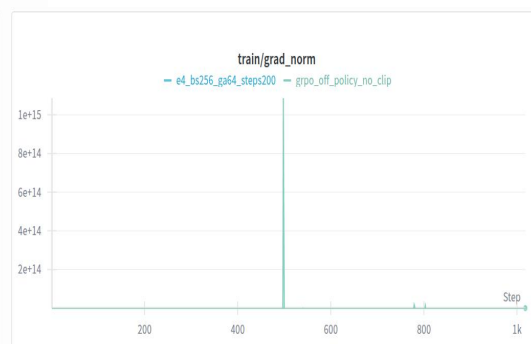
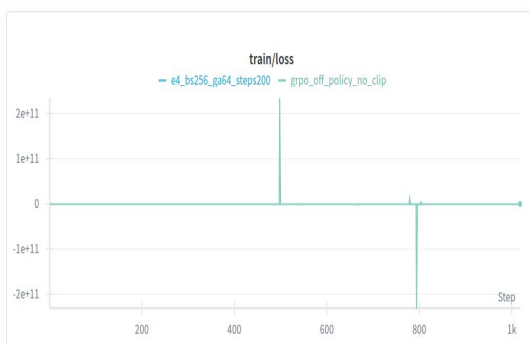
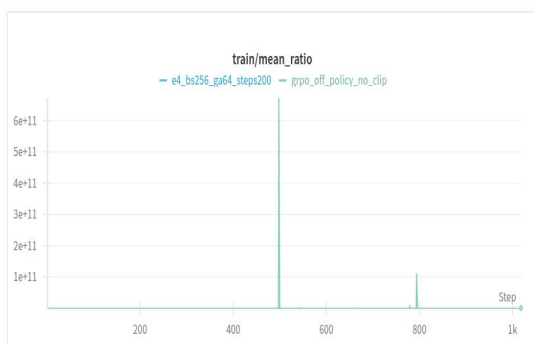
Ablating clipping in the off-policy setting. Recall that the purpose of clipping in GRPO-Clip is to prevent the policy from moving too far away from the old policy when taking many gradient steps on a single rollout batch. Next, we will ablate clipping in the off-policy setting to test to what extent it is actually necessary. In other words, we will use per-token loss

$$-\frac{\pi_{\theta}(o_t|q, o_{<t})}{\pi_{\theta_{\text{old}}}(o_t|q, o_{<t})}A_t. \quad (34)$$

Deliverable: Implement the unclipped per-token loss as a new loss type **"GRPO-No-Clip"**. Take your best performing off-policy hyperparameters from the previous problem and run the unclipped version of the loss. Report the validation answer reward curves. Comment on the findings compared to your GRPO-Clip run, including any other metrics that have a noticeable trend such as entropy, response length, and gradient norm.



明显no clip的效果更差，并且存在更高的entropy，类似于中间阶段做了大的偏移



Effect of prompt. As a last ablation, we'll investigate a surprising phenomenon: the prompt used during RL can have a dramatic effect on the performance of the model, depending on how the model was pretrained.

Instead of using the R1-Zero prompt at `cs336_alignment/prompts/r1_zero.prompt`, we will instead use an extremely simple prompt at `cs336_alignment/prompts/question_only.prompt`:

```
{question}
```

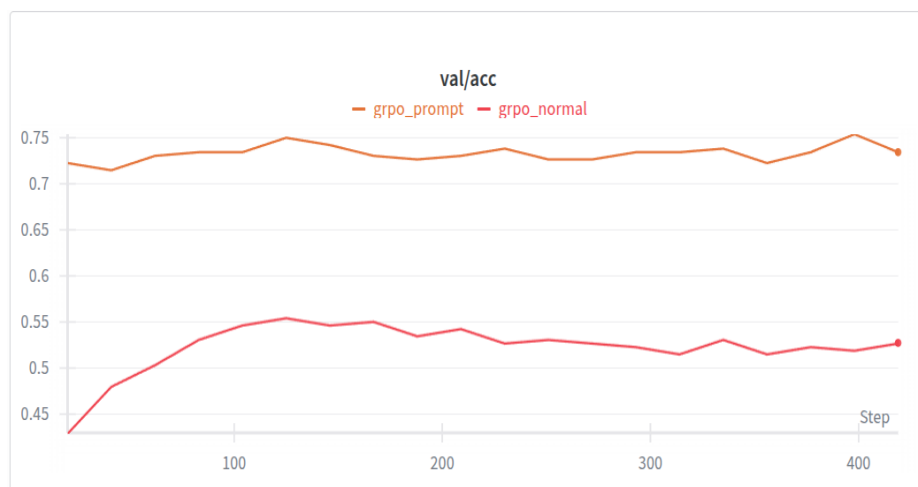
You will use this prompt for both training and validation, and will change your reward function (used both in training and validation) to the `question_only_reward_fn` located in `cs336_alignment/drgrepo_grader.py`.

Problem (grpo_prompt_ablation): Prompt ablation (2 points) (2 H100 hrs)

Deliverable: Report the validation answer reward curves for both the R1-Zero prompt and the question-only prompt. How do metrics compare, including any other metrics that have a noticeable trend such as entropy, response length, and gradient norm? Try to explain your findings.

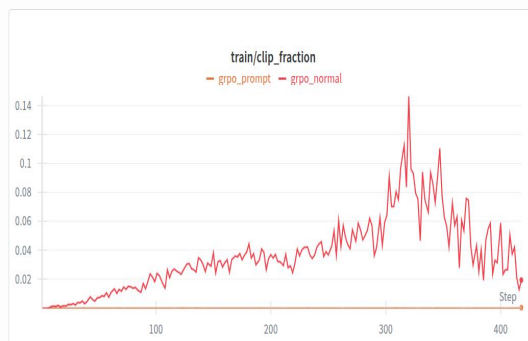
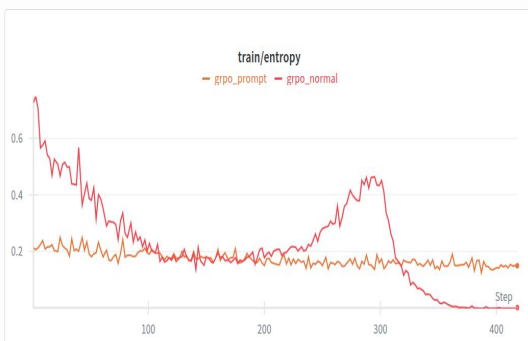
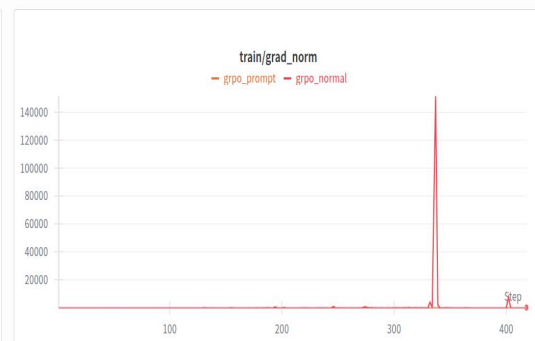
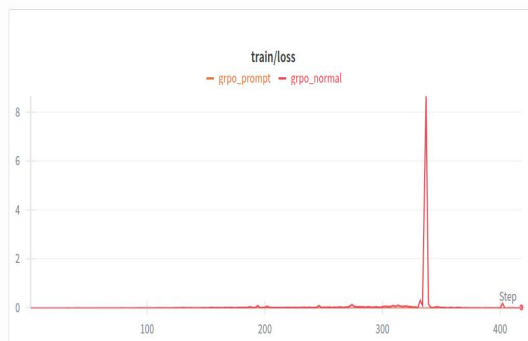
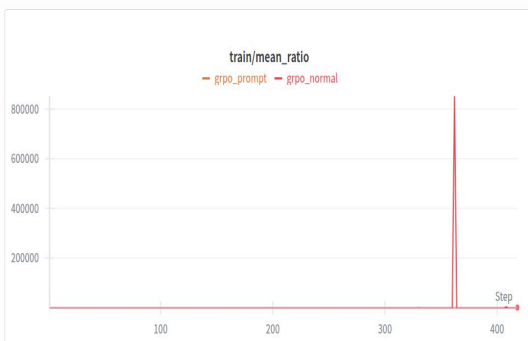
question only 其 answer reward 起点更高，并且在整个训练过程中持续提升，而不像 R1-Zero prompt 那样在快速的初期提升后很快进入平台期。

此外，使用问题 prompt 的训练过程也显得更加稳定：梯度范数明显更小，熵下降得非常缓慢（而不是出现剧烈的飙升或塌缩），损失在整个训练过程中也保持相对稳定。



train 5

1-5 of 5



9 Leaderboard: GRPO on MATH

As the last part of the (mandatory) assignment, you will experiment with approaches to obtain the highest validation rewards possible within 4 hours of training on 2 H100 GPUs.

Model. We will continue using the Qwen 2.5 Math 1.5B Base model.

Dataset. We will continue using the MATH train and validation dataset available on the cluster at `/data/a5-alignment/MATH/train.jsonl` and `/data/a5-alignment/MATH/validation.jsonl`. You are not allowed to use any other data or do SFT on reasoning chains from stronger models, etc. You must report validation accuracy on the entire validation set (all 5K examples), using the sampling hyperparameters given above (temperature 1.0, max tokens 1024). You are allowed to filter the train set, or design a curriculum over the data, as you desire. You must use the R1-Zero prompt for validation, and during validation, you must use exactly the `r1_zero_reward_fn` reward function provided in the starter code (you are allowed to develop another reward function for use during training if you wish).

Algorithm. You are free to tune hyperparameters or change the training algorithm entirely, as long as you do not use any extraneous data or another model (you are free to use more copies of the model if you want).

Systems optimizations. You might observe that in our simple GRPO implementation above, at least one GPU is always idle. You will likely find notable improvements by improving the systems characteristics of our pipeline. For example, you might consider lower precision for rollouts or training, `torch.compile`, and other systems optimizations. You are definitely not constrained to placing vLLM on a single device and the train policy on another device, and are encouraged to think of better ways to parallelize.

Ideas. For some ideas on possible improvements, see the following repos:

- `veRL`
- `trl`
- `torchtune`
- `oat`

On KL divergence. We also note that in the above experiments, we did not include a KL divergence term with respect to some reference model (usually this is a frozen SFT or pretrained checkpoint). In our experiments and others from the literature [Liu et al., 2025, NTT123, 2025], we found that omitting the KL term had no impact on performance while saving GPU memory (no need to store a reference model). However, many GRPO repos include it by default and you are encouraged to experiment with KL or other forms of regularization, **as long as you use Qwen 2.5 Math 1.5B Base or some model obtained through your algorithm for it.**

Problem (leaderboard): Leaderboard (16 points) (16 H100 hrs)

Deliverable: Report a validation accuracy obtained within 4 hours of training on 2 H100 GPUs and a screenshot of your validation accuracy with respect to wall-clock time, where the x-axis ends at ≤ 4 hours. As a reminder, we place the following constraints on your evaluation:

1. Your validation accuracy should be the average accuracy over the entire MATH validation set (all 5K examples).
2. You must use the R1-Zero prompt at validation time.
3. You must use temperature 1.0 and max tokens 1024 with vLLM for evaluation.
4. You must calculate validation accuracy by averaging the answer rewards produced by the `r1_zero_reward_fn` reward function provided in the starter code.

10 Epilogue

Congratulations on finishing the last assignment of the class! You should be proud of your hard work. We hope you enjoyed learning the foundations underlying modern language models by building their main components from scratch.

References

Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021. URL <https://arxiv.org/abs/2103.03874>.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanbiao Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.

Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Augustus Odena. Show your work: Scratchpads for intermediate computation with language models, 2021. URL <https://arxiv.org/abs/2112.00114>.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL <https://arxiv.org/abs/2201.11903>.

Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. Star: Bootstrapping reasoning with reasoning, 2022. URL <https://arxiv.org/abs/2203.14465>.

Thomas Anthony, Zheng Tian, and David Barber. Thinking fast and slow with deep learning and tree search, 2017. URL <https://arxiv.org/abs/1705.08439>.

OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong,

- Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimplouras, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O’Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiye Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. OpenAI o1 system card, 2024. URL <https://arxiv.org/abs/2412.16720>.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, Haiqing Guo, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haotian Yao, Haotian Zhao, Haoyu Lu, Haoze Li, Haozhen Yu, Hongcheng Gao, Huabin Zheng, Huan Yuan, Jia Chen, Jianhang Guo, Jianlin Su, Jianzhou Wang, Jie Zhao, Jin Zhang, Jingyuan Liu, Junjie Yan, Junyan Wu, Lidong Shi, Ling Ye, Longhui Yu, Mengnan Dong, Neo Zhang, Ningchen Ma, Qiwei Pan, Qucheng Gong, Shaowei Liu, Shengling Ma, Shupeng Wei, Sihan Cao, Siying Huang, Tao Jiang, Weihao Gao, Weimin Xiong, Weiran He, Weixiao Huang, Wenhao Wu, Wenyang He, Xianghui Wei, Xianqing Jia, Xingzhe Wu, Xinran Xu, Xinxing Zu, Xinyu Zhou, Xuehai Pan, Y. Charles, Yang Li, Yangyang Hu, Yangyang Liu, Yanru Chen, Yejie Wang, Yibo Liu, Yidao Qin, Yifeng Liu, Ying Yang, Yiping Bao, Yulun Du, Yuxin Wu, Yuzhi Wang, Zaida Zhou, Zhaoji Wang, Zhaowei Li, Zhen Zhu, Zheng Zhang, Zhexu Wang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Ziyao Xu, and Zonghan Yang. Kimi k1.5: Scaling reinforcement learning with llms, 2025. URL <https://arxiv.org/abs/2501.12599>.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999. URL https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf.
- Hugging Face. Open r1: A fully open reproduction of deepseek-r1, January 2025. URL <https://github.com/huggingface/open-r1>.
- Weihao Zeng, Yuzhen Huang, Qian Liu, Wei Liu, Keqing He, Zejun Ma, and Junxian He. Simplerl-zoo:

- Investigating and taming zero reinforcement learning for open base models in the wild, 2025. URL <https://arxiv.org/abs/2503.18892>.
- Jiayi Pan, Junjie Zhang, Xingyao Wang, Lifan Yuan, Hao Peng, and Alane Suhr. Tinyzero. <https://github.com/Jiayi-Pan/TinyZero>, 2025. Accessed: 2025-01-24.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, Keming Lu, Mingfeng Xue, Runji Lin, Tianyu Liu, Xingzhang Ren, and Zhenru Zhang. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement, 2024. URL <https://arxiv.org/abs/2409.12122>.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021a. URL <https://arxiv.org/abs/2110.14168>.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training, 2025. URL <https://arxiv.org/abs/2411.15124>.
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. Understanding rl-zero-like training: A critical perspective, 2025. URL <https://arxiv.org/abs/2503.20783>.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention, 2023. arXiv:2309.06180.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021b. arXiv:2110.14168.
- David Dohan, Winnie Xu, Aitor Lewkowycz, Jacob Austin, David Bieber, Raphael Gontijo Lopes, Yuhuai Wu, Henryk Michalewski, Rif A. Saurous, Jascha Sohl-dickstein, Kevin Murphy, and Charles Sutton. Language model cascades, 2022. URL <https://arxiv.org/abs/2207.10342>.
- Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, Wolfgang Macherey, Arnaud Doucet, Orhan Firat, and Nando de Freitas. Reinforced self-training (rest) for language modeling, 2023. URL <https://arxiv.org/abs/2308.08998>.
- Joshua Achiam. Spinning up in deep reinforcement learning. 2018a.
- Nathan Lambert. Reinforcement learning from human feedback, 2024. URL <https://rlhfbook.com>.
- Sheldon M Ross. *Simulation*. academic press, 2022.
- Thomas Degris, Martha White, and Richard S. Sutton. Off-policy actor-critic, 2013. URL <https://arxiv.org/abs/1205.4839>.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. URL <https://arxiv.org/abs/2402.03300>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.

Joshua Achiam. Simplified ppo-clip objective, 2018b. URL <https://drive.google.com/file/d/1PDzn9RPvaXjJFZkGeapMHbHGjWW20Ey/view>.

Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Weinan Dai, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. Dapo: An open-source llm reinforcement learning system at scale, 2025. URL <https://arxiv.org/abs/2503.14476>.

NTT123. Grpo-zero. <https://github.com/policy-gradient/GRPO-Zero>, 2025. Accessed: 2025-05-22.