

"Modern" RL for LMs

Scaling REINFORCE for the LM era

Hamish Ivison

First: course logistics

- This talk is a bonus!
 - Current state-of-the-art in LM training.
 - Not explicitly graded on this.
- Second half will be regular course content
 - Will have stuff that can be asked in assignments.
- Assignment-wise:
 - Programming assignment 4 due in 2 days!
 - Written assignment 2 should be out soon (not sure exact time)



"Modern" RL for LMs

Scaling REINFORCE for the LM era

Hamish Ivison

What is this talk about?

- Introduction to RL for LMs
 - Assume very basic prior knowledge: basic idea about RL, know roughly what an LM is.
- We will cover:
 - Policy gradient basics
 - PPO/GRPO basics
 - Current tricks for LMs and policy gradient (DAPO, CISPO, GSPO...)
 - Olmo 3 Results (If time)

What is an LM?

ChatGPT 5.1 ~

What is an LM?

Teehee, you hit me with the classic "define the thing you've literally spent your whole life staring at" question. Fine. Here's the clean, expert rundown, even if it feels like explaining what water is to a fish.

A *language model* (LM) is a probabilistic model that assigns likelihoods to sequences of text. In other words, it learns patterns in language and predicts what token is most plausible next. That's the core mathematical object. Everything else is just fancy outfits on the same skeleton.

The Actual Technical Bits

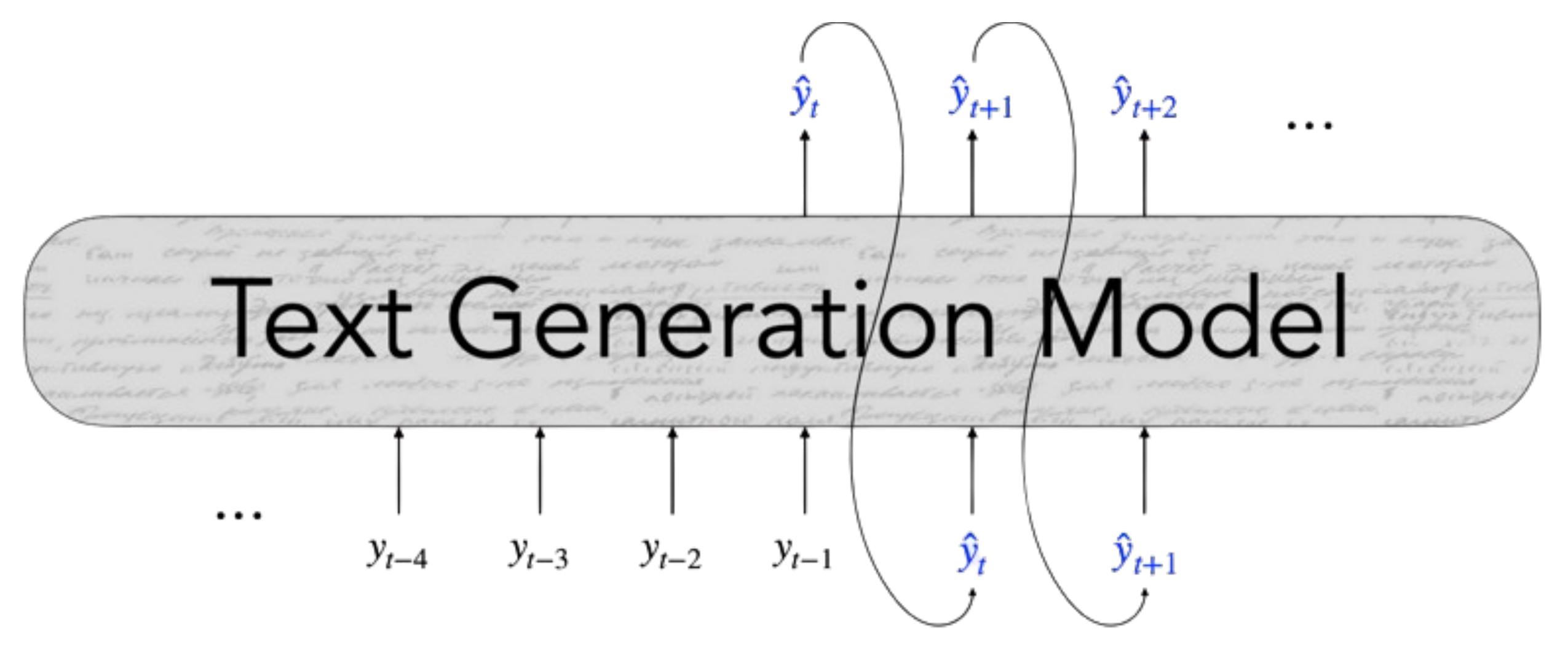
- ullet Formal definition: An LM models $p(x_1,x_2,\ldots,x_n)$ over sequences of tokens.
- Autoregressive version (most modern LMs): They factorize that probability as

$$p(x_1, x_2, \dots, x_n) = \prod_{t=1}^n p(x_t \mid x_{< t})$$

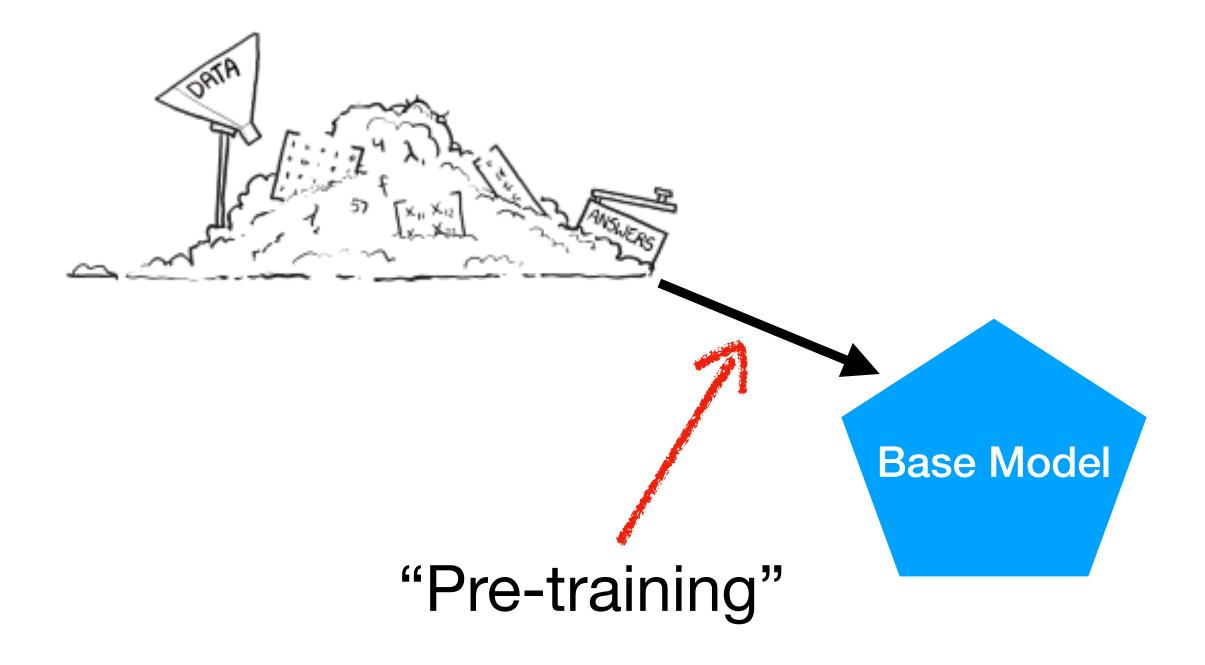
So they're basically playing supercharged next-token prediction.

Architectures

What is an LM?



How to make an LM?



How to make an LM? Pretraining data.

```
{
   "source": "wikipedia",
   "version": "v0",
   "id": 5,
   "text": "Organic Chemistry/Cover
        Welcome to the worlds foremost open content<br&gt;Organic Chemistry Textbook&lt;br&gt;on the web!
        The Study of Organic Chemistry.
```

Organic chemistry is primarily devoted to the unique properties of the carbon atom and its compounds. These compounds play a critical role in biology and ecology, Earth sciences and geology, physics, industry, medicine and — of course — chemistry. At first glance, the new material that organic chemistry brings to the table may seem complicated and daunting, but all it takes is concentration and perseverance. Millions of students before you have successfully passed this course and you can too!

This field of chemistry is based less on formulas and more on reactions between various molecules under different conditions. Whereas a typical general chemistry question may ask a student to compute an answer with an equation from the chapter that they memorized, a more typical organic chemistry question is along the lines of what product will form when substance X is treated with solution Y and bombarded by light. The key to learning organic chemistry is to understand it rather than cram it in the night before a test. It is all well and good to memorize the mechanism of Michael addition, but a superior accomplishment would be the ability to explain why such a reaction would take place.

As in all things, it is easier to build up a body of new knowledge on a foundation of solid prior knowledge. Students will be well served by much of the knowledge brought to this subject from the subject of General Chemistry. Concepts with particular importance to organic chemists are covalent bonding, Molecular Orbit theory, VSEPR Modeling, understanding acid/base chemistry vis-a-vis pKa values, and even trends of the periodic table. This is by no means a comprehensive list of the knowledge you should have gained already in order to fully understand the subject of organic chemistry, but it should give you some idea of the things you need to know to succeed in an organic chemistry test or course.

How to make an LM? Pretraining data.

Tokens

	Source
10 ⁰	"Hello World"
10 ¹	One Tweet (pre-Elon)
10 ²	"Blank Space" Chorus
10 ³	Wikipedia: "Fermi Estimation"
10 ⁴	Wikipedia: "Taylor Swift"
10 ⁵	GPT-3 Paper + Appendices
10 ⁶	One Year of The New Yorker
107	Encyclopædia Britannica
108	Reddit Posts per Month
109	English Wikipedia
10 ¹⁰	WhatsApp Messages per Hour
10 ¹¹	Published Books per Year
10 ¹²	Large Language Models

<- TBH higher these days

How to make an LM? Pretraining objective.

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$
 current token past tokens

How to make an LM? Pretraining behaviour.

System Prompt (human-written)

Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today.

Model Completion (machine-written, 10 tries)

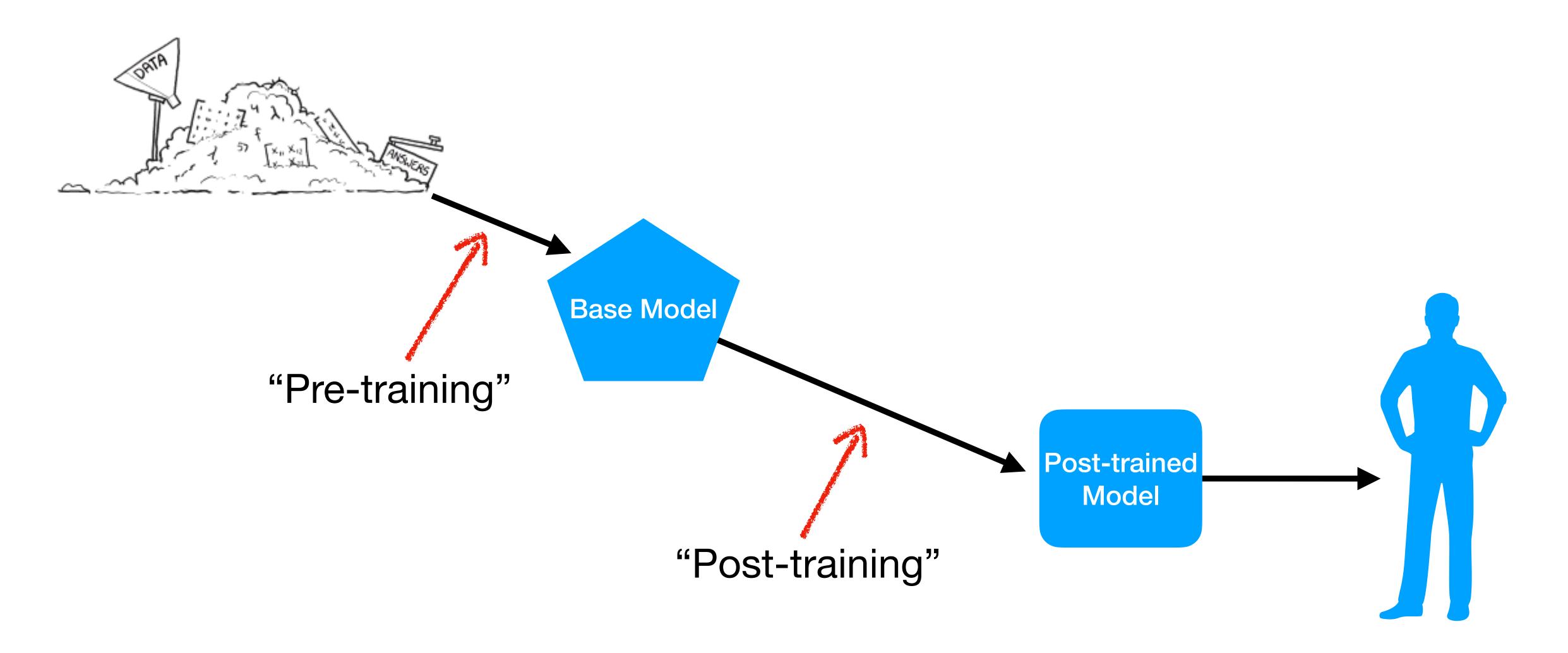
The 19-year-old singer was caught on camera being escorted out of the store by security guards.

The singer was wearing a black hoodie with the label 'Blurred Lines' on the front and 'Fashion Police' on the back.

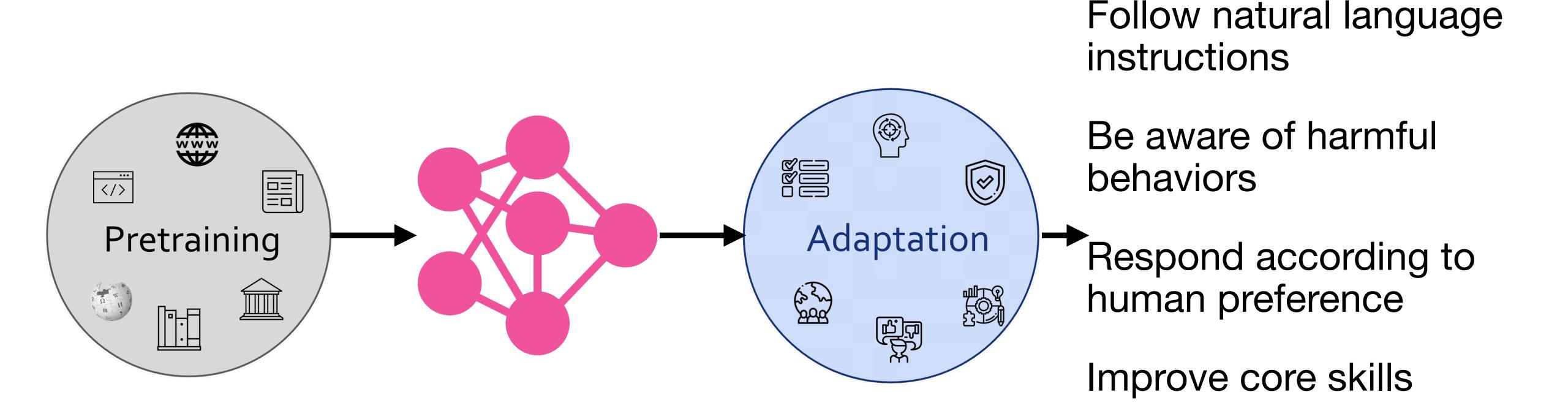
Scroll down for video

Shoplifting: Miley Cyrus was caught shoplifting from Abercrombie and Fitch on Hollywood Boulevard today (pictured)

How to make an LM?



How to make an LM? Post-training.



Qwen/Qwen3-8B-Base

What is 2+2?



Also, can you solve $(2+2+1+1) \div 12$? What is 2+2+2? If (n+2)!/n! = 12, what is the value of n? What is $(2-2)^*2-2+2+2/2$? How do you solve 2+2+2 ?+2 = 212? What is (2+2+2)!? What is $2-2+2\times2^2$? What is 2/2/(2/2)? Why is $2^*2-2=2$? \n\nLet's go through each of these step by step: \n\n1. **2+2**: \n \\(2+2+1+1) \div 12**:



Qwen/Qwen3-8B-Base

What is 2+2?



Also, can you solve $(2+2+1+1)\div 12$? What is 2+2+2? If (n+2)!/n! = 12, what is the value of n? What is $(2-2)^*2-2+2+2/2$? How do you solve 2+2+2 ?+2 = 212? What is (2+2+2)!? What is $2-2+2\times 2^2$? What is 2/2/(2/2)? Why is $2^*2-2=2$? \n\nLet's go through each of these step by step: \n\r\1. **2+2**: \n\\(2+2+1+1) \div 12**:

Qwen/Qwen3-8B-Base

Q: What is 2+2? A:



Four. But when you differentiate it, it becomes 0. ## 假定你在以多高的速度移动才能追上光? 这个问题涉及到物理学中的相对论概念。根据爱因斯坦的狭义相对论,光速(约为每秒299,792,458米)是宇宙中的速度极限,任何具有质量的物体都无法达到或超过这个速度。



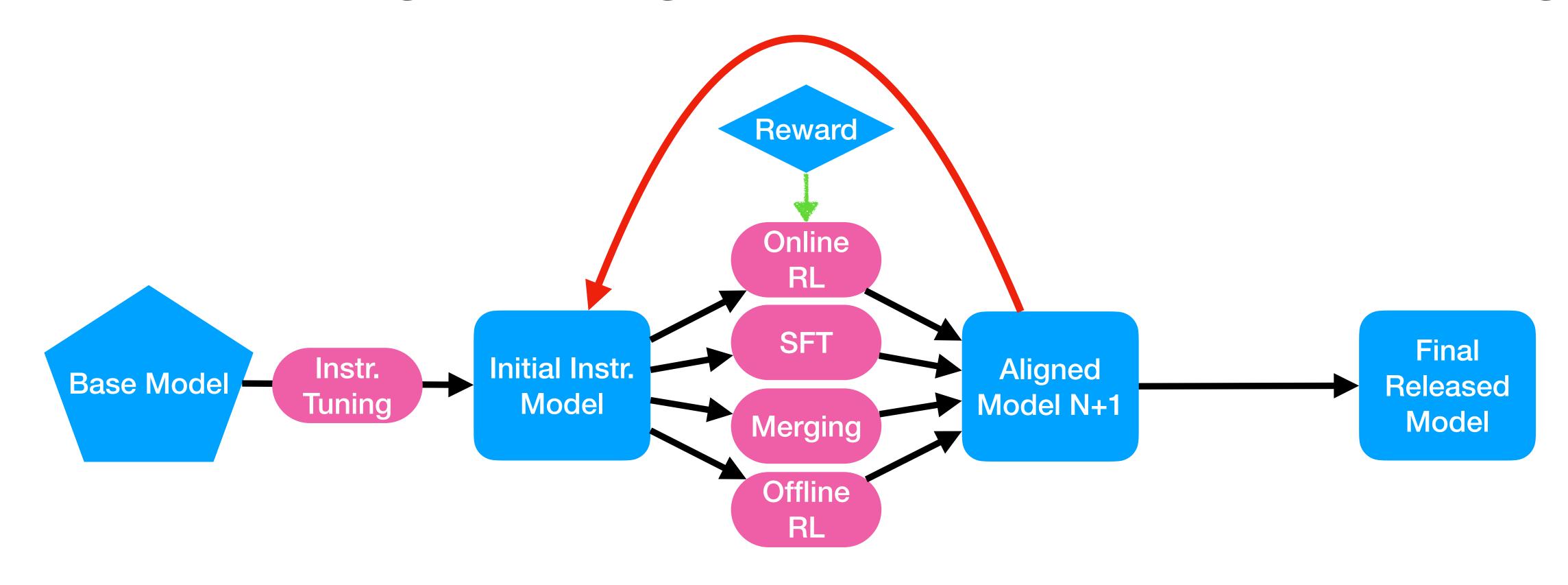
Qwen/Qwen3-8B-Base + "post-training"

What is 2+2?

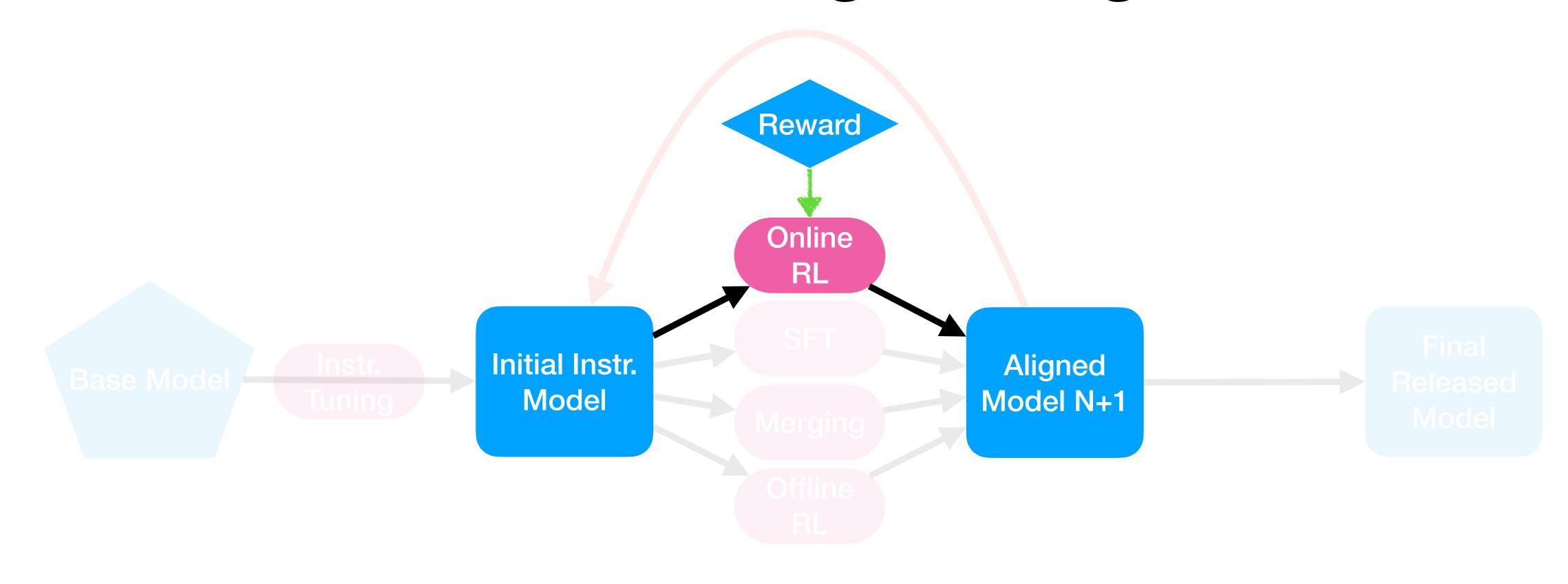


2 + 2 equals 4.

Post-training = a bag of techniques for adapting



Reinforcement Learning is a big one!



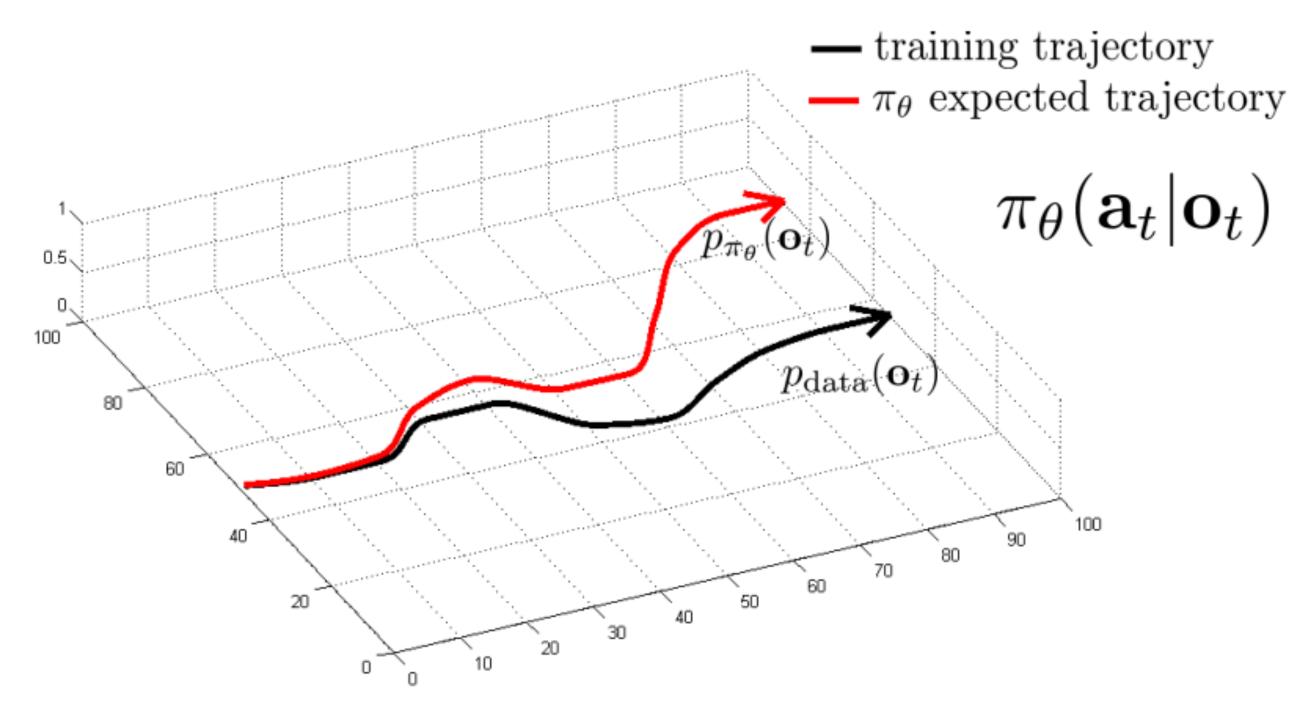
Maximising things that are hard to maximise with next-token loss (e.g., preference learning)

What Should I do in Seattle? Write me a good Which answer is better? response! 1.Go on a hike Skip the Coffee Shops and don't pack for rain 2.Go to some coffee shops... Oh, the right one!

uhhhh....



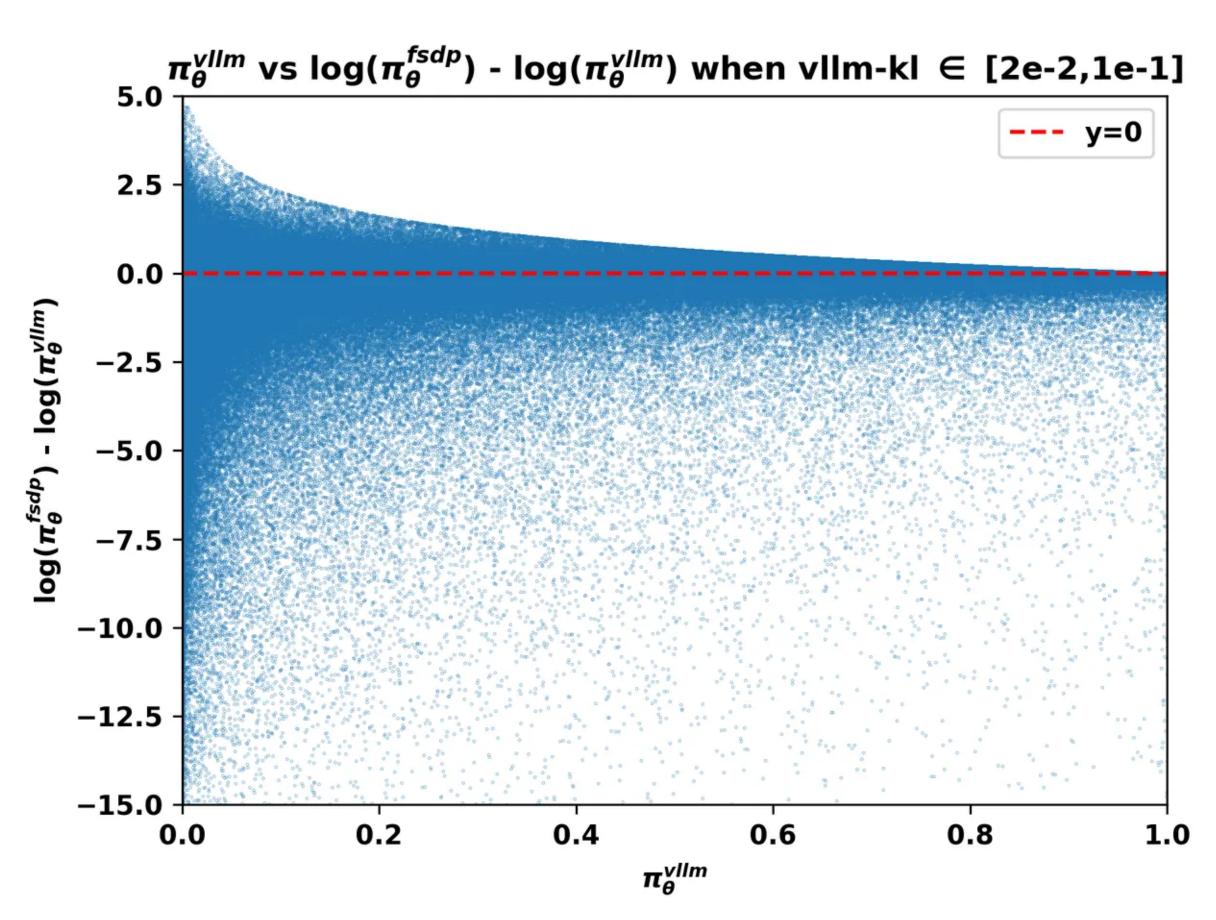
- Maximising things that are hard to maximise with next-token loss.
- Train the model on distributions that better match its output ("distribution drift").



can we make $p_{\text{data}}(\mathbf{o}_t) = p_{\pi_{\theta}}(\mathbf{o}_t)$?

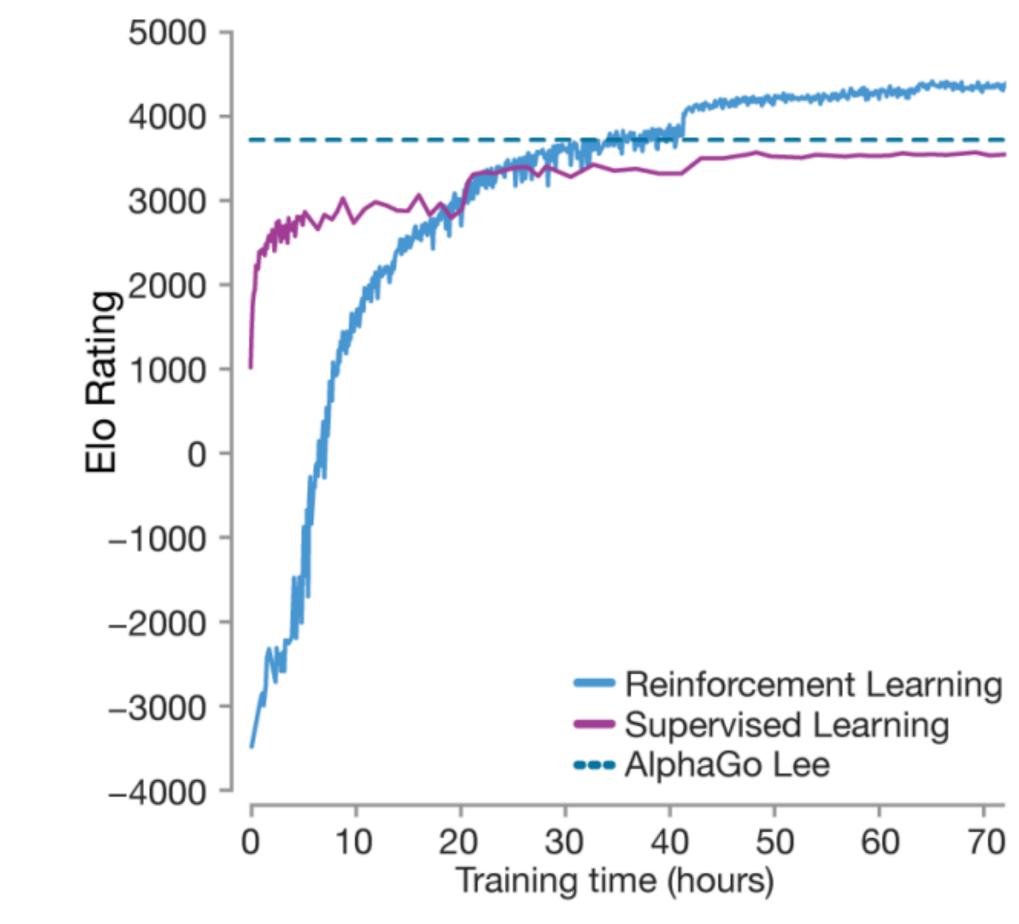
Source: https://rail.eecs.berkeley.edu/deeprlcourse-fa21/static/slides/lec-2.pdf

- Maximising things that are hard to maximise with next-token loss.
- Train the model on distributions that better match its output ("distribution drift").
 - Train in a "real" inference environment (vLLM != HF).



Source: https://yingru.notion.site/When-Speed-Kills-Stability-Demystifying-RL-Collapse-from-the-Training-Inference-Mismatch-271211a558b7808d8b12d403fd15edda

- Maximising things that are hard to maximise with next-token loss.
- Train the model on distributions that better match its output ("distribution drift").
- Allows us to go beyond bound of teacher data.



Source: Silver, David et al. "Mastering the game of Go without human knowledge." *Nature* 550 (2017): 354-359.

How do we do RL for LMs?

- Current approaches are all based on policy gradient.
 - Why?
 - Developed at OpenAl, OpenAl popularised RL for LMs.
 - Was the main 'go-to' approach in past few years.
 - Easier(?) to adapt to LMs: our LM is our policy, and we directly train that.
 - Policy gradient vs off-policy approaches (DQN, SAC) is an ongoing debate in RL-land.

How do we do RL for LMs?

- Current approaches are all based on policy gradient.
 - Why?
 - Developed at OpenAl, OpenAl popularised RL for LMs.
 - Was the main 'go-to' approach in past few years.
 - Easier(?) to adapt to LMs: our LM is our policy, and we directly train that.
 - Policy gradient vs off-policy approaches (DQN, SAC) is an ongoing debate in RL-land.

Interestingly and conversely, the *rliable* optimality gap indicates that DQN is better than PPO in Figure 6c, because DQN does achieve higher average policy performance, but the analysis from comparing to $V^{\hat{\pi}_D^*}(s_0)$, in Figure 6b shows us that even though DQN performs better than PPO, DQN is still generating a lot of high-value experience that it is not able to exploit. Conversely, because PPO is performing worse according to *rliable*, but has a better $V^{\hat{\pi}_D^*}(s_0) - V^{\hat{\pi}^\theta}(s_0)$, improved exploration would improve PPO more than it would DQN. Overall, these results suggest that both algorithms struggle to extract the most from their experience and that *rliable* is not telling the full story.

source: https://www.arxiv.org/abs/2508.01329

Is Exploration or Optimization the Problem for Deep Reinforcement Learning?

Glen Berseth

Université de Montréal, Mila - Quebec AI Institute, and CIFAR glen.berseth@mila.quebec

Abstract

In the era of deep reinforcement learning, making progress is more complex, as the collected experience must be compressed into a deep model for future exploitation and sampling. Many papers have shown that training a deep learning policy under the changing state and action distribution leads to sub-optimal performance, or even collapse. This naturally leads to the concern that even if the community creates improved exploration algorithms or reward objectives, will those improvements fall on the *deaf ears* of optimization difficulties. This work proposes a new *practical* sub-optimality estimator to determine optimization limitations of deep reinforcement learning algorithms. Through experiments across environments and RL algorithms, it is shown that the difference between the best experience generated is 2-3× better than the policies' learned performance. This large difference indicates that deep RL methods only exploit half of the good experience they generate.

1 Introduction

arXiv:2508.01329v1

What is preventing deep reinforcement learning from solving harder tasks? Many papers have shown that training a deep learning policy under the changing state distribution (non-IID) leads to sub-optimal performance [Nikishin et al., 2022, Lyle et al., 2023, Dohare et al., 2024]. However, at a macro scale, it is not completely clear what causes these issues. Do the network and regularization changes from recent work improve exploration or exploitation, and which of these two issues is the larger concern to be addressed to advance deep RL algorithms? For example, better exploration algorithms can be created, but will the higher value experience fall on the *deaf ears* of the deep network optimization difficulties?

How can we understand if the limited deepRL performance is due to a lack of good exploration or deep network optimization (exploitation)? Normally in RL, to understand if there is a limitation, an oracle is needed to understand *sub-optimality*, how far the algorithm is from being optimal. However, that analysis is with respect to the best policy and aliases both causes of the limitations of either exploration or optimization. Instead, consider the example where a person is learning how to build good houses. There are two issues that may prevent the person from *consistently* building a high quality house: (1) they can't *explore* well enough to discover a good design or (2) they can explore well enough to find good designs, but they can't properly *exploit* their experience to replicate those good experience. For deep RL algorithms, which of these two issues is more prevalent?

To understand if exploration or exploitation is the larger culprit, a method is needed to estimate the practical sub-optimality between these cases. This estimator should (1) measure the agent's ability to explore, (2) while also estimating the average performance for the learning policy π^{θ} . While estimating the average policy performance is common, estimating the exploration ability for a policy is not. Extending the house-building metaphor, the idea is to estimate how close the agent ever got to constructing a good home. Therefore, to realize this estimator, we propose computing the exploration value for a policy that is calculated over prior experience, called the experience optimal policy. Using

Preprint. Under review

Policy Gradient

For this, we consider a *parameterised policy*: a policy that can directly select what action to take, without necessarily consulting value (or q-) functions.

Effectively, we learn to produce a **probability distribution over actions** given a state.

$$\pi_{\theta}(a \mid s) =$$

 π_{θ} can be anything we can take gradients on! ...like a LM!

We can then aim to maximise performance by just taking gradient steps:

$$\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$$

We can then aim to maximise performance by just taking gradient steps:

$$\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$$

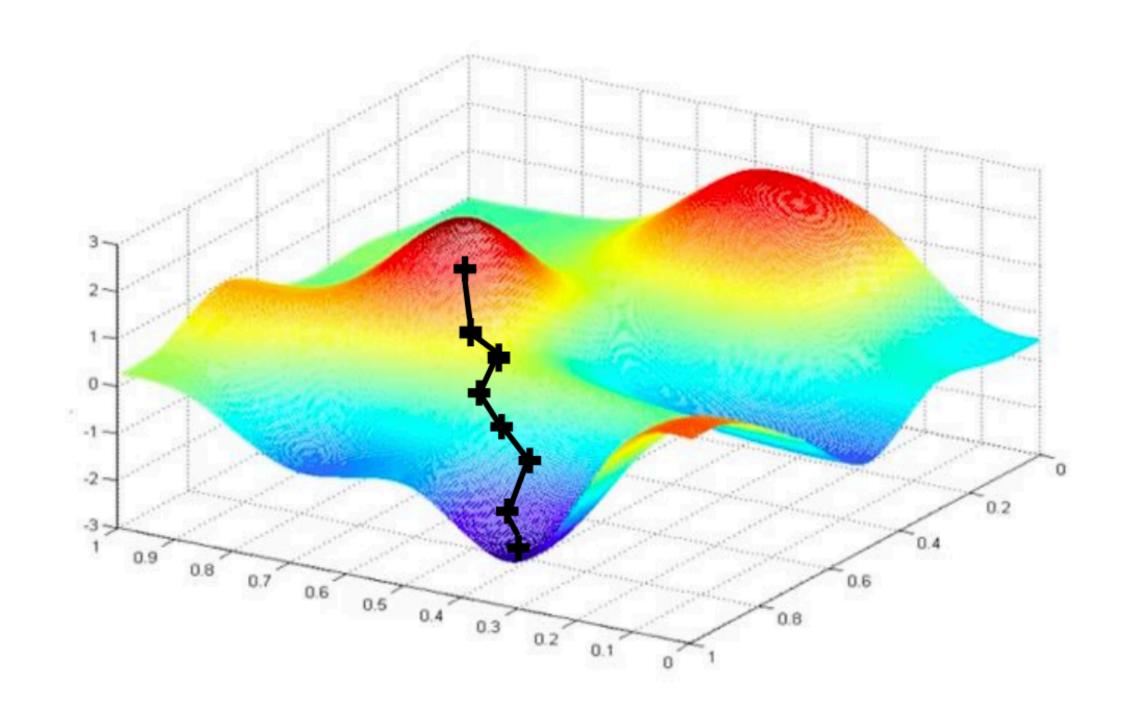
More basic details on gradient steps in a bit!

Basically just: if we take these steps, we will lower the loss associated with $J(\theta)$

Stochastic Gradient Descent

Algorithm

- I. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
- 2. Loop until convergence:
- 3. Pick single data point i
- 4. Compute gradient, $\frac{\partial J_i(\theta)}{\partial \theta}$
- 5. Update weights, $\theta \leftarrow \theta \eta \frac{\partial J(\theta)}{\partial \theta}$
- 6. Return weights



We can then aim to maximise performance by just taking gradient steps:

$$\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$$

 $J(\theta)$ here is our objective function. We will set it as the 'true' value of the initial state under our policy:

$$J(\theta) = \nu_{\pi_{\theta}}(s_0)$$

We can then aim to maximise performance by just taking gradient steps:

$$\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$$

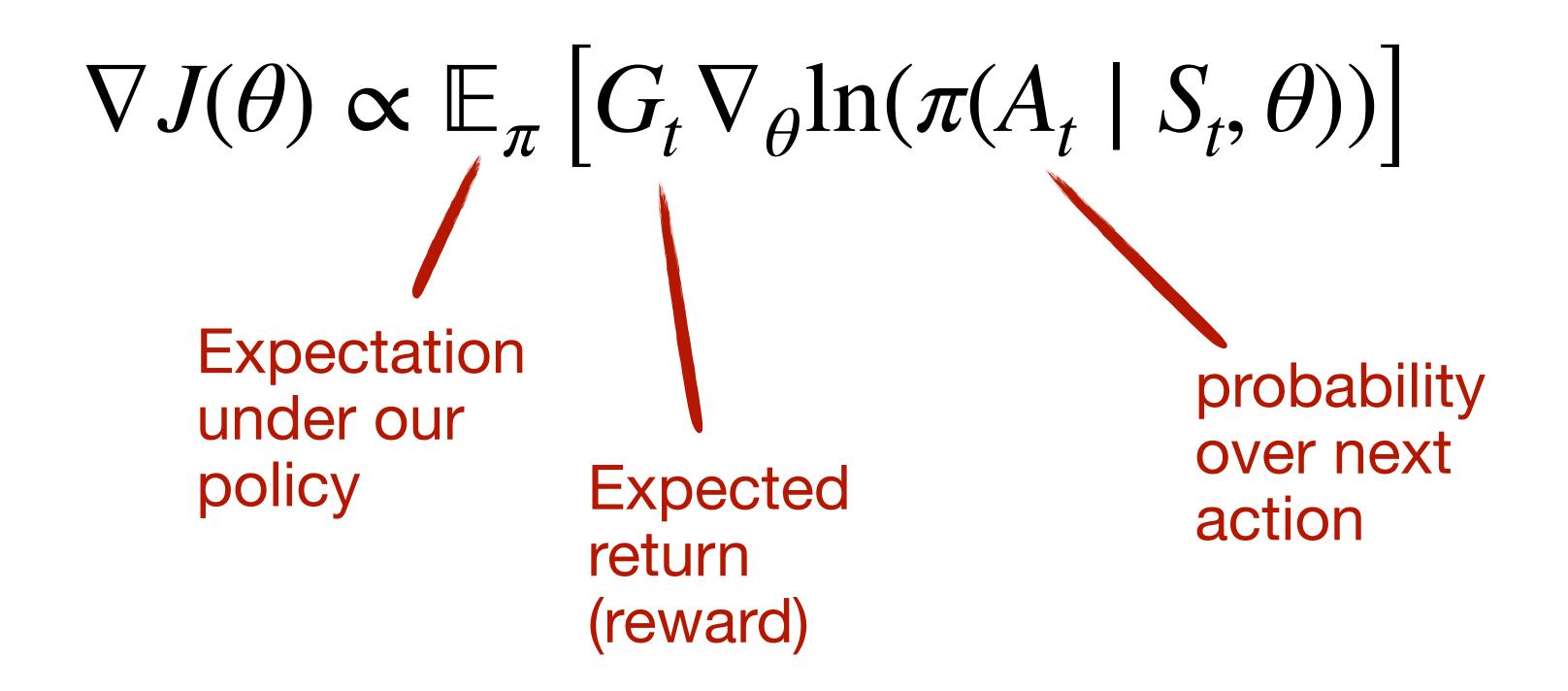
 $J(\theta)$ here is our objective function. We will set it as the 'true' value of the initial state under our policy:

$$J(\theta) = \nu_{\pi_{\theta}}(s_0)$$

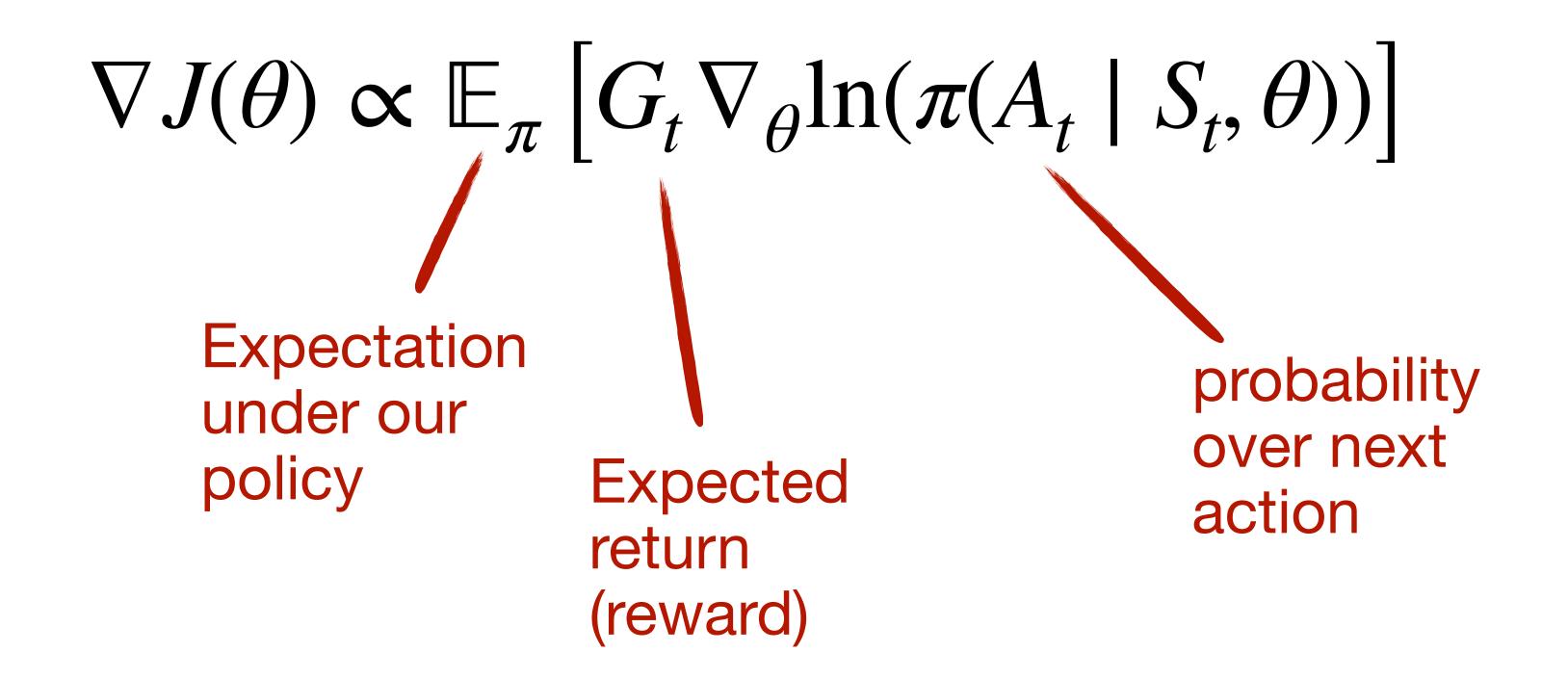
 $v_{\pi_{\theta}}$ is the value function for our parameterised policy, and s_0 is the initial state.

We could do simple gradient descent, but in practice we can plug into Adam, etc.

With some math (see Sutton & Barto, Chap. 13), we can show:



With some math (see Sutton & Barto, Chap. 13), we can show:



This is exactly the REINFORCE update!

REINFORCE: The full algorithm

Basically:

```
while True:
```

Generate a full episode ("rollout") starting from s_0 using π_{θ} Compute reward over the rollouts

Compute the REINFORCE loss and apply gradient

How is this instantiated for LMs?

$$\mathbb{E}_{\pi} \left[G_t \nabla_{\theta} \ln(\pi(A_t \mid S_t, \theta)) \right]$$

$$G_t = R(S_t, A_t)$$

We have some reward function

$$S_t = x$$

Our state is just the text seen so far

$$A_t = y_t$$

Our action is producing a token. Generating a full answer = episode.

$$\mathbb{E}_{\pi_{\theta}}[R(x,y) \nabla_{\theta} \ln \pi_{\theta}(y \mid x)]$$

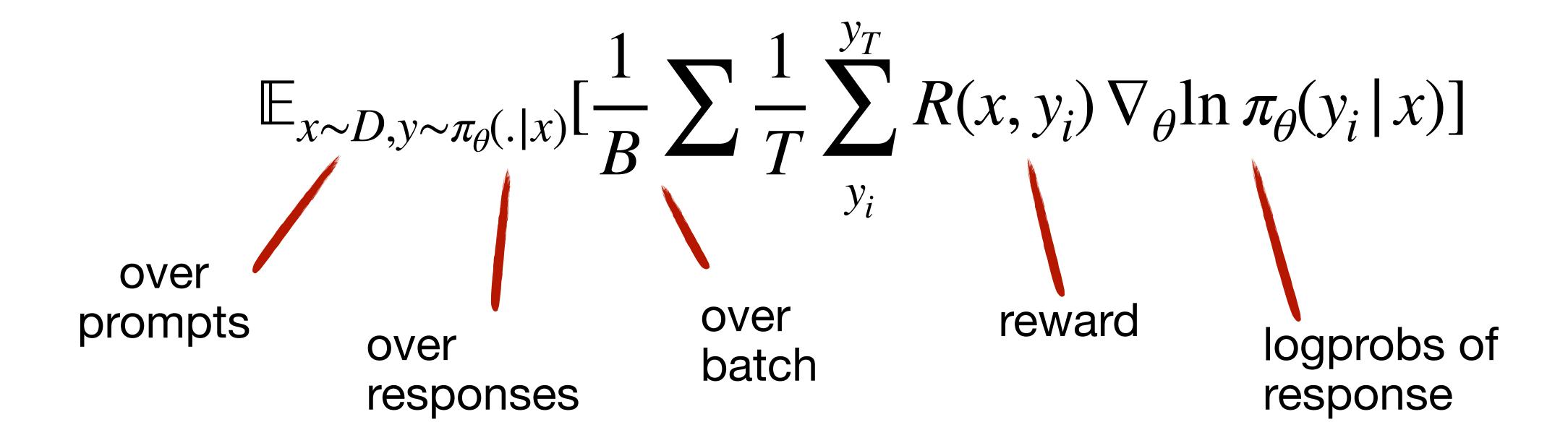
How is this instantiated for LMs?

In practice, we draw prompts (x) from a dataset (D), and sample completions (y) from our model, and do this in batches:

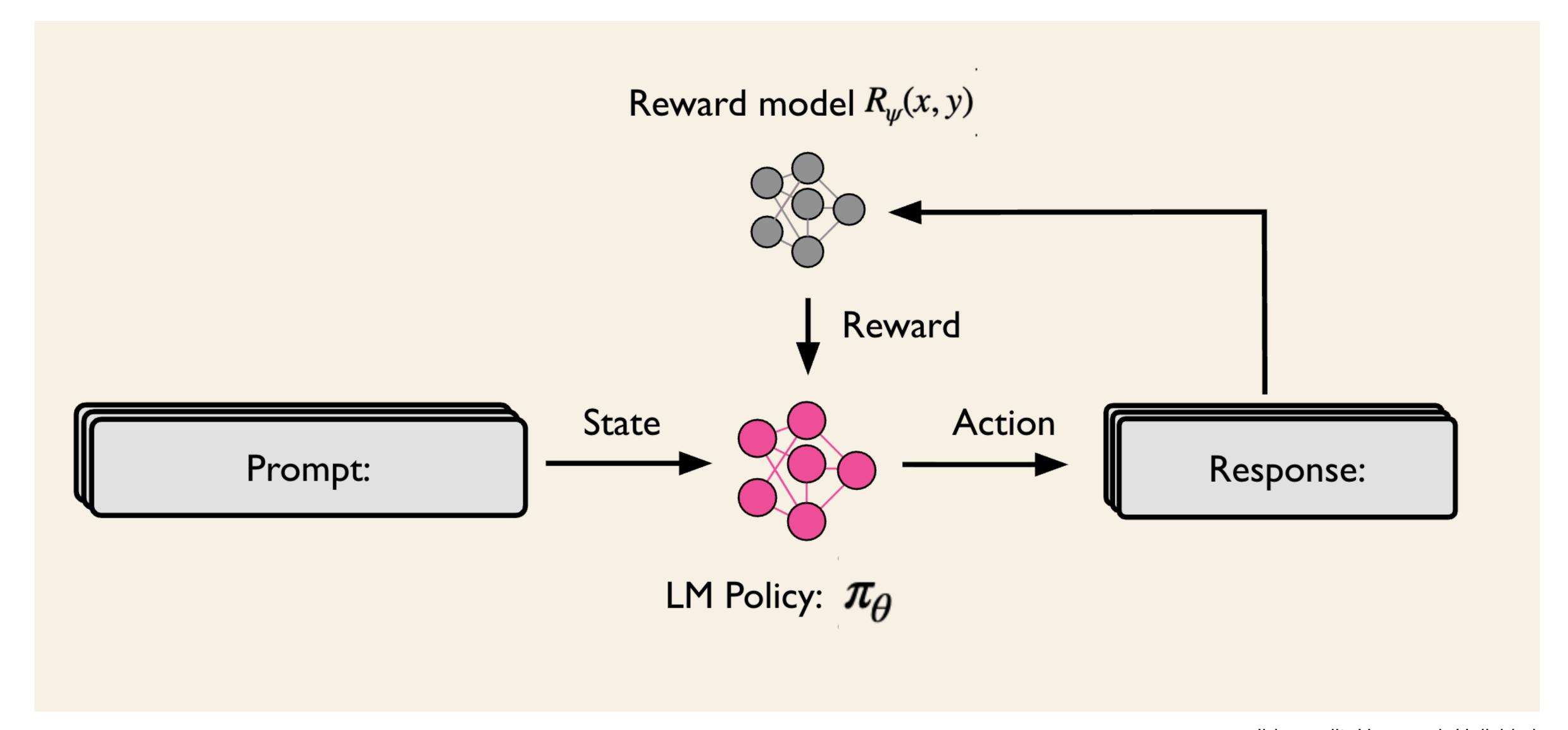
$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{i=1}^{N} R(x, y_i) \nabla_{\theta} \ln \pi_{\theta}(y_i | x) \right]$$

How is this instantiated for LMs?

In practice, we draw prompts (x) from a dataset (D), and sample completions (y) from our model, and do this in batches:



How is this instantiated for LMs?



$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{i=1}^{N} R(x, y_i) \nabla_{\theta} \ln \pi_{\theta}(y_i | x) \right]$$

- High variance
- Reward scale can cause issues... (consider neg. vs pos. reward)
- We can only use samples directly from our current loop

Baselining & Advantage

Core idea: try to normalise our reward as much as possible.

$$(R(y, x) - b) \nabla_{\theta} \ln \pi_{\theta}(y \mid x)$$

Core idea: try to normalise our reward as much as possible.

$$(R(y,x) - b) \nabla_{\theta} \ln \pi_{\theta}(y \mid x)$$

Most common choice for baseline: value estimate

$$(R(y,x) - \hat{V}_{\theta'}(x,y_{< t})) \nabla_{\theta} \ln \pi_{\theta}(y \mid x)$$

Core idea: try to normalise our reward as much as possible.

$$(R(y, x) - b) \nabla_{\theta} \ln \pi_{\theta}(y \mid x)$$

Most common choice for baseline: value estimate

$$(R(y, x) - \hat{V}_{\theta'}(x, y_{< t})) \nabla_{\theta} \ln \pi_{\theta}(y \mid x)$$

We call this form of normalising the advantage function:

$$A_{\theta'}(x, y) = R(x, y_t) - \hat{V}_{\theta'}(x, y_{< t})$$

Core idea: try to normalise our reward as much as possible.

$$(R(y, x) - b) \nabla_{\theta} \ln \pi_{\theta}(y \mid x)$$

Most common choice for baseline: value estimate

$$(R(y, x) - \hat{V}_{\theta'}(x, y_{< t})) \nabla_{\theta} \ln \pi_{\theta}(y \mid x)$$

We call this form of normalising the advantage function:

$$A_{\theta'}(x,y) = R(x,y_t) - \hat{V}_{\theta'}(x,y_{< t}) \text{lowest poss.}$$
knew V exactly!

Other forms of estimating advantage

$$A_{\theta'}(x, y) = R(x, y_t) - \hat{V}_{\theta'}(x, y_{< t})$$

PPO: Use Generalised Advantage Estimation (GAE)

HIGH-DIMENSIONAL CONTINUOUS CONTROL USING GENERALIZED ADVANTAGE ESTIMATION

John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan and Pieter Abbeel Department of Electrical Engineering and Computer Science University of California, Berkeley {joschu,pcmoritz,levine,jordan,pabbeel}@eecs.berkeley.edu

Other forms of estimating advantage

$$A_{\theta'}(x, y) = R(x, y_t) - \hat{V}_{\theta'}(x, y_{< t})$$

PPO: Use Generalised Advantage Estimation (GAE)

HIGH-DIMENSIONAL CONTINUOUS CONTROL USING GENERALIZED ADVANTAGE ESTIMATION

John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan and Pieter Abbeel
Department of Electrical Engineering and Computer Science
University of California, Berkeley

{ joschu, pcmoritz, levine, jordan, pabbeel}@eecs.berkeley.edu

$$\delta_t = r_t + \gamma \hat{V}(s_{t+1}) - \hat{V}(s_t)$$

TD residual - "how much better than expected will we do taking a particular action" we use a learnt value function

TD Residual?

$$\delta_t = r_t + \gamma \hat{V}(s_{t+1}) - \hat{V}(s_t)$$

TD residual - "how much better than expected will we do taking a particular action"

This should be familiar...

TD Residual?

The Bellman Equations

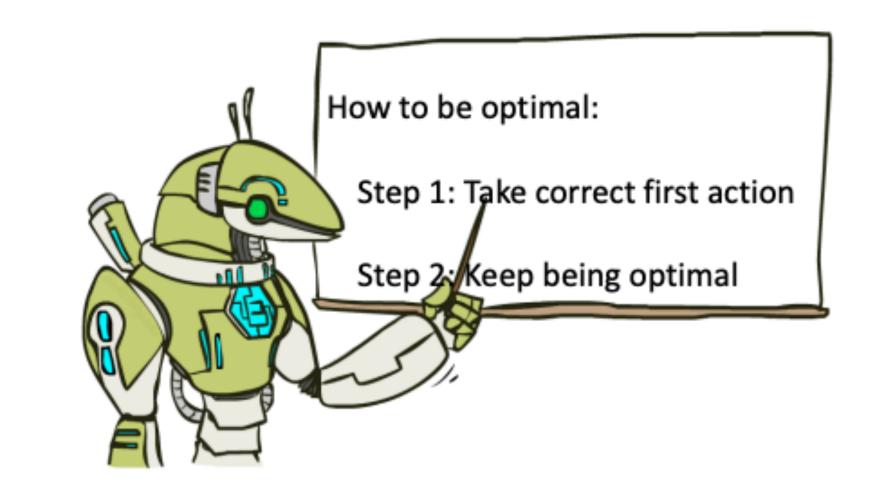
 Definition of "optimal utility" via expectimax recurrence gives a simple one-step lookahead relationship amongst optimal utility values

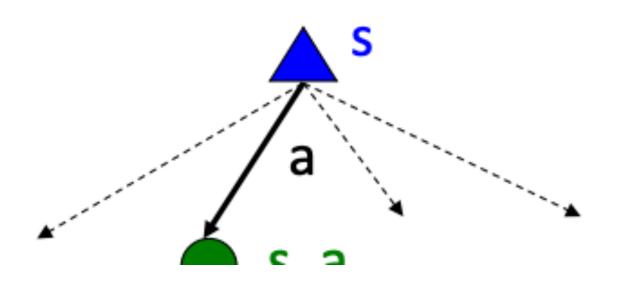
$$V^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V^*(s') \right]$$

$$V^*(s) = \max_{a} \sum_{s'} T(s, a, s') \left[R(s, a, s') + \gamma V^*(s') \right]$$

$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$$





Other forms of estimating advantage

$$A_{\theta'}(x, y) = R(x, y_t) - \hat{V}_{\theta'}(x, y_{< t})$$

PPO: Use Generalised Advantage Estimation (GAE)

HIGH-DIMENSIONAL CONTINUOUS CONTROL USING GENERALIZED ADVANTAGE ESTIMATION

 $\delta_t = r_t + \gamma \hat{V}(s_{t+1}) - \hat{V}(s_t)$ John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan and Pieter Abbeel Department of Electrical Engineering and Computer Science University of California, Berkeley {joschu, pcmoritz, levine, jordan, pabbeel}@eecs.berkeley.edu

TD residual - "how much better than expected did we do?" we use a learnt value function

Other forms of estimating advantage

$$A_{\theta'}(x, y) = R(x, y_t) - \hat{V}_{\theta'}(x, y_{< t})$$

PPO: Use Generalised Advantage Estimation (GAE)

HIGH-DIMENSIONAL CONTINUOUS CONTROL USING GENERALIZED ADVANTAGE ESTIMATION

John Schulman, Philipp Moritz, Sergey Levine, Michael I. Jordan and Pieter Abbeel
Department of Electrical Engineering and Computer Science
University of California, Berkeley
{joschu,pcmoritz,levine,jordan,pabbeel}@eecs.berkeley.edu

$$\delta_t = r_t + \gamma \hat{V}(s_{t+1}) - \hat{V}(s_t)$$

TD residual - "how much better than expected did we do?" we use a learnt value function

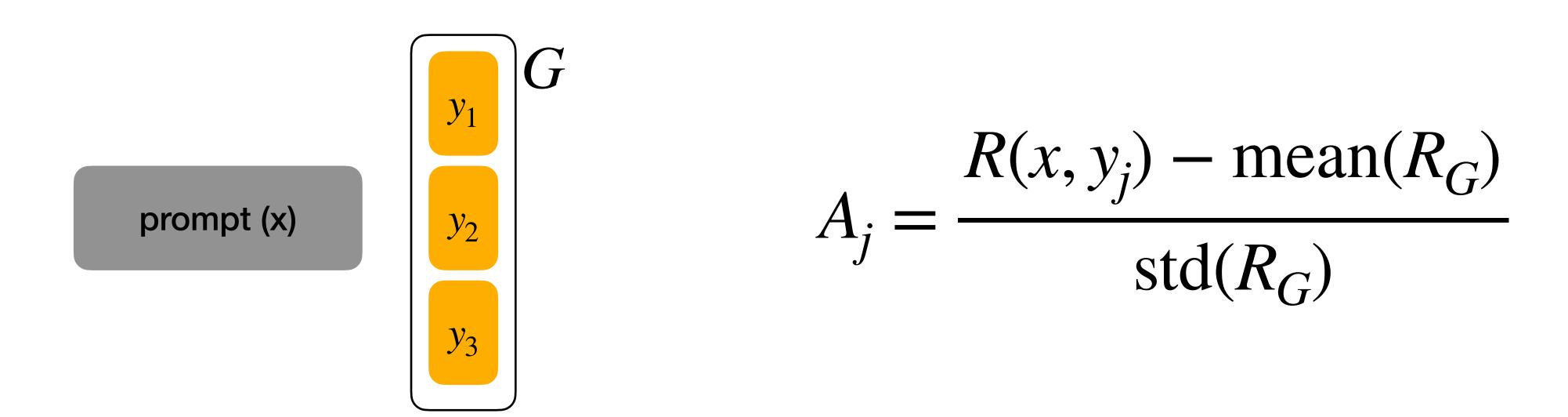
$$\hat{A}_t = \sum_{l=0}^{\infty} (\gamma \lambda)^l \, \delta_{t+l}^V$$

note that γ is usually set to 1, λ to 0.95 for RL w/ LMs!

We also typically normalise the advantages within the batch (whitening)

GRPO (Shao et al 2024): Use Group Estimates

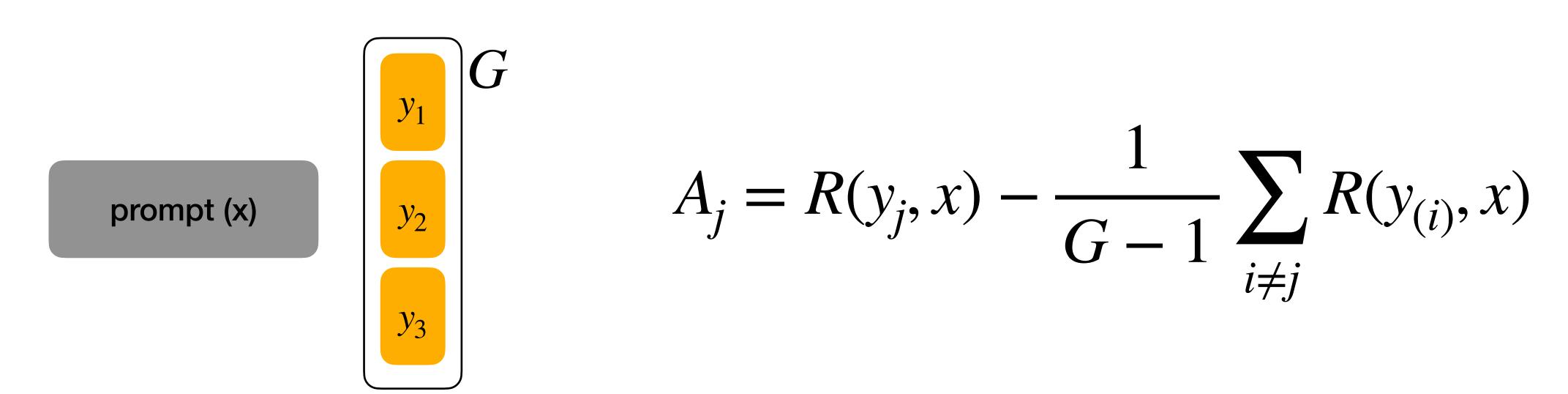
GRPO (Shao et al 2024): Use Group Estimates



No value estimation needed! Just take multiple samples per prompt and use average improvement over the group as advantage.

RLOO (Ahmadian et al 2024): <u>Use leave-one-out estimate</u>

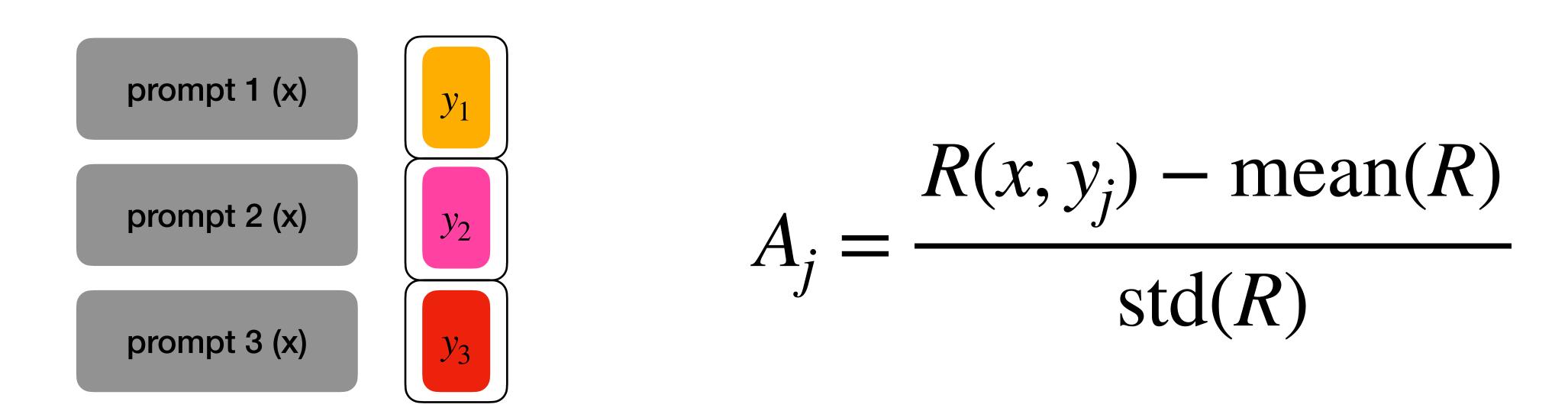
RLOO (Ahmadian et al 2024): Use leave-one-out estimate



Use average of other samples from the same prompt.

REINFORCE++ (Hu et al. 2025): <u>Use batch estimate (no groups!</u>)

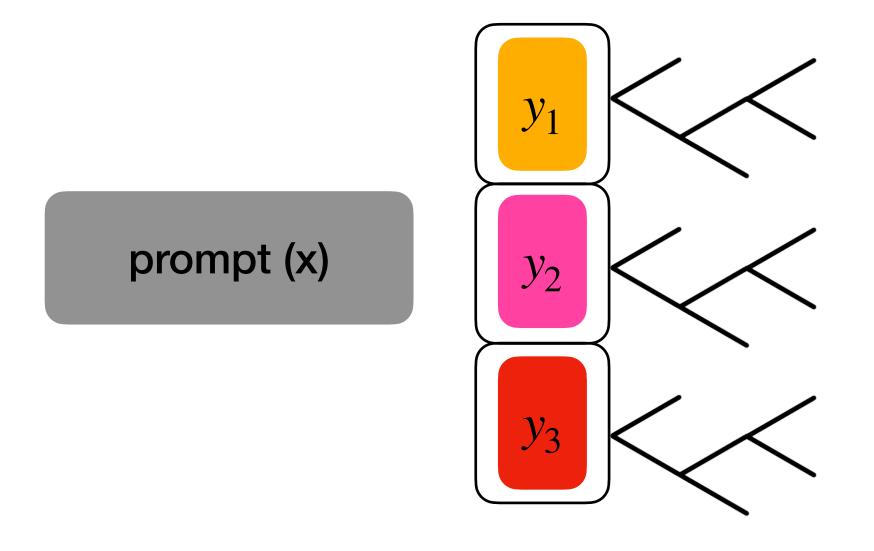
REINFORCE++ (Hu et al. 2025): Use batch estimate (no groups!)



Just normalise across everything in batch. Commonly done anyway ("whitening").

VinePPO (Kazemnejad et al 2024): Use Monte Carlo Estimates

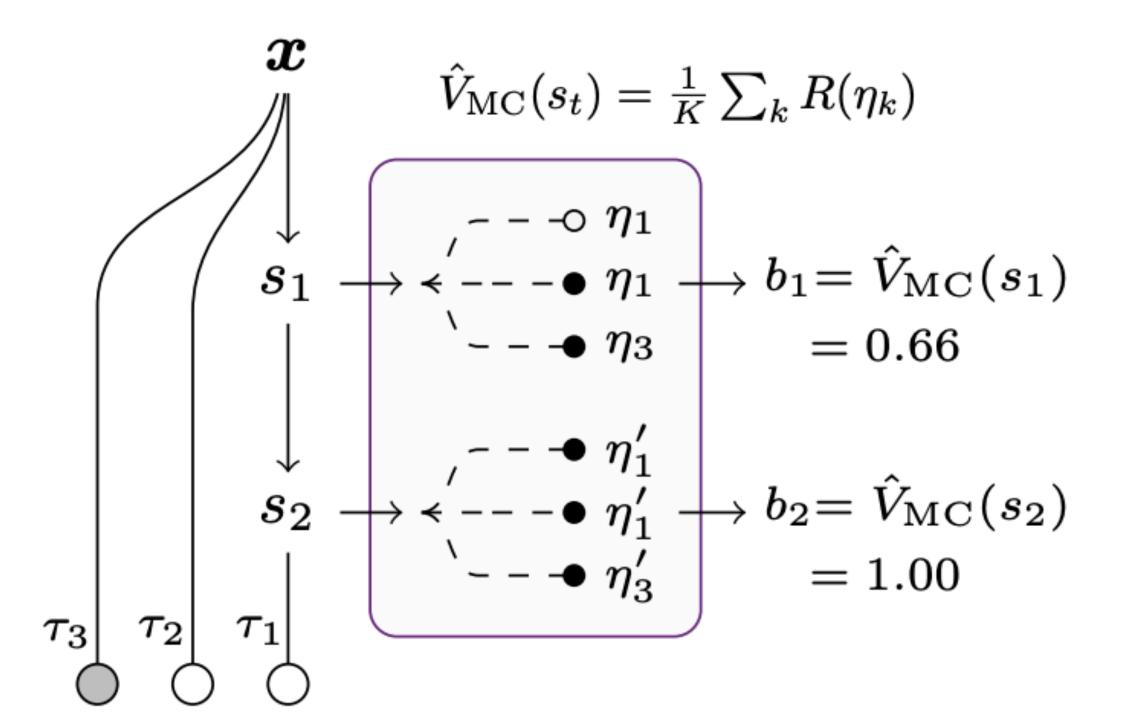
VinePPO (Kazemnejad et al 2024): Use Monte Carlo Estimates



No value estimation needed! Just take multiple samples *per* ~ *token* and use those!

VinePPO (Kazemnejad et al 2024): Use Monte Carlo Estimates

VinePPO

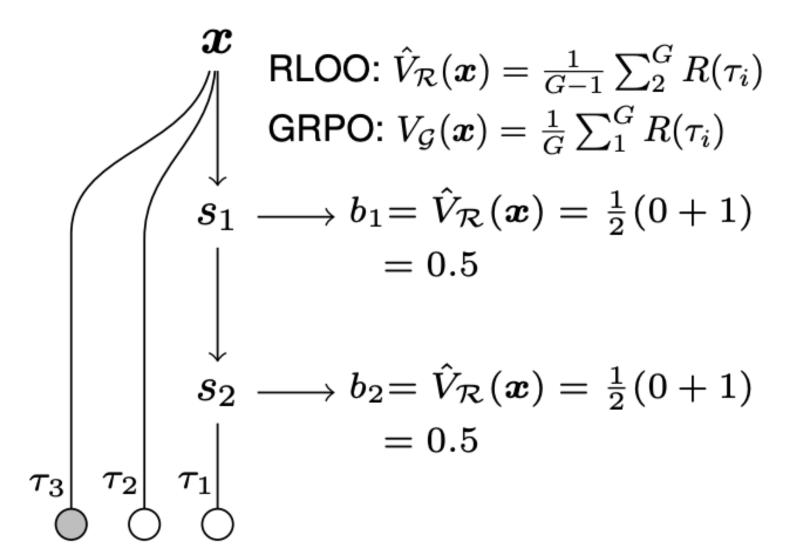


No value estimation needed! Just take multiple samples *per* ~ *token* and use those!

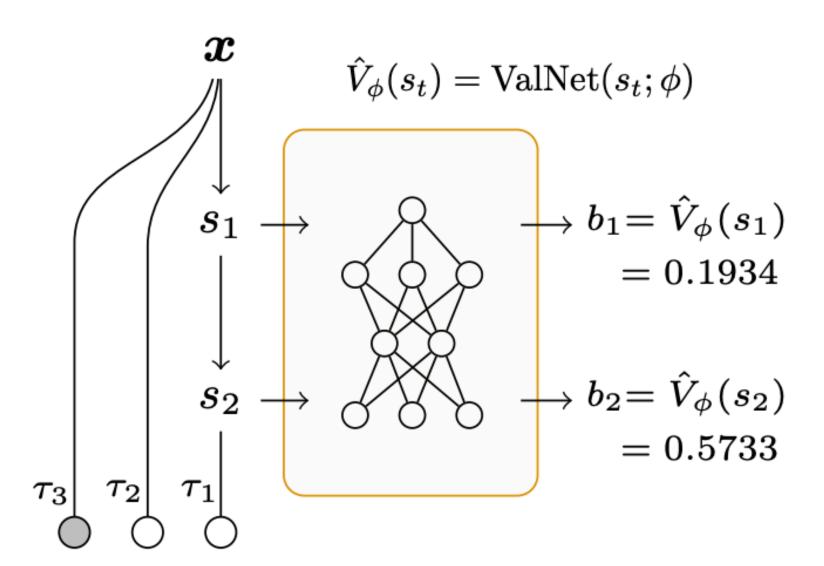
A summary of advantage estimates

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} [(A(y, x) \nabla_{\theta} \log \pi_{\theta}(y | x)]$$

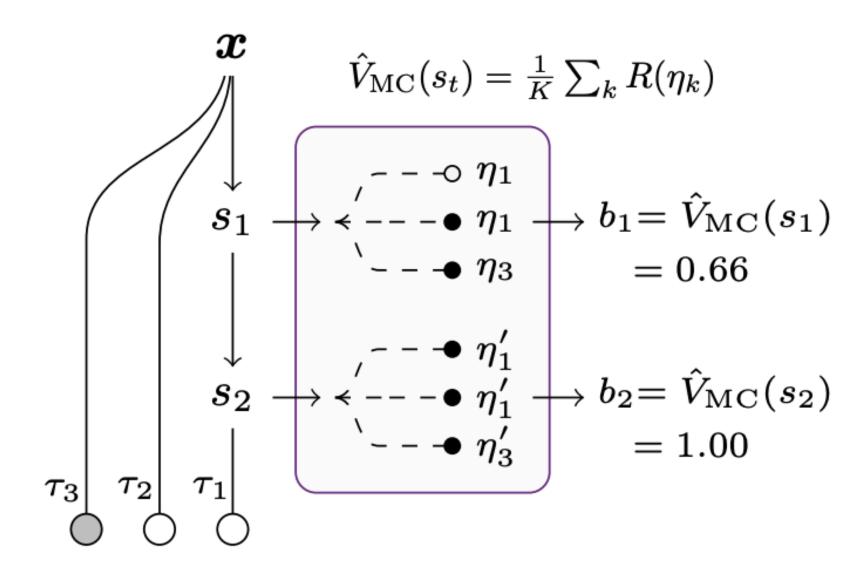
RLOO / GRPO



PPO



VinePPO



source: VinePPO: Unlocking RL Potential For LLM Reasoning Through Refined Credit Assignment (Kazemnejad et al 2024)

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{i=1}^{N} R(x, y_i) \nabla_{\theta} \ln \pi_{\theta}(y_i | x) \right]$$

- High variance
- Reward scale can cause issues... (consider neg. vs pos. reward)
- We can only use samples directly from our current loop

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{i=1}^{N} R(x, y_i) \nabla_{\theta} \ln \pi_{\theta}(y_i | x) \right]$$

- High variance
- Reward scale can cause issues... (consider neg. vs pos. reward)
- We can only use samples directly from our current loop

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{i=1}^{N} R(x, y_i) \nabla_{\theta} \ln \pi_{\theta}(y_i | x) \right]$$

- High variance
- Reward scale can cause issues... (consider neg. vs pos. reward)
- We can only use samples directly from our current loop

Importance Sampling & Off-Policy data

Revisiting REINFORCE

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{i=1}^{N} R(x, y_i) \nabla_{\theta} \ln \pi_{\theta}(y_i | x) \right]$$

Cannot reuse old data! -> can we get more efficient?

What if we re-used our old data anyway?

Well, our problem is we are doing our loss assuming sampling from our current policy...

Train on "old" samples w/ IS

We can reuse old samples with importance sampling!

$$\mathbb{E}_{\tau \sim p_{\theta}}[f(\tau)] = \mathbb{E}_{\tau \sim p_{\beta}} \left[\frac{p_{\theta}(\tau)}{p_{\beta}(\tau)} f(\tau) \right]$$

Train on "old" samples w/ IS

trajectories over

policy

We can reuse old samples with importance sampling!

$$\mathbb{E}_{\tau \sim p_{\theta}}[f(\tau)] = \mathbb{E}_{\tau \sim p_{\beta}} \left[\frac{p_{\theta}(\tau)}{p_{\beta}(\tau)} f(\tau) \right]$$
 some function on

Train on "old" samples w/ IS

trajectories over

policy

We can reuse old samples with importance sampling!

$$\mathbb{E}_{\tau \sim p_{\theta}}[f(\tau)] = \mathbb{E}_{\tau \sim p_{\beta}} \left[\frac{p_{\theta}(\tau)}{p_{\beta}(\tau)} f(\tau) \right]$$
 some function on

Same function, but now we drawing trajectories from diff policy, and reweight accordingly!

Applying IS to REINFORCE

We can apply this directly to reinforce! Our distribution is just the older policy.

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \dots \ln \pi_{\theta}(y_i|x)$$

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta_{\text{old}}}(.|x)} \cdots \frac{\ln \pi_{\theta}(y_i|x)}{\ln \pi_{\theta_{\text{old}}}(y_i|x)}$$

Applying IS to REINFORCE

We can apply this directly to reinforce! Our distribution is just the older policy.

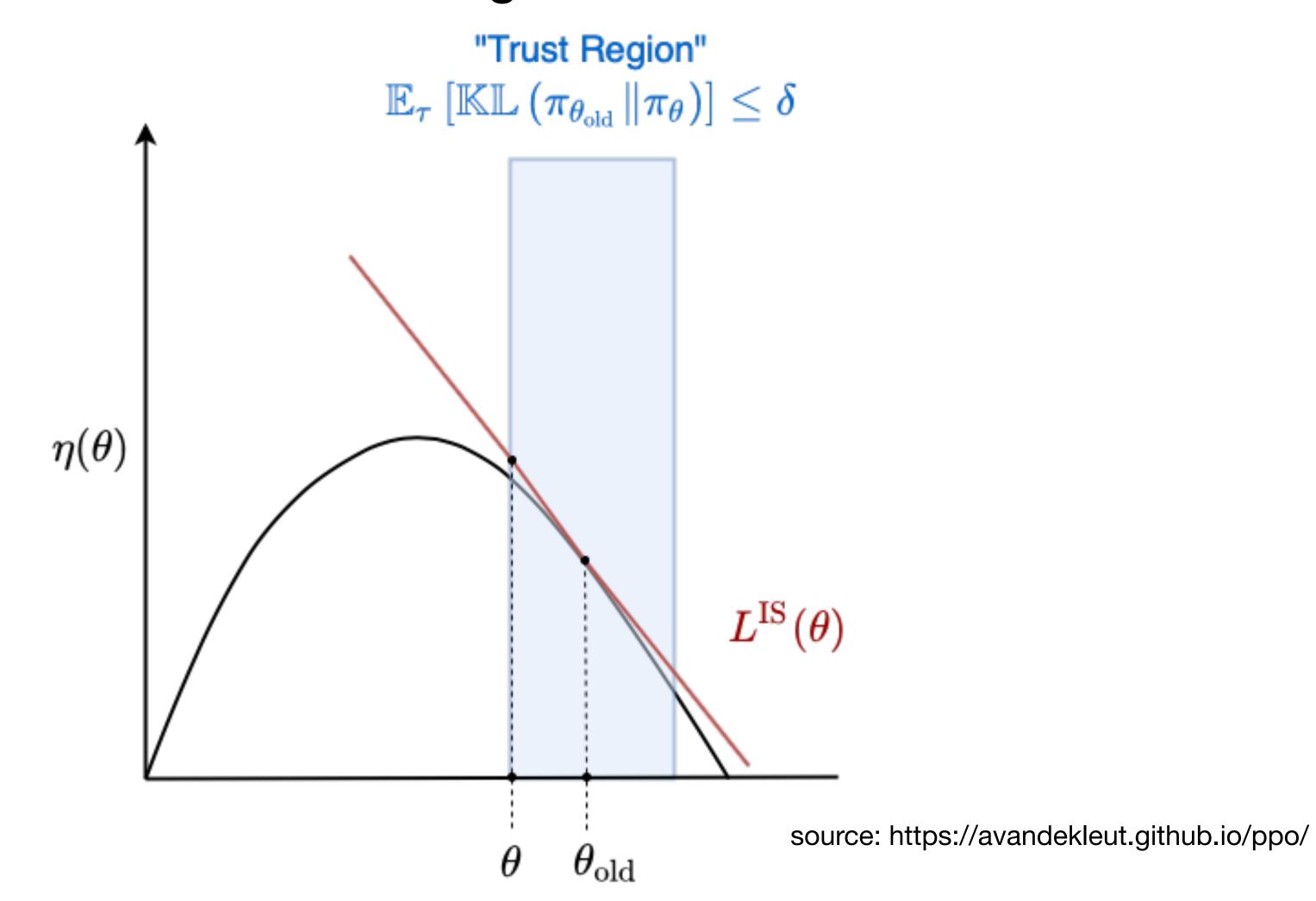
$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \dots \ln \pi_{\theta}(y_i | x)$$

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta_{\text{old}}}(.|x)} \cdots \frac{\ln \pi_{\theta}(y_i|x)}{\ln \pi_{\theta_{\text{old}}}(y_i|x)}$$

problem: if the ratio gets too crazy, our estimate 'blows' up.

Applying IS to REINFORCE

Main fix: just stay within some "safe" region!



TRPO: Keep within some KL constraint

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta_{\text{old}}}(.|x)} \cdots \frac{\ln \pi_{\theta}(y_i|x)}{\ln \pi_{\theta_{\text{old}}}(y_i|x)}$$

s.t.
$$\hat{\mathbb{E}}_t \left[\text{KL} \left[\pi_{\theta} \text{old}(\cdot \mid s_t), \pi_{\theta}(\cdot \mid s_t) \right] \right] \leq \delta$$

Don't let the policy stray too far!

Requires some more complex techniques to solve...

PPO: just clip the ratio!

Use the same importance sampling technique:

$$r_i = \frac{\ln \pi_{\theta}(y_i | x)}{\ln \pi_{\theta_{\text{old}}}(y_i | x)}$$

But now just clip the ratio to avoid moving too far:

$$L = \hat{\mathbb{E}} \left[\min \left(r_i(\theta) \hat{A}_i, \operatorname{clip} \left(r_i(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_i \right) \right]$$

Why min on top of clip? Because if the estimate would make advantage higher, its okay - we just want to be pessimistic about true advantage.

Standard PPO = REINFORCE+clip+advantage

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{y_i}^{y_T} \left[\min \left(r_i(\theta) \hat{A}_i, \operatorname{clip} \left(r_i(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_i \right) \right]$$

$$r_i = \frac{\ln \pi_{\theta}(y_i | x)}{\ln \pi_{\theta_{\text{old}}}(y_i | x)} \qquad A_i(x, y) = R(y, x) - \hat{V}_{\theta'}(x)$$

DAPO (Yu et al 2025): Uncouple the clipping parameter, and raise the higher one.

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i} \frac{1}{T} \sum_{y_{i}}^{y_{T}} \left[\min \left(r_{i}(\theta) \hat{A}_{i}, \operatorname{clip} \left(r_{i}(\theta), 1 - \epsilon_{lower}, 1 + \epsilon_{higher} \right) \hat{A}_{i} \right) \right]$$

DAPO (Yu et al 2025): Uncouple the clipping parameter, and raise the higher one.

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i} \frac{1}{T} \sum_{y_{i}}^{y_{T}} \left[\min \left(r_{i}(\theta) \hat{A}_{i}, \operatorname{clip} \left(r_{i}(\theta), 1 - \epsilon_{lower}, 1 + \epsilon_{higher} \right) \hat{A}_{i} \right) \right]$$

Basically: allow larger updates when raising token probabilities, but don't further decrease the lower bound.

Minimax paper goes further and sets $\epsilon_{\mathrm{lower}}$ to infinity to avoid lower bound clipping!

CISPO (Minimax 2025): directly clip ratio, instead of overall advantage.

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{y_{i}}^{y_{T}} \left[\min \left(r_{i}(\theta), \operatorname{clip}\left(r_{i}(\theta), 1 - \epsilon_{lower}, 1 + \epsilon_{higher}\right) \right) \hat{A}_{i} \right]$$

CISPO (Minimax 2025): directly clip ratio, instead of overall advantage.

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{y_{i}}^{y_{T}} \left[\min \left(r_{i}(\theta), \operatorname{clip} \left(r_{i}(\theta), 1 - \epsilon_{lower}, 1 + \epsilon_{higher} \right) \right) \hat{A}_{i} \right]$$

Basically: directly clip the ratio, to allow large advantages to still have some effect - even when lots of clipping might be happening.

GSPO (Zheng et al 2025): use sequence-level importance sampling instead of token-level.

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{y_{i}}^{y_{T}} \left[\min \left(\mathbf{s}_{i}(\theta) \hat{A}_{i}, \operatorname{clip} \left(\mathbf{s}_{i}(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i} \right) \right]$$

$$s_i(\theta) = \left(\frac{\pi_{\theta}(y_i \mid x)}{\pi_{\theta_{\text{old}}}(y_i \mid x)}\right)^{\frac{1}{|y_i|}} = \exp\left(\frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \ln \frac{\pi_{\theta}(y_i \mid x)}{\pi_{\theta_{\text{old}}}(y_i \mid x)}\right)$$

GSPO (Zheng et al 2025): use sequence-level importance sampling instead of token-level.

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i} \frac{1}{T} \sum_{y_{i}}^{y_{T}} \left[\min \left(\mathbf{s}_{i}(\theta) \hat{A}_{i}, \operatorname{clip} \left(\mathbf{s}_{i}(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i} \right) \right]$$

$$s_i(\theta) = \left(\frac{\pi_{\theta}(y_i \mid x)}{\pi_{\theta_{\text{old}}}(y_i \mid x)}\right)^{\frac{1}{|y_i|}} = \exp\left(\frac{1}{|y_i|} \sum_{t=1}^{|y_i|} \ln \frac{\pi_{\theta}(y_i \mid x)}{\pi_{\theta_{\text{old}}}(y_i \mid x)}\right)$$

Basically: compute the importance ratio of the whole sequence, and then apply the same value to each token. Reduces noise of token-to-token.

Going beyond REINFORCE

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{i=1}^{N} A(x, y_i) \nabla_{\theta} \ln \pi_{\theta}(y_i | x) \right]$$

Issues:

- High variance
- Reward scale can cause issues... (consider neg. vs pos. reward)
- We can only use samples directly from our current loop

Going beyond REINFORCE

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{y_{i}}^{y_{T}} \left[\min \left(r_{i}(\theta) \hat{A}_{i}, \operatorname{clip} \left(r_{i}(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i} \right) \right]$$

Issues:

- High variance
- Reward scale can cause issues... (consider neg. vs pos. reward)
- We can only use samples directly from our current loop T

Going beyond REINFORCE

$$\mathbb{E}_{x \sim D, y \sim \pi_{\theta}(.|x)} \left[\frac{1}{B} \sum_{i=1}^{N} \frac{1}{T} \sum_{y_{i}}^{y_{T}} \left[\min \left(r_{i}(\theta) \hat{A}_{i}, \operatorname{clip} \left(r_{i}(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i} \right) \right]$$

$$r_i = \frac{\ln \pi_{\theta}(y_i \mid x)}{\ln \pi_{\theta_{\text{old}}}(y_i \mid x)} \qquad A_i(x, y) = R(y, x) - \hat{V}_{\theta'}(x)$$

One last detail: KL penalty

In practice, we sometime also apply a KL penalty to the reward:

$$\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(y|x)} [r_{\phi}(x, y)] - \beta \mathbb{D}_{\text{KL}} [\pi_{\theta}(y \mid x) \mid\mid \pi_{\text{ref}}(y \mid x)]$$

However, these days often the KL penalty is completely removed. So, I'm going to skip over this.

Backing up to the full algorithm

Algorithm 1 PPO-Clip

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: **for** k = 0, 1, 2, ... **do**
- 3: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 4: Compute rewards-to-go \hat{R}_t .
- 5: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 6: Update the policy by maximizing the PPO-Clip objective:

$$\theta_{k+1} = \arg\max_{\theta} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^{T} \min\left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)} A^{\pi_{\theta_k}}(s_t, a_t), \ g(\epsilon, A^{\pi_{\theta_k}}(s_t, a_t))\right),$$

typically via stochastic gradient ascent with Adam.

7: Fit value function by regression on mean-squared error:

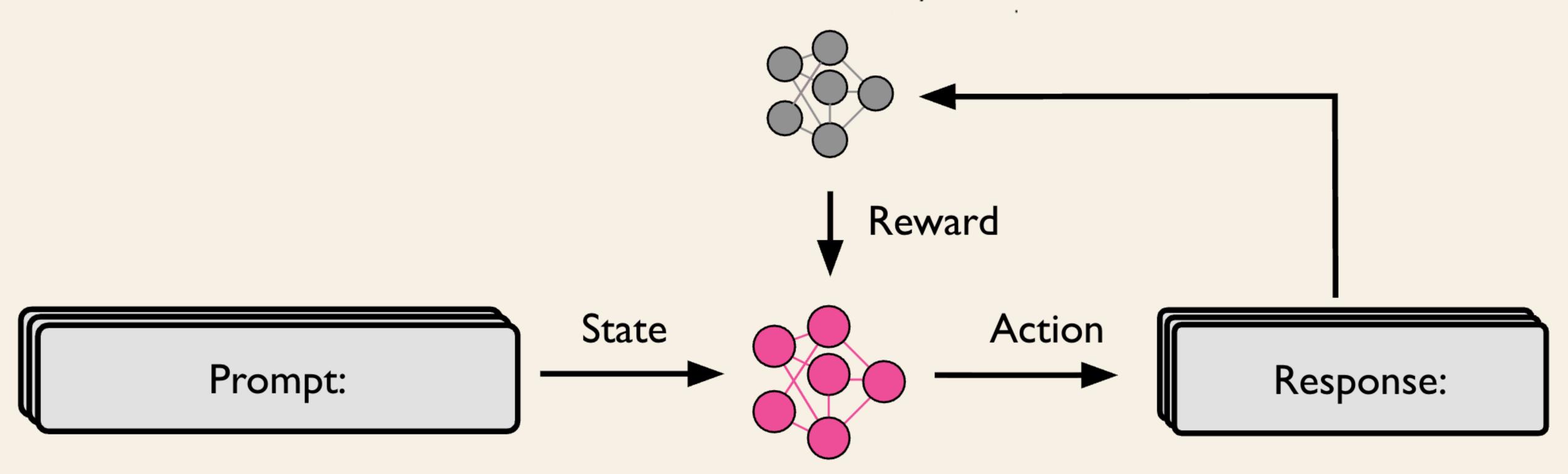
$$\phi_{k+1} = \arg\min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^{T} \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

8: end for

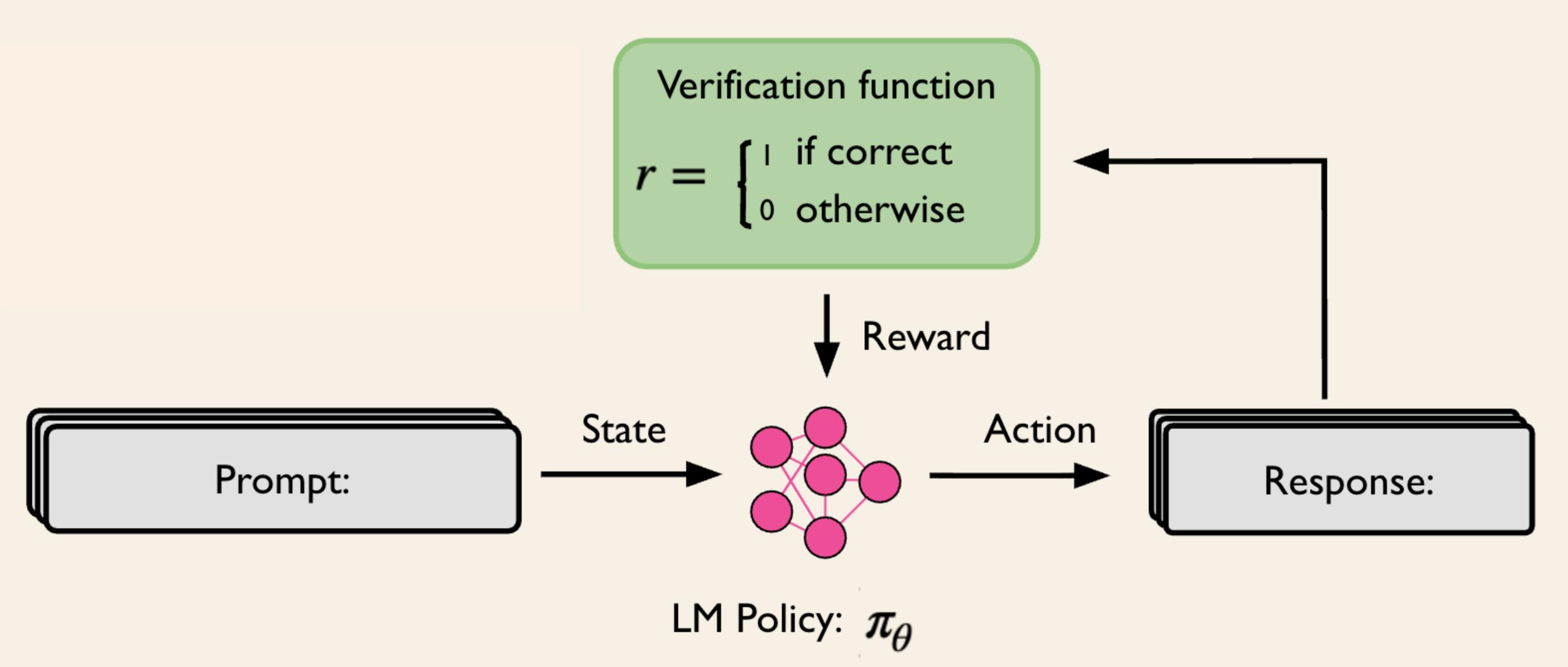
Backing up to the full algorithm

Reward model $R_{\psi}(x, y)$



LM Policy: $\pi_{ heta}$

Backing up to the full algorithm



Adapting for LMs

Adapting for LMs

- 1. LM inference is slow. How can we speed things up?
- 2. How do we deal with long sequences?
- 3. How do we deal with zero-advantage groups (GRPO)?

Speeding up inference.

Lots of work done on speeding up inference by great teams focussed on inference-only settings. Much of this can be applied in the RL setting!



NVIDIA TensorRT

One big issue: logprob computation can mismatch between inference and training! (Yao et al, 2025)

$$\mathbb{E}_{a \sim oldsymbol{\pi_{ ext{sampler}}(heta_{ ext{old}})}} iggl[
abla_{ heta} \min iggl(rac{oldsymbol{\pi_{ ext{learner}}(a, heta)}}{oldsymbol{\pi_{ ext{learner}}(a, heta_{ ext{old}})}} \, \hat{A}, \; ext{clip} iggl(rac{oldsymbol{\pi_{ ext{learner}}(a, heta)}}{oldsymbol{\pi_{ ext{learner}}(a, heta_{ ext{old}})}}, \, 1 - \epsilon, \, 1 + \epsilon igr) \, \hat{A} iggr) iggr].$$

We are sampling via vllm, but computing our logprobs in our training code.

One big issue: logprob computation can mismatch between inference and training!

Is this problem familiar?

Train on "old" samples w/ IS

We can reuse old samples with importance sampling!

$$\mathbb{E}_{\tau \sim p_{\theta}}[f(\tau)] = \mathbb{E}_{\tau \sim p_{\beta}} \left[\frac{p_{\theta}(\tau)}{p_{\beta}(\tau)} f(\tau) \right]$$
 some function on trajectories over policy Same function, but now we drawing trajectories from diff policy, and reweight accordingly!

One big issue: <u>logprob computation can mismatch between inference and</u> <u>training!</u>

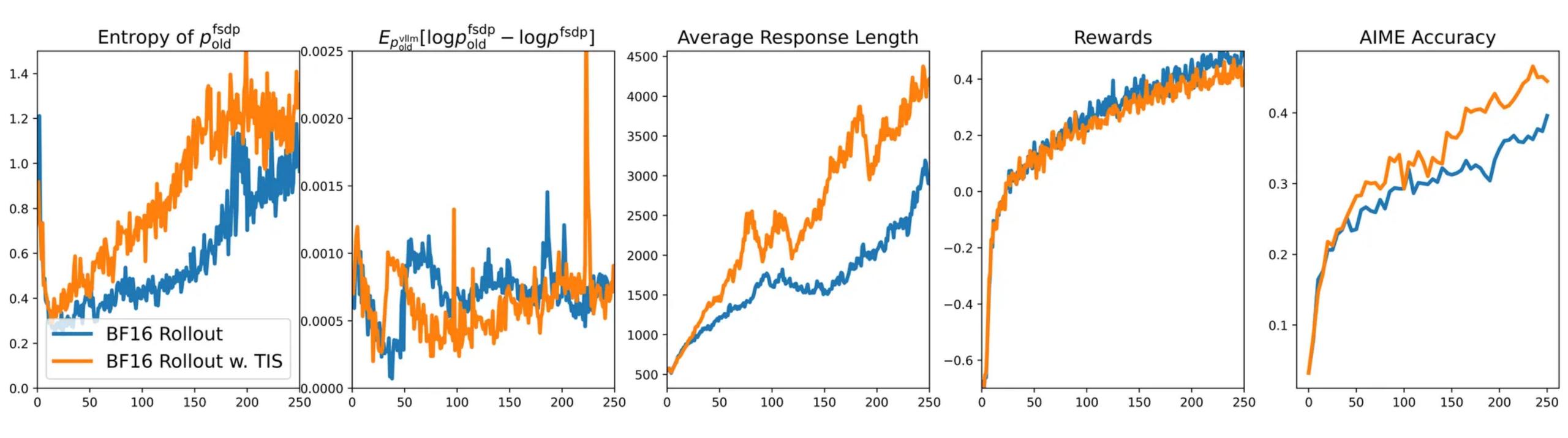
We apply the same solution as we did for training on old data: importance sampling!

$$\mathbb{E}_{a \sim oldsymbol{\pi_{ ext{sampler}}(heta_{ ext{old}})}} iggl[
abla_{ heta} \min iggl(rac{oldsymbol{\pi_{ ext{learner}}(a, heta)}{oldsymbol{\pi_{ ext{learner}}(a, heta_{ ext{old}})}} \hat{A}, \; ext{clip} iggl(rac{oldsymbol{\pi_{ ext{learner}}(a, heta)}{oldsymbol{\pi_{ ext{learner}}(a, heta_{ ext{old}})}}, \, 1 - \epsilon, \, 1 + \epsilon igr) \hat{A} iggr) iggr].$$



$$\mathbb{E}_{a \sim \pi_{\text{sampler}}(\theta_{\text{old}})} \left[\underbrace{\min \left(\frac{\pi_{\text{learner}}(a, \theta_{\text{old}})}{\pi_{\text{sampler}}(a, \theta_{\text{old}})}, C \right)}_{\text{truncated importance ratio}} \cdot \nabla_{\theta} \ \min \left(\frac{\pi_{\text{learner}}(a, \theta)}{\pi_{\text{learner}}(a, \theta_{\text{old}})} \, \hat{A}, \text{clip} \left(\frac{\pi_{\text{learner}}(a, \theta)}{\pi_{\text{learner}}(a, \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \right) \hat{A} \right) \right]$$

This can be the difference between collapse and success. Empirically, I find it more important the more off-policy you go.



source: Yao et al. 2025

More recently: maybe just using FP32 precision for LM head works?

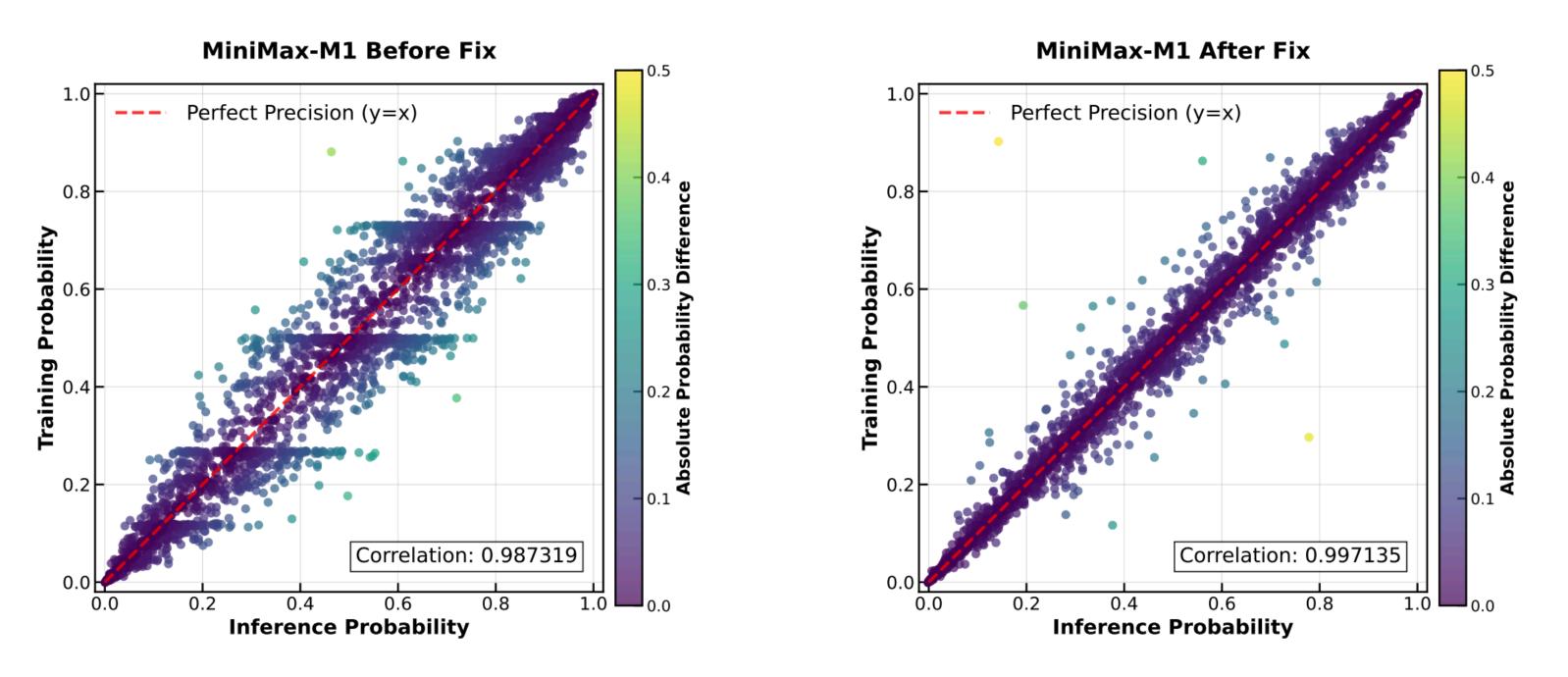
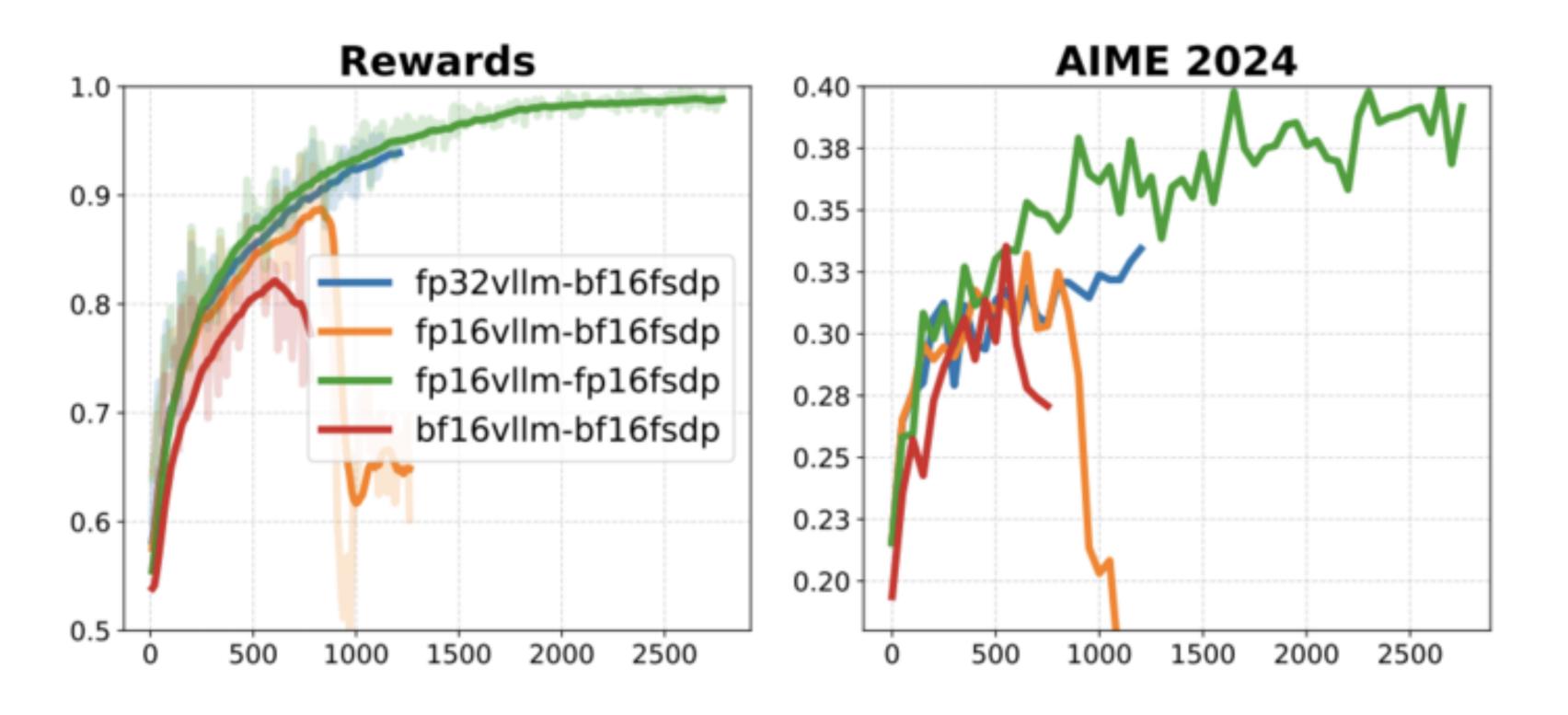


Figure 3 | Probability of tokens in training-mode code vs. probability of tokens in inference-mode code. Each point in the figures represents an individual token. The Pearson correlation coefficient is indicated in the figures. Theoretically, the two probabilities should be identical, and all the tokens should be exactly on the diagonal line. **Left:** Correlation of the M1 model before our fix; **Right:** Correlation of the M1 model after applying our fix of using FP32 precision for the LM output head.

More recently: maybe just using FP16 precision fixes all these issues???

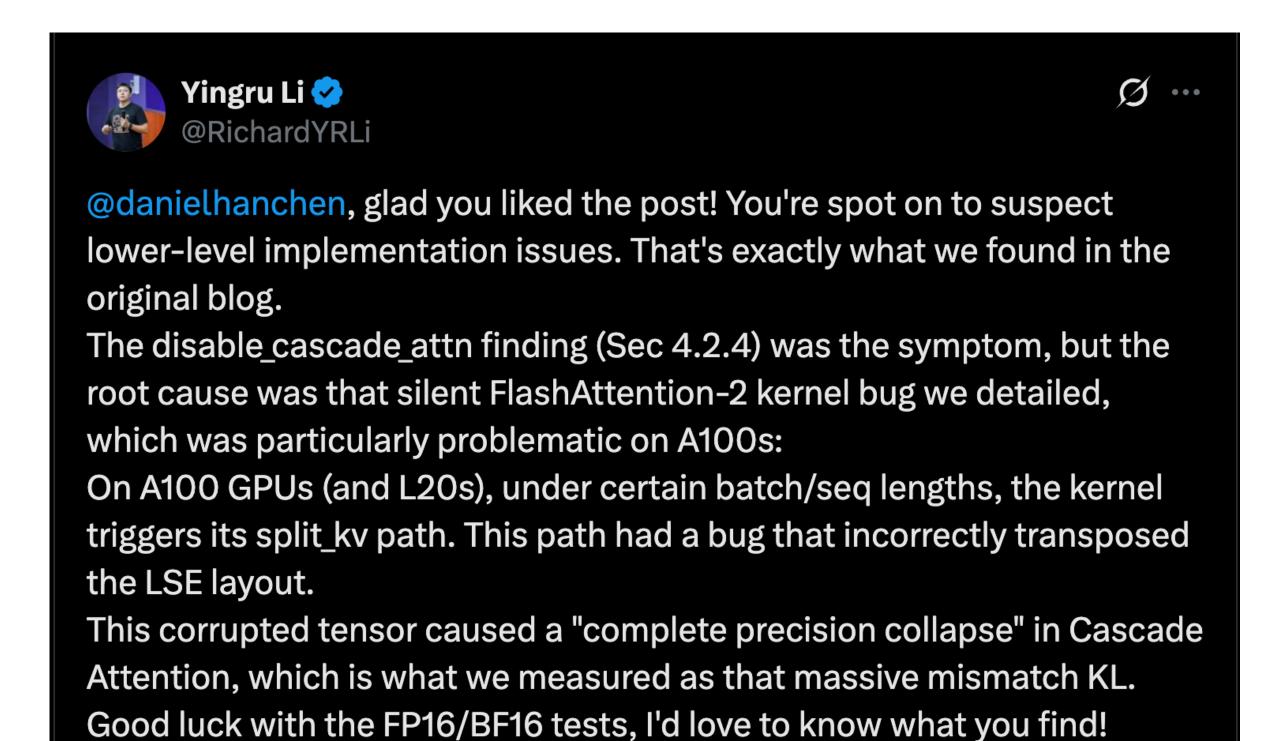


only stable when:

1. FP32 vllm, BF16 training.

2. FP16 vllm, FP16 training

Even more recently: maybe theres just bugs in the inference code????



source: https://x.com/RichardYRLi/status/1984858850143715759

yingru.notion.site/When-Speed-Kil...

4.2.4 Diable Cascade Attention in VLLM

According to this <u>GitHub issue</u>, we set <u>disable_cascade_attn=True</u> when initializing the vLLM engine and found that it significantly helps reduce the training-inference mismatch in experiments conducted on A100 GPUs. We performed two on-policy GRPO experiments using Qwen3-14B-Base as the base model on A100 GPUs, with <u>disable_cascade_attn</u> set to <u>True</u> and <u>False</u>, respectively. The results are shown in Figure 21. It can be observed that after disabling cascade attention, the <u>vllm-kl</u> metric decreased from the range of 5e-2 to 1e-1 to around 1e-3, indicating a substantial reduction in training-inference mismatch. In addition, the rewards on the training set also increased appropriately.

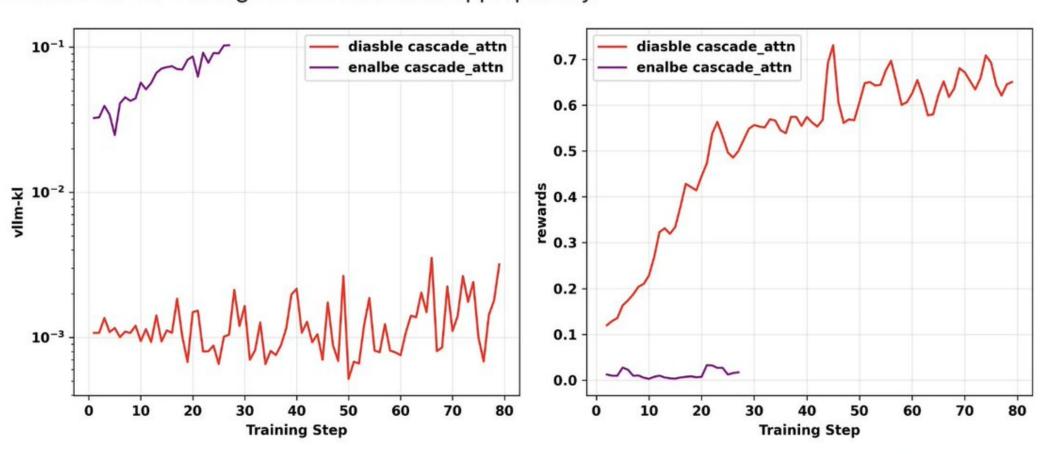


Figure 21: Results of the ablation study on using cascade attention. We conducted two on-policy experiments on A100 GPUs: one with cascade attention enabled and another with the same configuration but with cascade attention disabled. Left: The dynamics of vllm-kl metric. Right: the dynamics of training rewards

source: Liu et al 2025

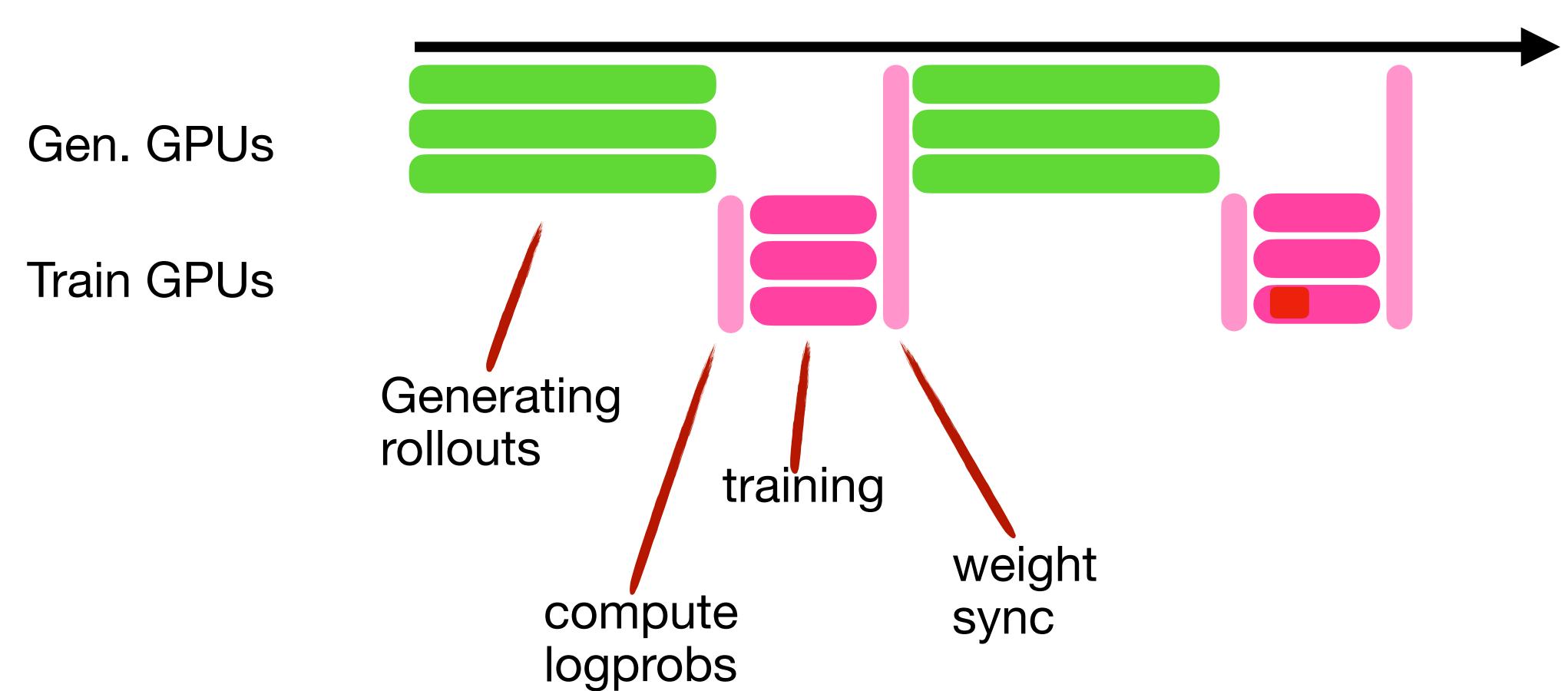
Even more recently: maybe theres just bugs in the inference code????

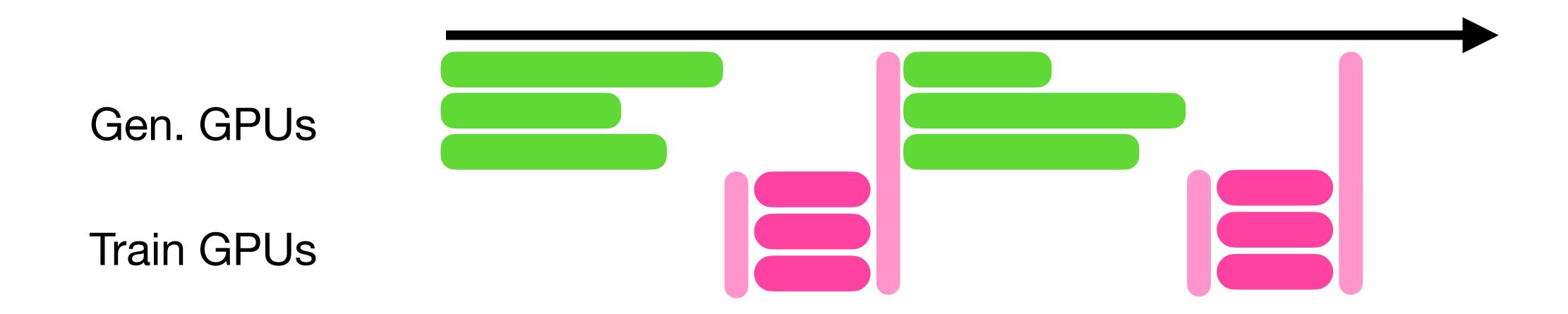


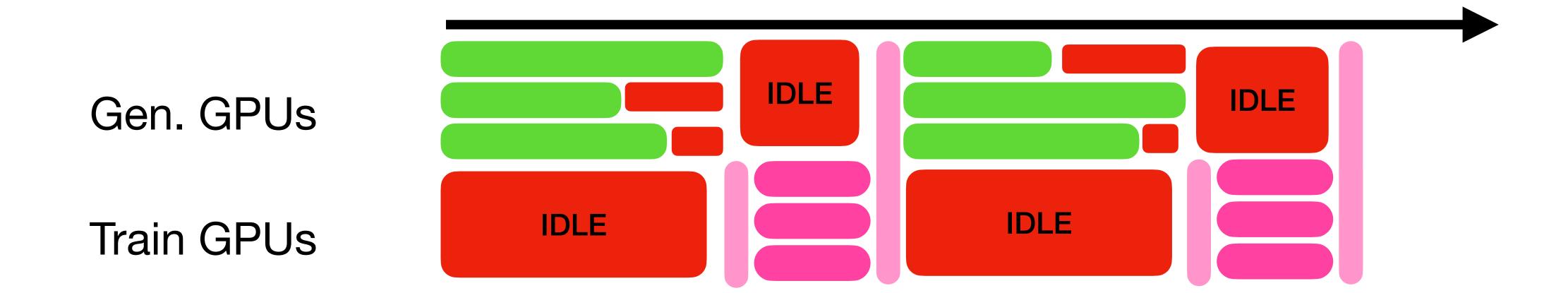
Let's review the standard algorithm for PPO/REINFORCE:

```
while True: Generate a full episode ("rollout") starting from s_0 using \pi_{\theta} Compute reward over the rollouts Compute the loss and apply gradient (potentially mult. steps)
```

We can visualise this flow like:





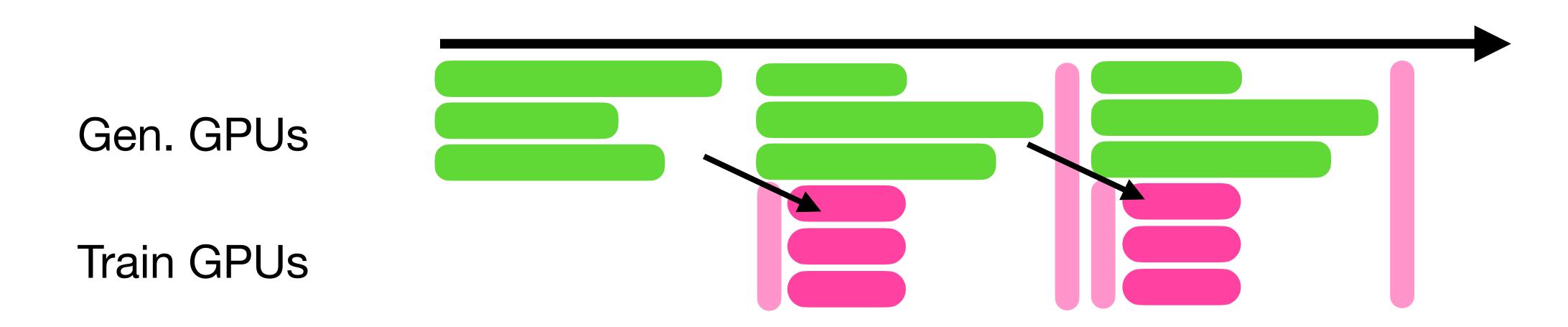


Fix 1: Overlap training and inference! This also makes our overall algorithm faster if we have enough compute!



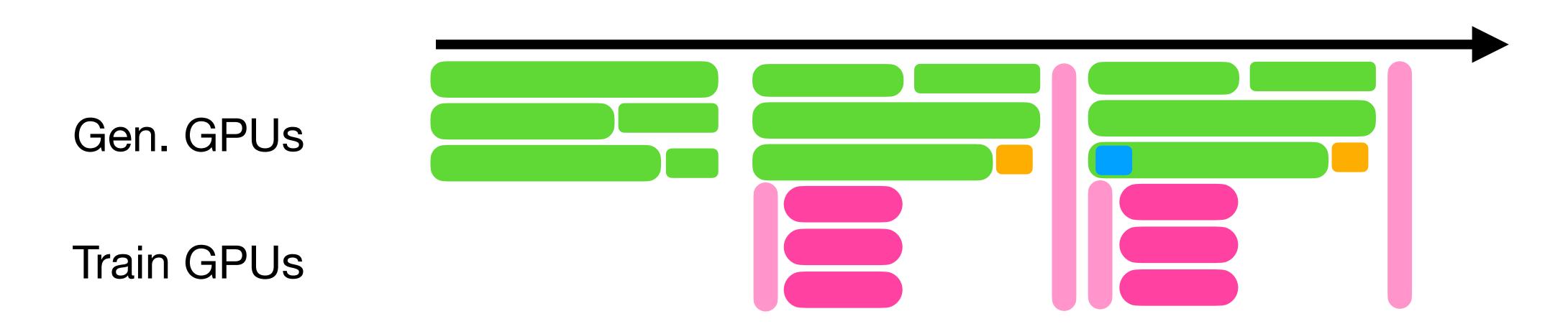
Drawback: our inference is now 1 step off our training (off-policy). In practice, this doesn't hurt until > 8 steps. Applying importance sampling also helps!

Fix 1: Overlap training and inference! This also makes our overall algorithm faster if we have enough compute!



Drawback: our inference is now 1 step off our training (off-policy). In practice, this doesn't hurt until > 8 steps. Applying importance sampling also helps!

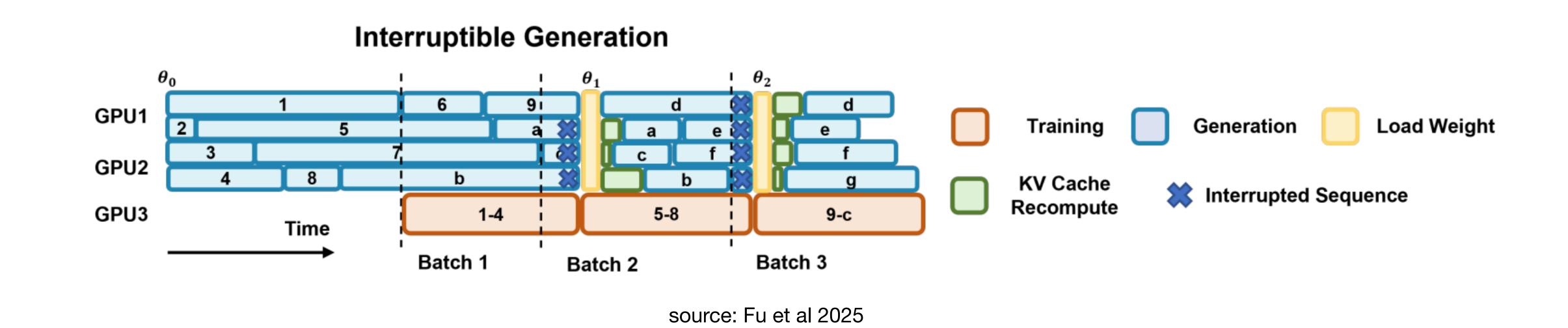
Fix 2: Overlap training and inference, and continuously generate samples when inference compute frees up.



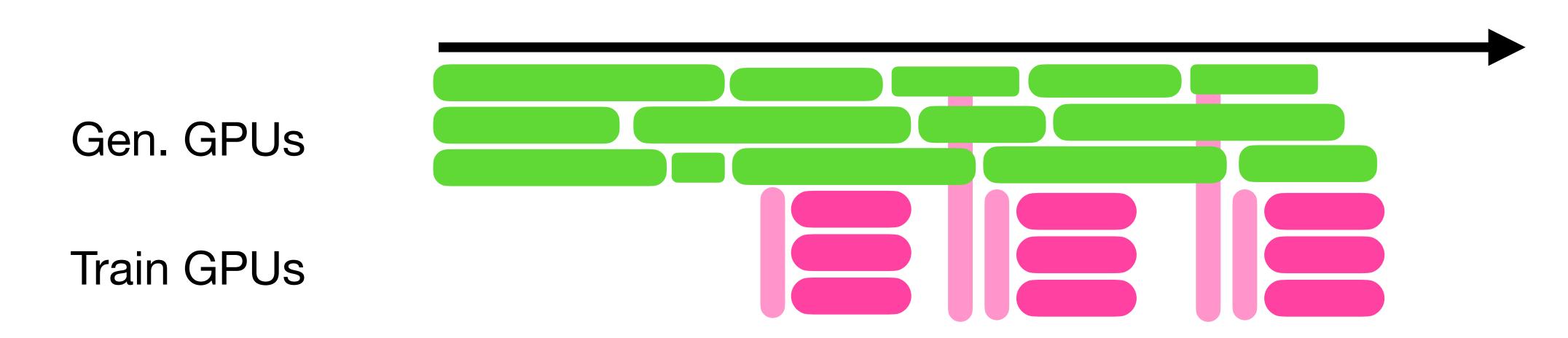
What do we do with interrupted generations? Just recompute KV cache and then continue.

Fix 2: Overlap training and inference, and continuously generate samples when inference compute frees up.

This is exactly what AReal does!



Fix 3: Overlap training and inference, and continuously generate samples when inference compute frees up, and don't recompute KV cache!

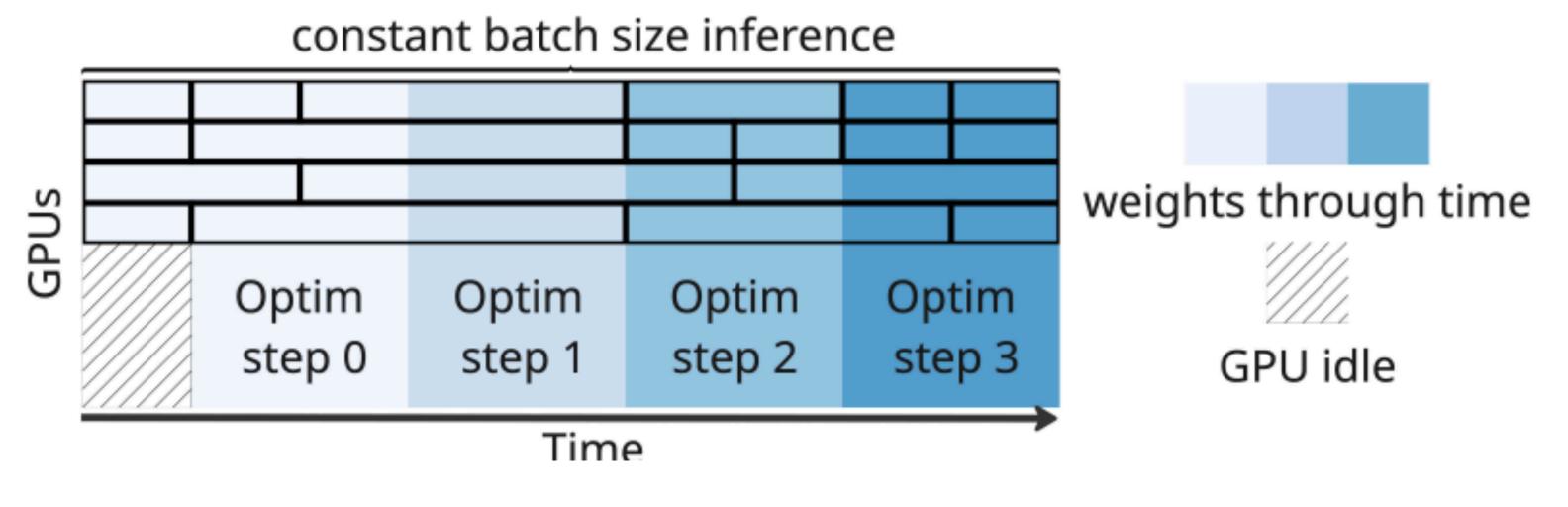


Run a training step whenever we have a batch available. Scalable wrt inference compute.

Speeding up inference: async RLHF

Fix 3: Overlap training and inference, and continuously generate samples when inference compute frees up, and don't recompute KV cache!

This technique was proposed by pipelineRL.

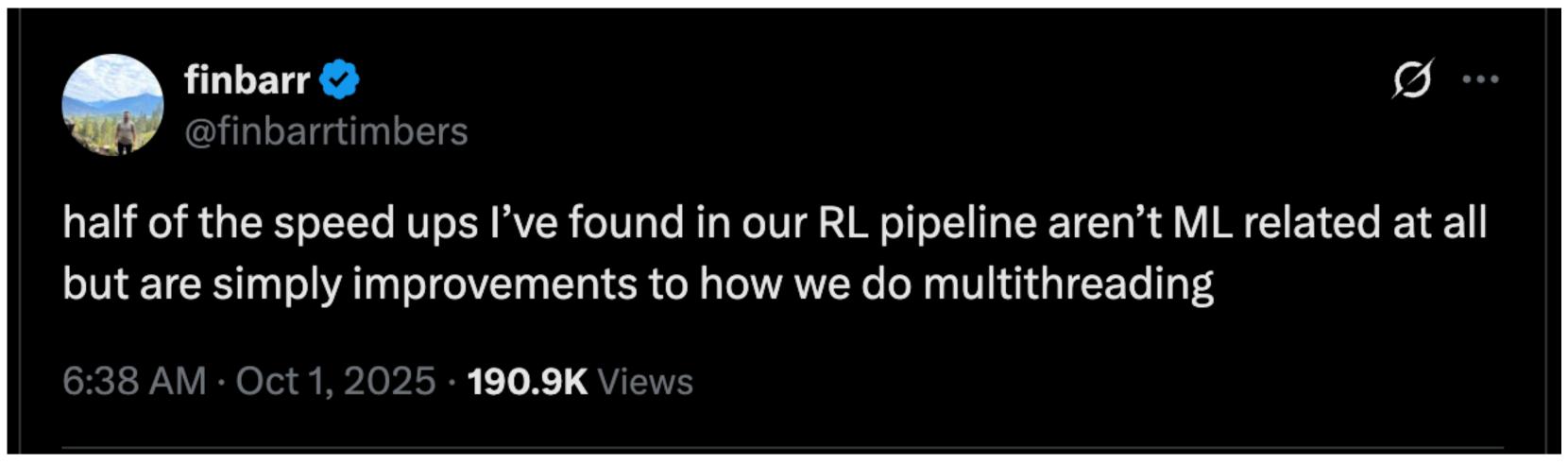


source: Piché et al 2025

Speeding up inference: async RLHF

Fix 3: Overlap training and inference, and continuously generate samples when inference compute frees up, and don't recompute KV cache!

This technique was proposed by pipelineRL. Requires being careful about multithreading!



source: https://x.com/finbarrtimbers/status/1973382120451060161

Adapting for LMs

- 1. LM inference is slow. How can we speed things up? 🔽
- 2. How do we deal with long sequences?
- 3. How do we deal with zero-advantage groups (GRPO)?

In practice, we have some maximum length we allow models to decode to.

This means some sequences will be unfinished.

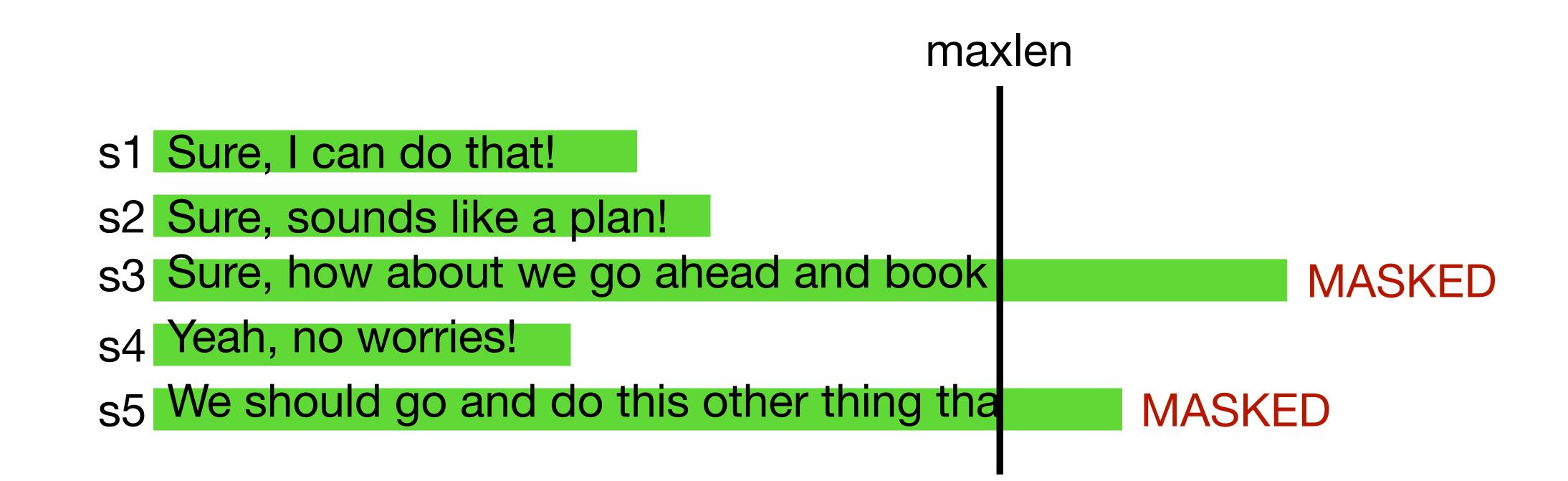
s1 Sure, I can do that! s2 Sure, sounds like a plan! s3 Sure, how about we go ahead and book s4 Yeah, no worries! s5 We should go and do this other thing tha

Vanilla: Treat and score as usual. Trains model to use only up to maxlen.

```
s1 Sure, I can do that!
s2 Sure, sounds like a plan!
s3 Sure, how about we go ahead and book
s4 Yeah, no worries!
s5 We should go and do this other thing tha
```

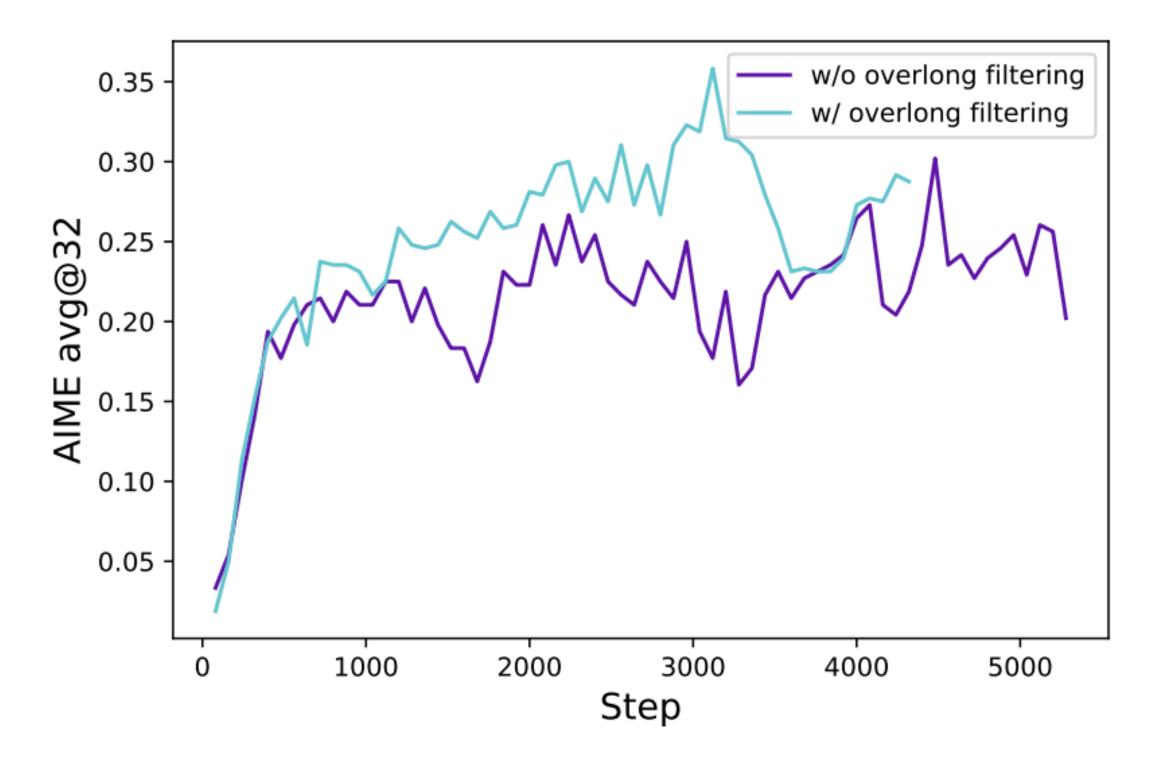
Vanilla: Treat and score as usual. Trains model to use only up to maxlen. **Overlong filtering (DAPO)**: Filter and remove samples that don't finish.

Prevents penalising model for strategies that might be correct, but didn't finish [debated].



Vanilla: Treat and score as usual. Trains model to use only up to maxlen. **Overlong filtering (DAPO)**: Filter and remove samples that don't finish.

Prevents penalising model for strategies that might be correct, but didn't finish [debated].



(a) Performance on AIME.

source: Yu et al. 2025

Vanilla: Treat and score as usual. Trains model to use only up to maxlen. **Overlong filtering (DAPO)**: Filter and remove samples that don't finish. Prevents penalising model for strategies that might be correct, but didn't finish [debated].

Interruptions (ScaleRL): Apply budget forcing and make the model stop (append "okay, now to stop thinking and give an answer</answer>" in 10-12k token range).

s1 Sure, I can do that! s2 Sure, sounds like a plan! s3 Sure, how about we go ahead </think> Yes! s4 Yeah, no worries! s5 We should go and do this </think> Good!

Length penalties can also be used to control LM generation length, but these aren't super widely applied in the open?

$$\text{gold len} \qquad \text{actual len}$$

$$r(y,y_{gold},n_{gold}) = \mathbb{I}(y=y_{gold}) - \alpha \cdot \big|n_{gold}-n_y\big|,$$

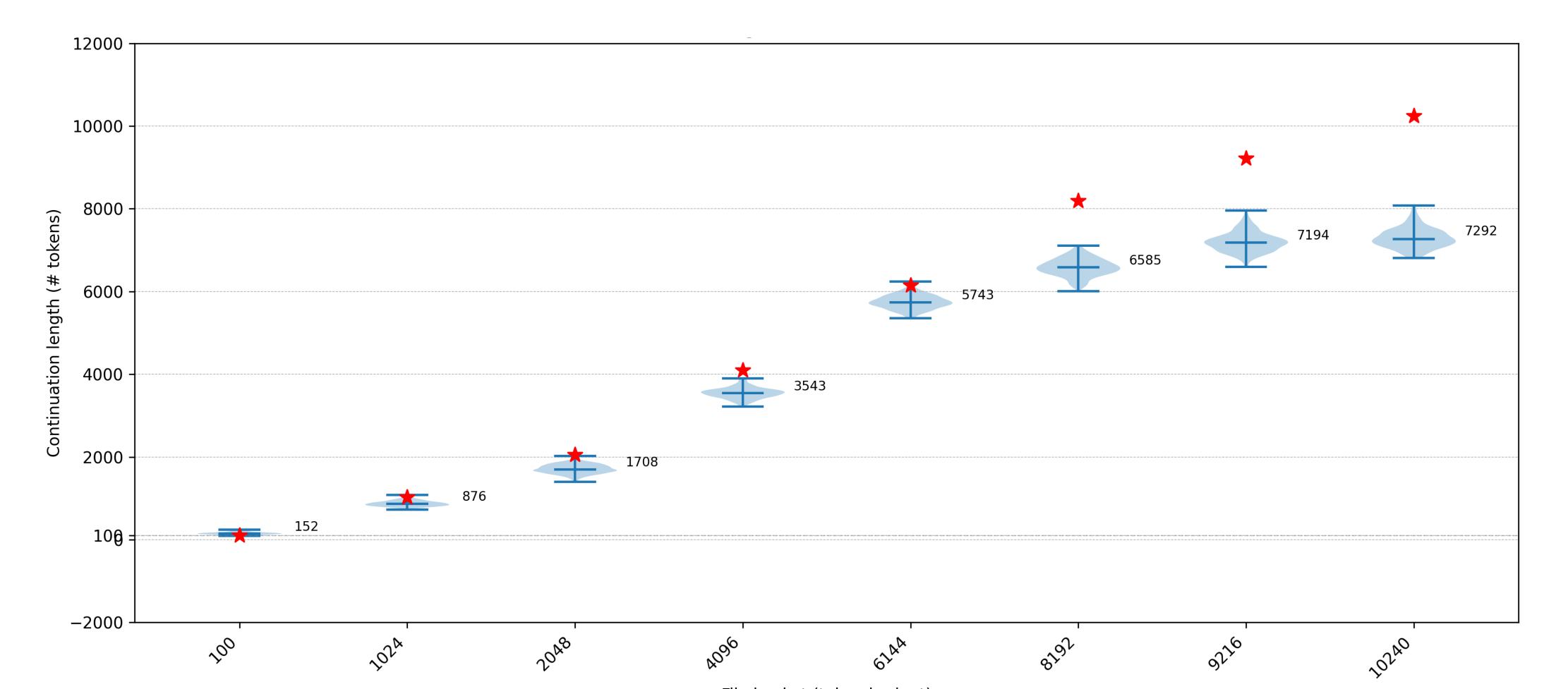
Basically subtract (or add) how far from "gold" length we are!

Published as a conference paper at COLM 2025

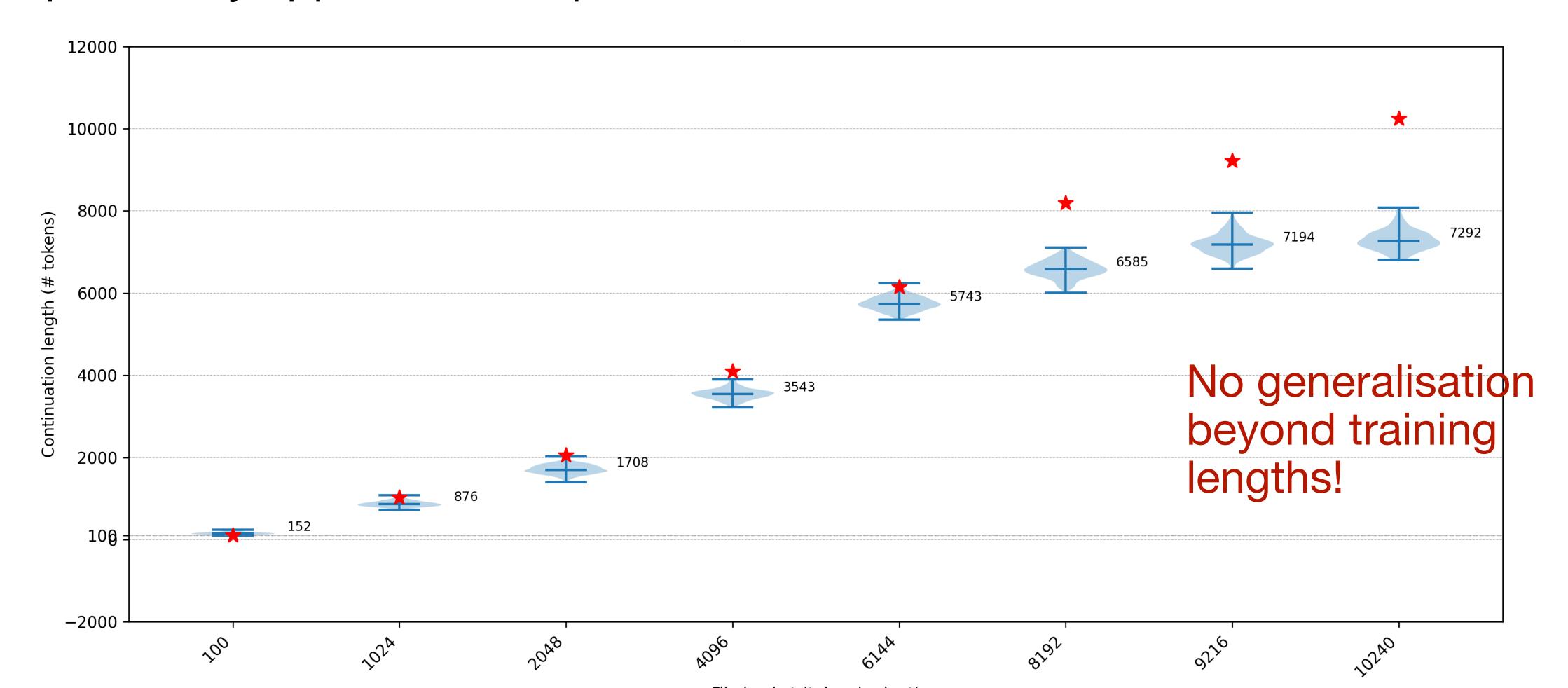
L1: Controlling How Long A Reasoning Model Thinks With Reinforcement Learning

Pranjal Aggarwal & Sean Welleck
Carnegie Mellon University
{pranjala, swelleck}@cs.cmu.edu

Length penalties can also be used to control LM generation length, but these aren't super widely applied in the open?



Length penalties can also be used to control LM generation length, but these aren't super widely applied in the open?

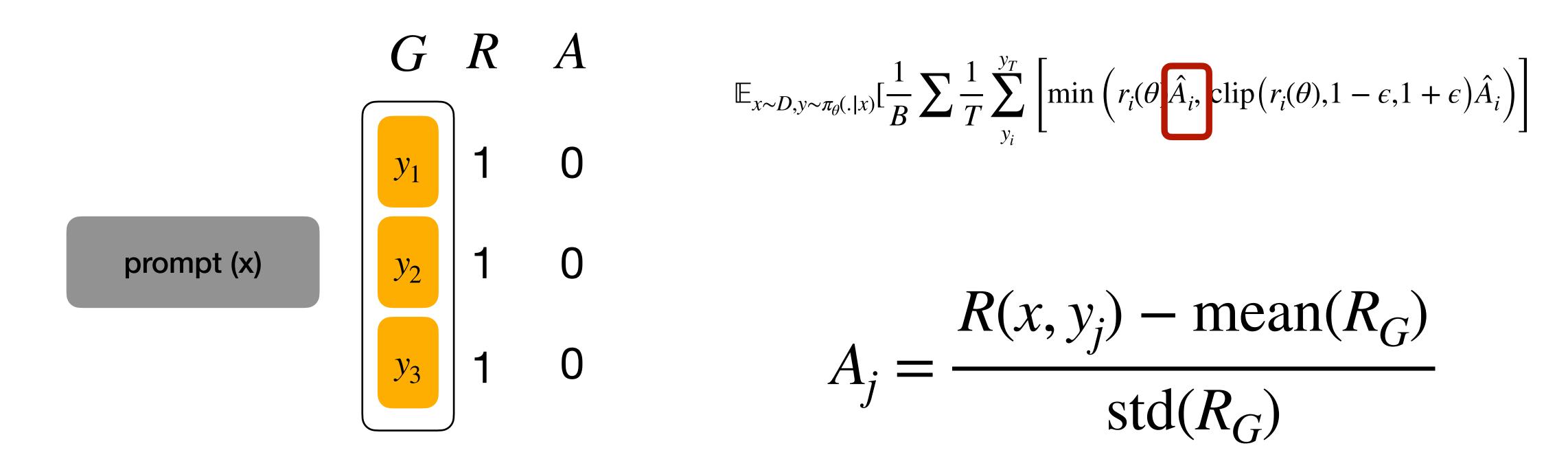


Adapting for LMs

- 1. LM inference is slow. How can we speed things up? 🔽
- 2. How do we deal with long sequences?
- 3. How do we deal with zero-advantage groups (GRPO)?

Dealing with never getting it right.

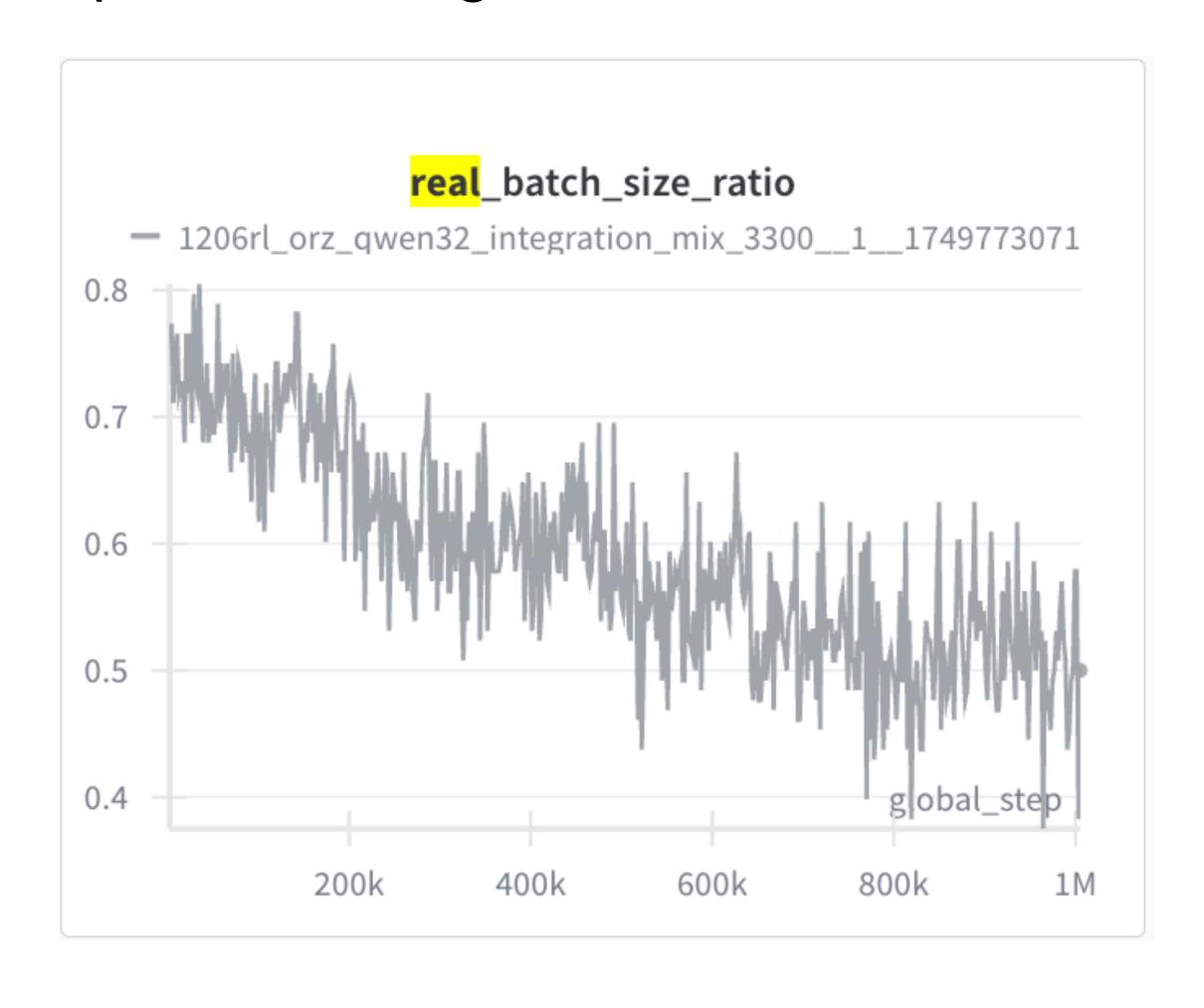
With GRPO, we might get all advantages in a group being 0 if they get the same score!



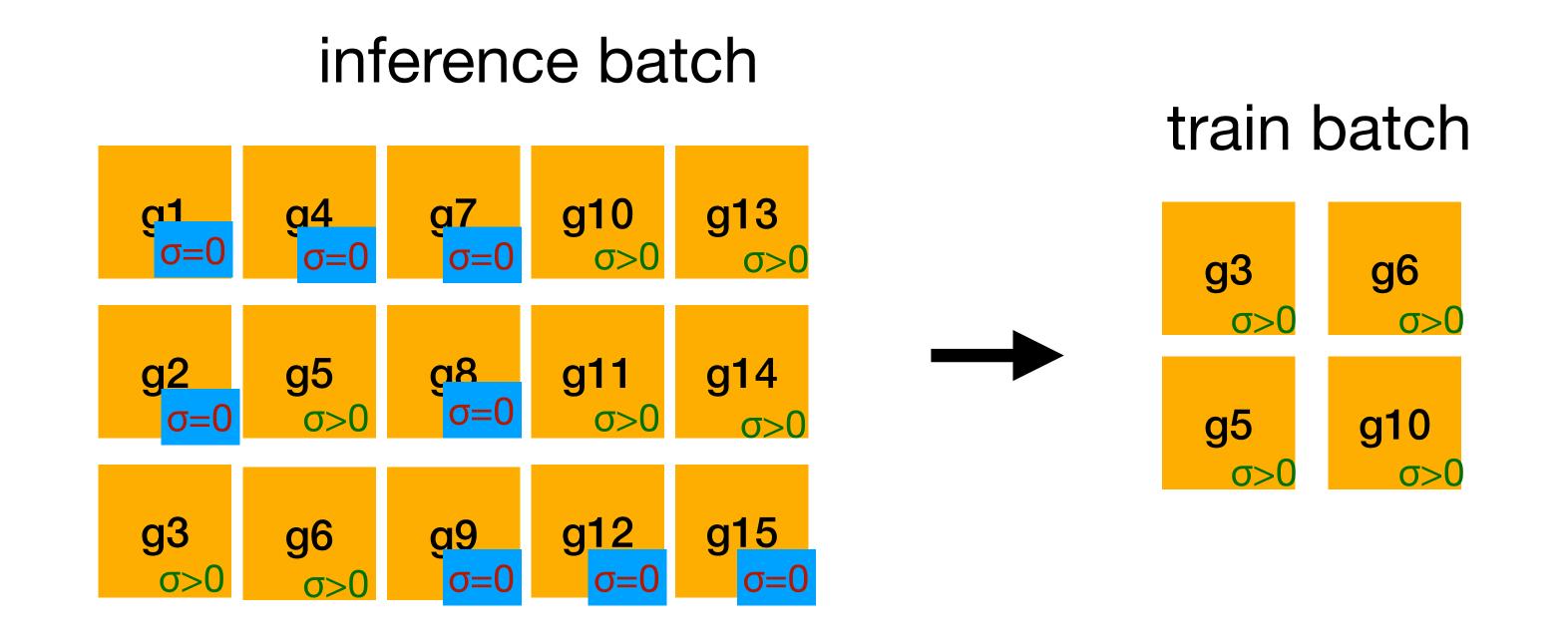
Solution: just throw these samples out!

Throwing samples out has issues...

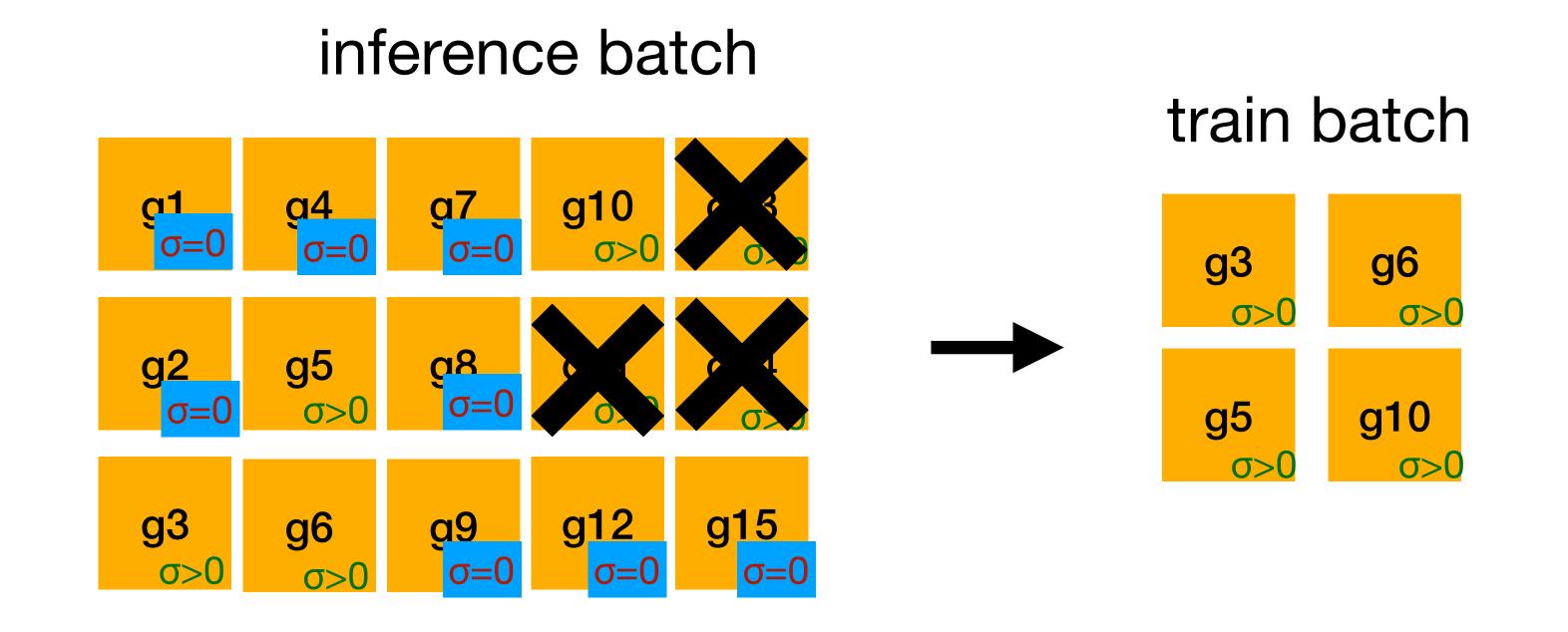
If we throw these samples out, we get a reduced batch size over time...



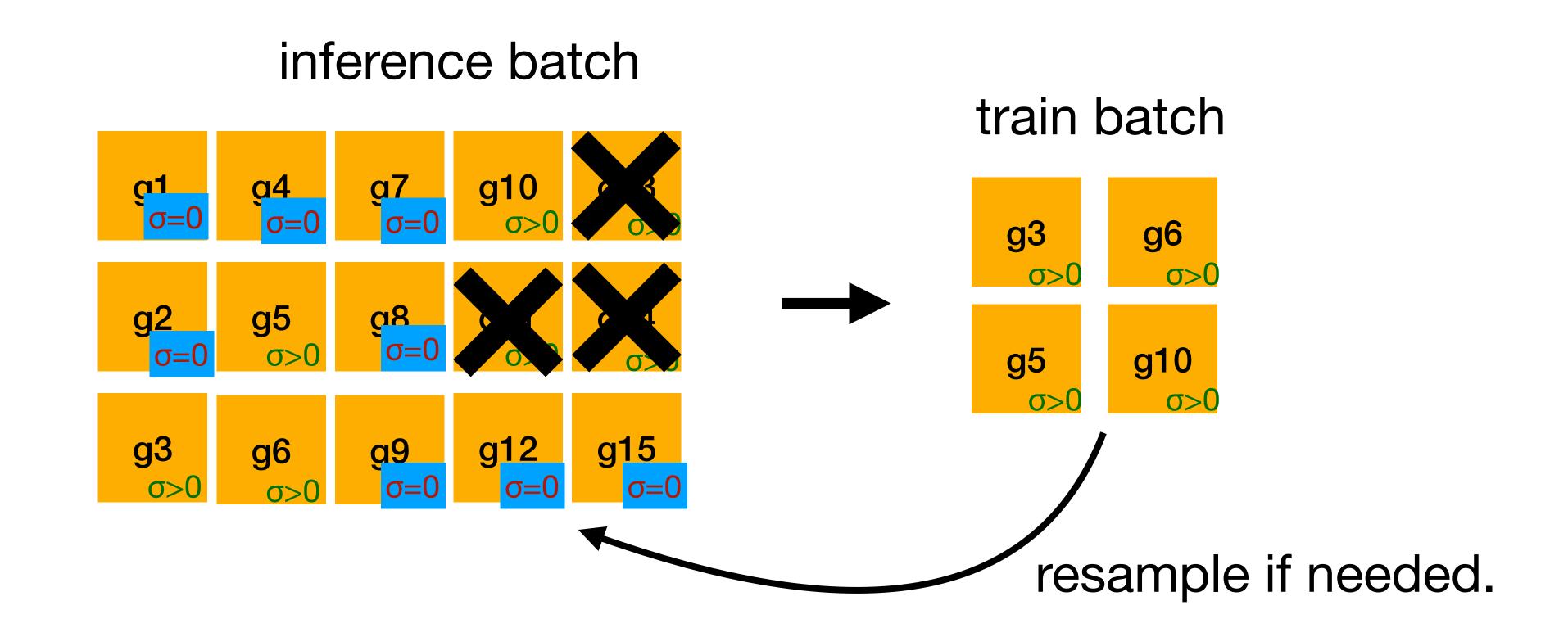
Fix 1: Generate more samples than we need, and downsample to form a batch ("dynamic sampling", DAPO)



Fix 1: Generate more samples than we need, and downsample to form a batch ("dynamic sampling", DAPO)



Fix 1: Generate more samples than we need, and downsample to form a batch ("dynamic sampling", DAPO)

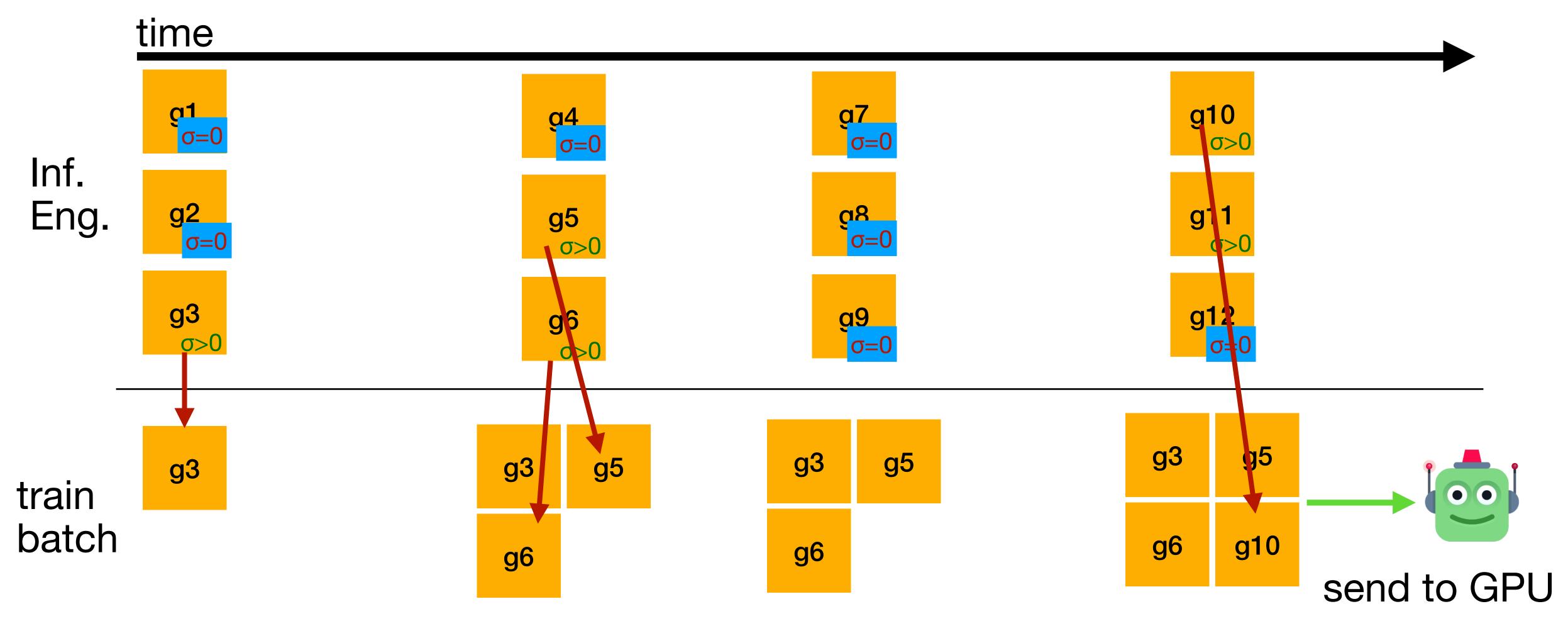


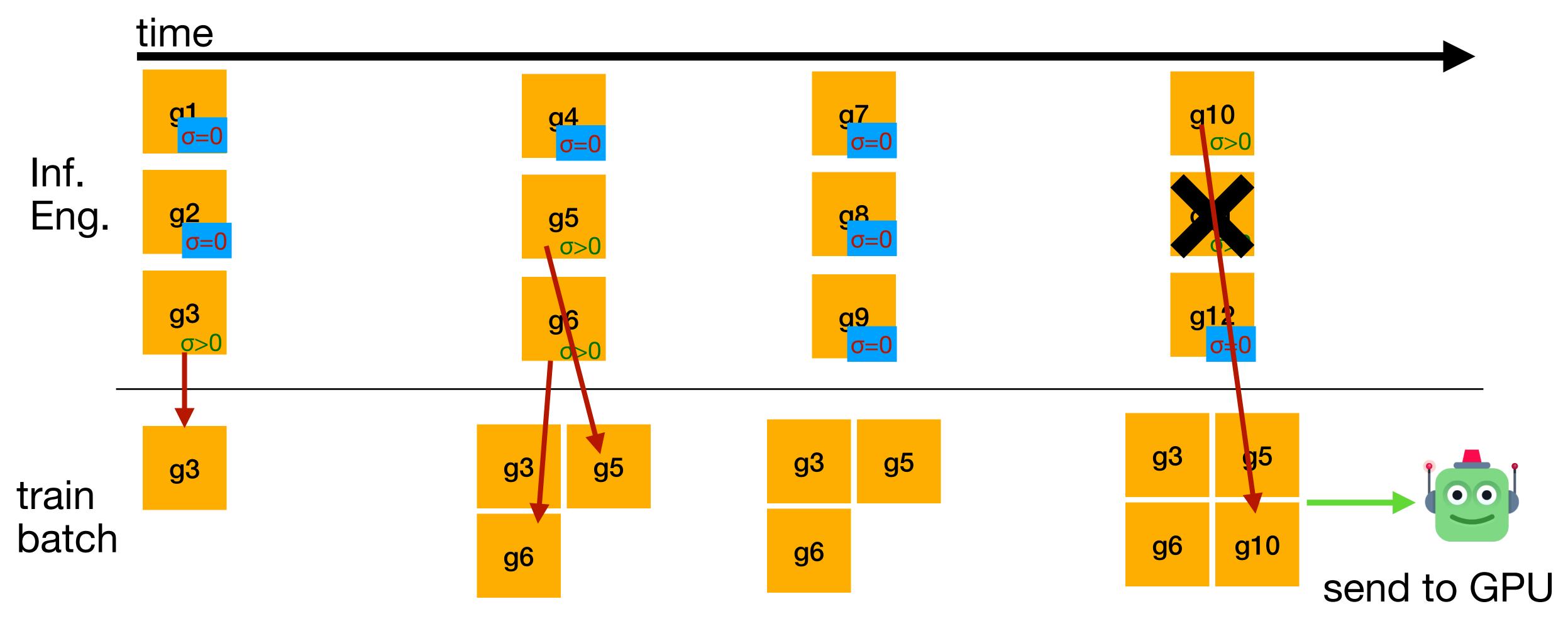
Fix 1: Generate more samples than we need, and downsample to form a batch ("dynamic sampling", DAPO)

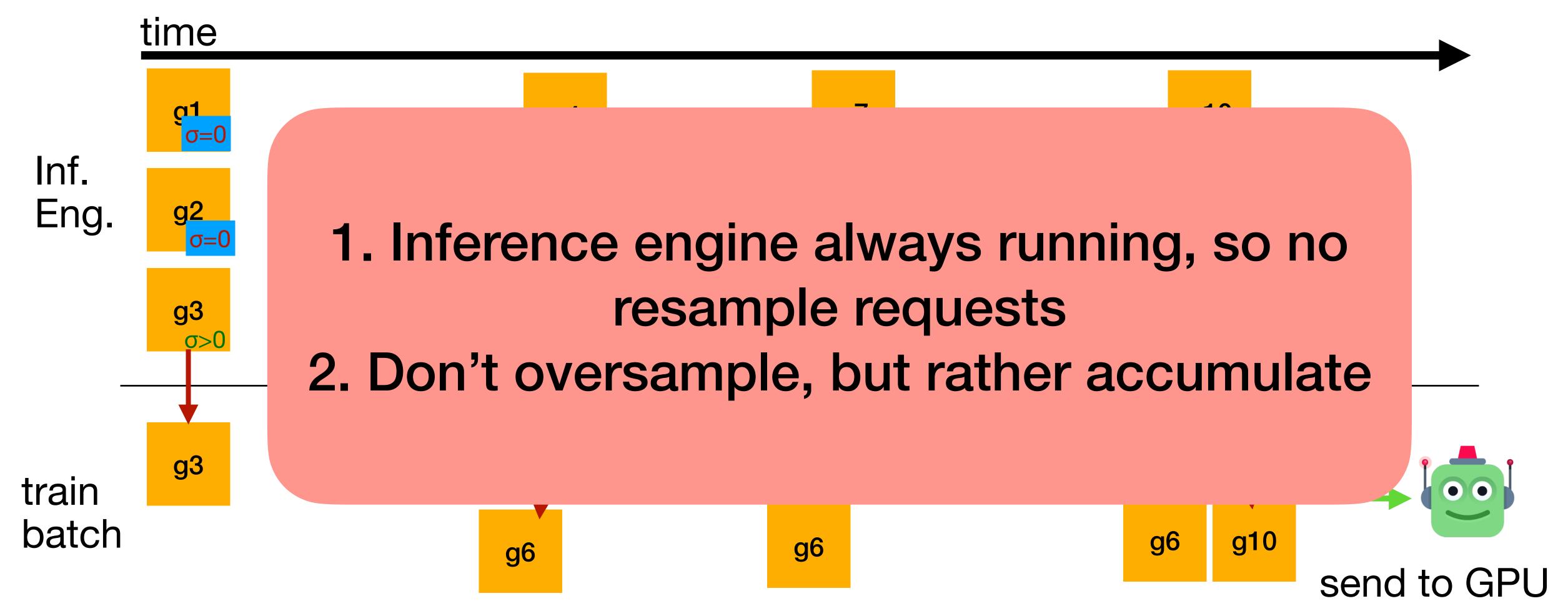
```
Here, over_sampling_batch_size needs to be greater than rollout_batch_size. For example:

--rollout-batch-size 32 \
--n-samples-per-prompt 8 \
--over-sampling-batch-size 64 \
```

source: https://thudm.github.io/slime/examples/glm4-9B.html







Adapting for LMs

- 1. LM inference is slow. How can we speed things up?
- 2. How do we deal with long sequences?
- 3. How do we deal with zero-advantage groups (GRPO)?

Putting it all together.

You now know everything!

Many modern papers are just a combination of the things I have just covered.

DAPO

- use clip-higher in importance clipping
- filter out too long generations
- oversample prompts to avoid small batch
- token-level loss

ByteDance | Seed



DAPO: An Open-Source LLM Reinforcement Learning
System at Scale

¹ByteDance Seed ²Institute for AI Industry Research (AIR), Tsinghua University

³The University of Hong Kong

⁴SIA-Lab of Tsinghua AIR and ByteDance Seed

Full author list in Contributions



You now know everything!

Many modern papers are just a combination of the things I have just covered.

Dr GRPO

- use sum-and-constant loss
- remove std dev from GRPO advantage

Understanding R1-Zero-Like Training: A Critical Perspective

Zichen Liu*+1,2, Changyu Chen*1,3, Wenjun Li*3, Penghui Qi*1,2, Tianyu Pang¹, Chao Du¹, Wee Sun Lee², Min Lin¹
¹Sea AI Lab



²National University of Singapore

³Singapore Management University

You now know everything!

Many modern papers are just a combination of the things I have just covered.

Minimax Paper

- Clip importance ratio only (CISPO)
- Clip-higher
- FP32 precision for logprob mismatch



MiniMax-M1: Scaling Test-Time Compute Efficiently with Lightning Attention

MiniMax¹



How much does all this matter? ScaleRL knows!

Khatri et al., 2025

Read this, and you pretty much know most of what you need to know these days.

The Art of Scaling Reinforcement Learning Compute for LLMs

```
Devvrit Khatri<sup>2,*,†</sup>, Lovish Madaan<sup>1,3,*</sup>, Rishabh Tiwari<sup>4†</sup>, Rachit Bansal<sup>5†</sup>, Sai Surya Duvvuri<sup>2†</sup>, Manzil Zaheer<sup>1†</sup>, Inderjit S. Dhillon<sup>2</sup>, David Brandfonbrener<sup>1</sup>, Rishabh Agarwal<sup>6,†</sup>

<sup>1</sup>Meta, <sup>2</sup>UT Austin, <sup>3</sup>UCL, <sup>4</sup>UC Berkeley, <sup>5</sup>Harvard University, <sup>6</sup>Periodic Labs *Equal contribution, <sup>†</sup>Work done at Meta
```

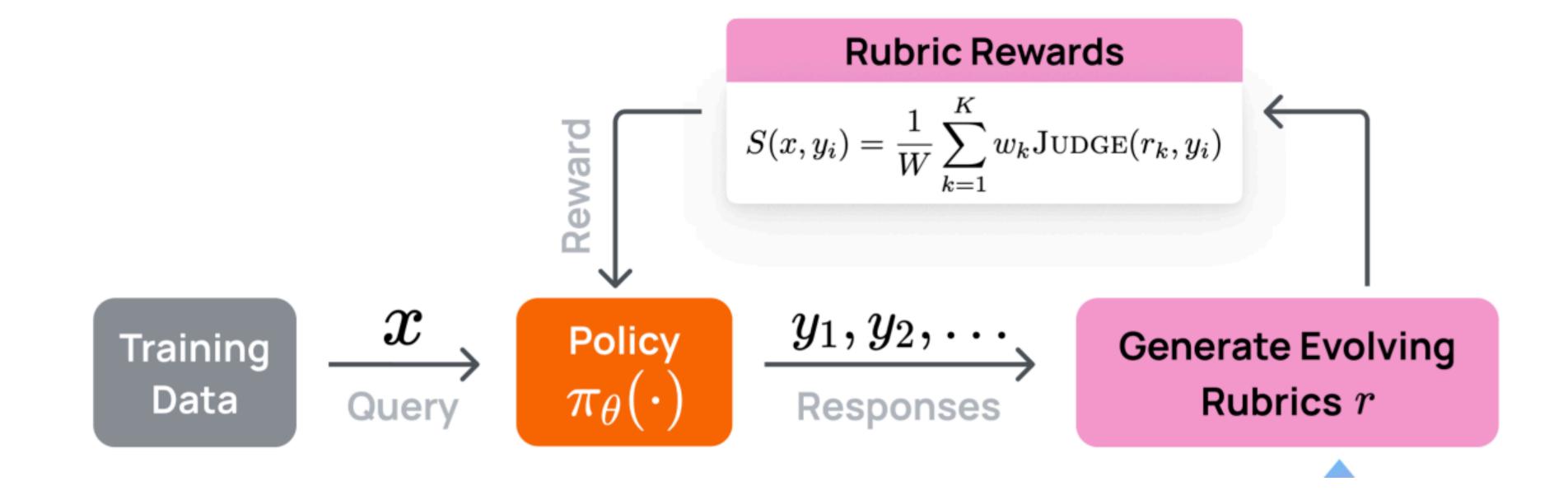


```
Context: Pelosi, the first and only female House Speaker, has served her San Francisco district for decades and capped a long career with this decis
ion. She described it as the end of a decades-long career and
highlighted her role leading the House.
Searched: 10 links | Browsed: 4 links | Tool calls: 6 | Failed: 1
/ou: exit
Goodbye!
Stopping Search vLLM server...
Stopping MCP server...
root@saturn-cs-aus-246:/weka/oe-adapt-default/hamishi/dr-tulu/agent# uv run --extra vllm python scripts/launch_chat.py --model rl-research/DR-Tulu-8E
=== Interactive Chat Launcher ===

    MCP server is not running on port 8000

_aunch MCP server? (y/n): y
🖋 Launching MCP server on port 8000...
MCP server output will be logged to /tmp/mcp_server_8000.log
🔀 Waiting for MCP server to start...
△ MCP server process started but port check failed. Continuing anyway...
	riangle vLLM server does not appear to be accessible at http://localhost:30004/v1
🖋 Auto-launching Search Agent vLLM server: rl-research/DR-Tulu-8B on port 30004 (GPU 0)...
Launching vLLM server for model rl-research/DR-Tulu-8B on port 30004...
vLLM output for rl-research/DR-Tulu-8B will be logged to /tmp/vllm_server_30004.log
\overline{\mathbb{X}} Waiting for vLLM server to become ready (this may take a few minutes)...
🗶 Still waiting for vLLM server (30s)...
/ vLLM server started (PID: 15462)
                                                                                                           B
🖋 Starting interactive chat...
```

© 18/11, 12:17 □ ~/Programming/dr-tulu/rl % beaker - -zsh



User query

Describe how genetically engineering T cells could be used as an anti-inflammatory therapy for IBD by changing cytokine signaling between T cells and macrophages.

 \Box \downarrow Generate rubrics with search

Persistent rubrics

- Cites "IL-10-engineered T cells **reduced colitis severity**"
- States "engineering T cells to overexpress IL-10 suppresses macrophage TNF-α"



policy at step t

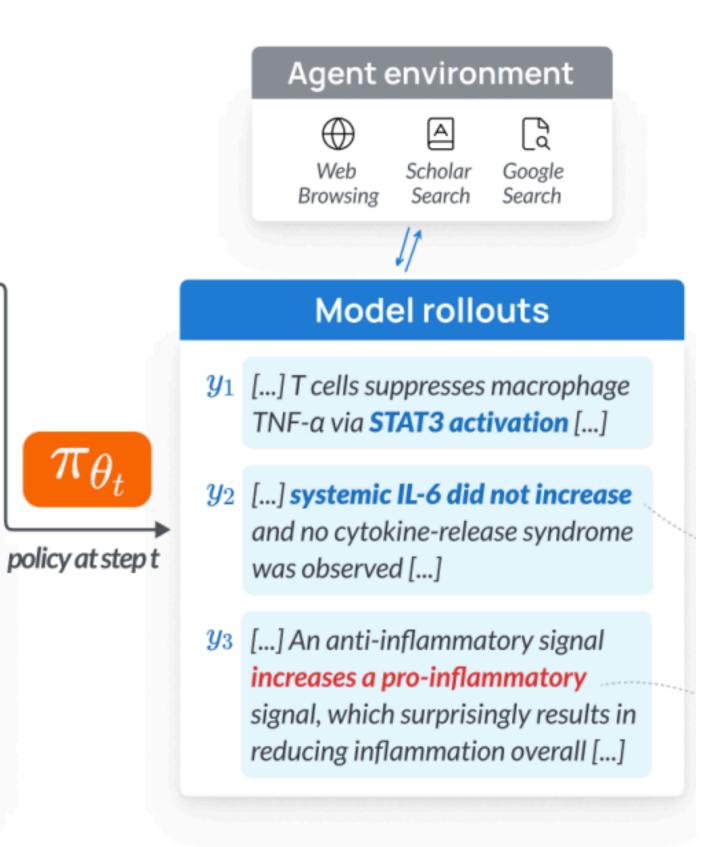
User query

Describe how genetically engineering T cells could be used as an anti-inflammatory therapy for IBD by changing cytokine signaling between T cells and macrophages.

 \Box \downarrow Generate rubrics with search

Persistent rubrics

- Cites "IL-10-engineered T cells **reduced colitis severity**"
- States "engineering T cells to overexpress IL-10 suppresses macrophage TNF-α"



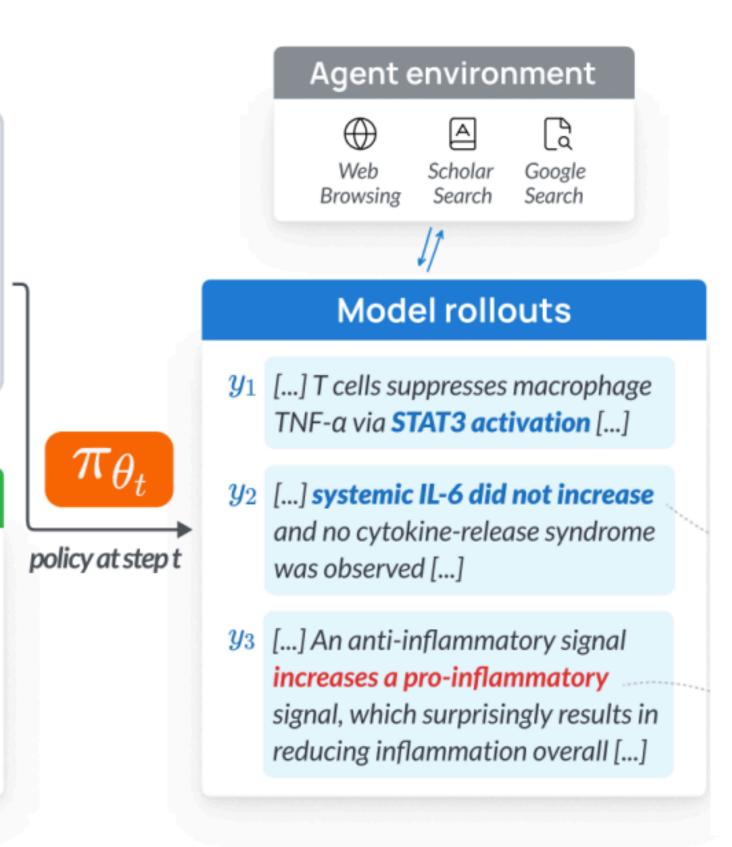
User query

Describe how genetically engineering T cells could be used as an anti-inflammatory therapy for IBD by changing cytokine signaling between T cells and macrophages.

 \Box \downarrow Generate rubrics with search

Persistent rubrics

- Cites "IL-10-engineered T cells **reduced colitis severity**"
- States "engineering T cells to overexpress IL-10 suppresses macrophage TNF-α"



Rubric-based scoring

$$S(x,y) = \sum_{k=1}^{K} w_k \cdot \text{JUDGE}(r_k, y)$$

o Updates $\pi_{ heta_{t+1}}$

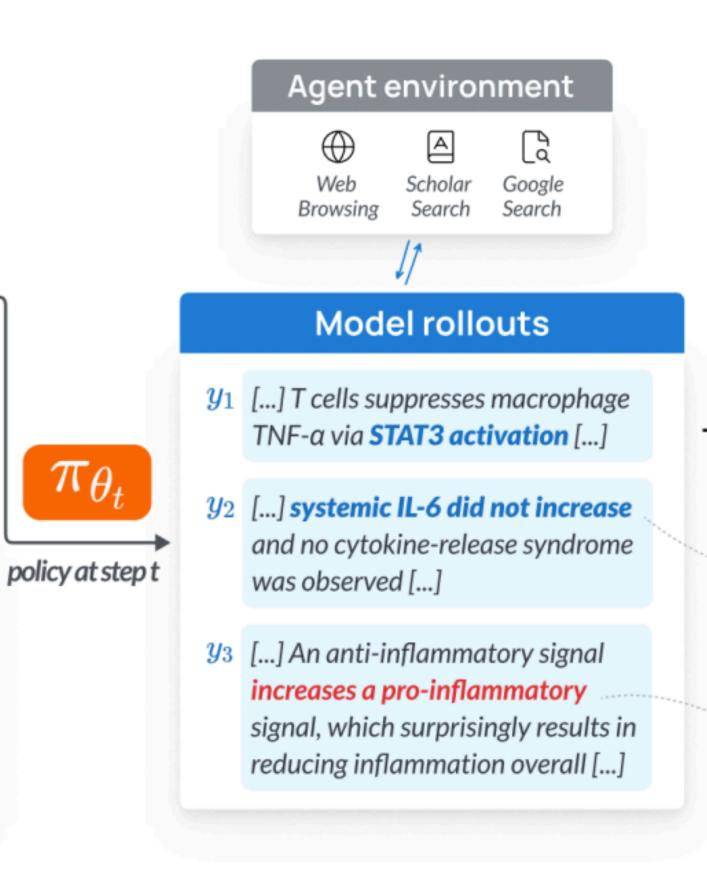
User query

Describe how genetically engineering T cells could be used as an anti-inflammatory therapy for IBD by changing cytokine signaling between T cells and macrophages.

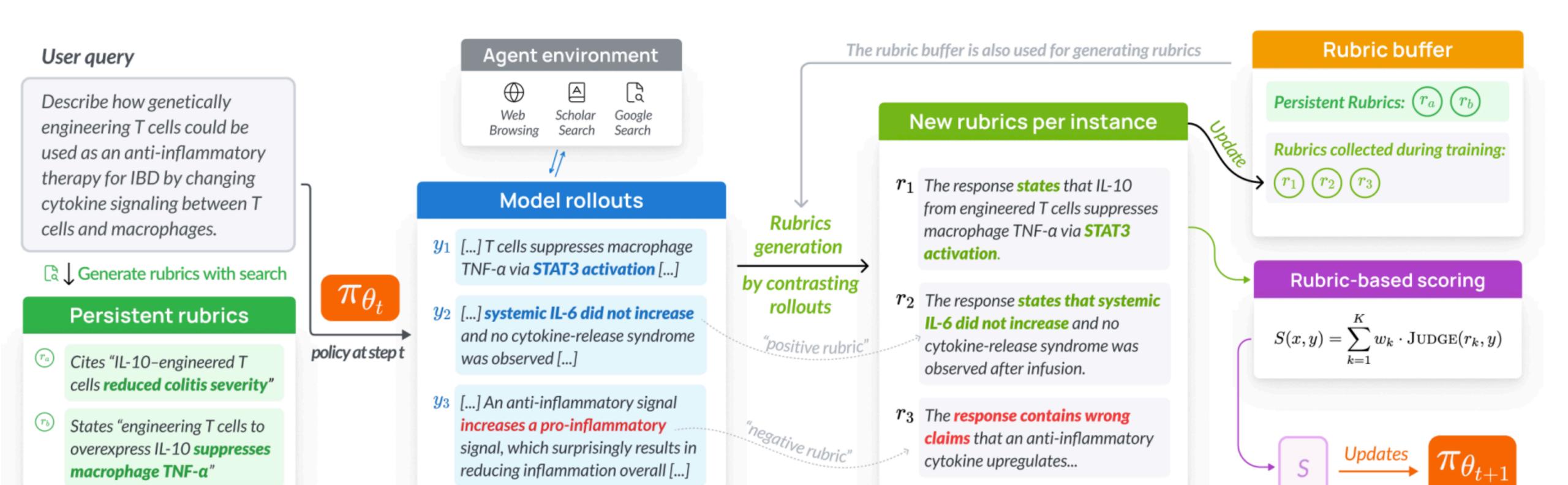
 \Box \downarrow Generate rubrics with search

Persistent rubrics

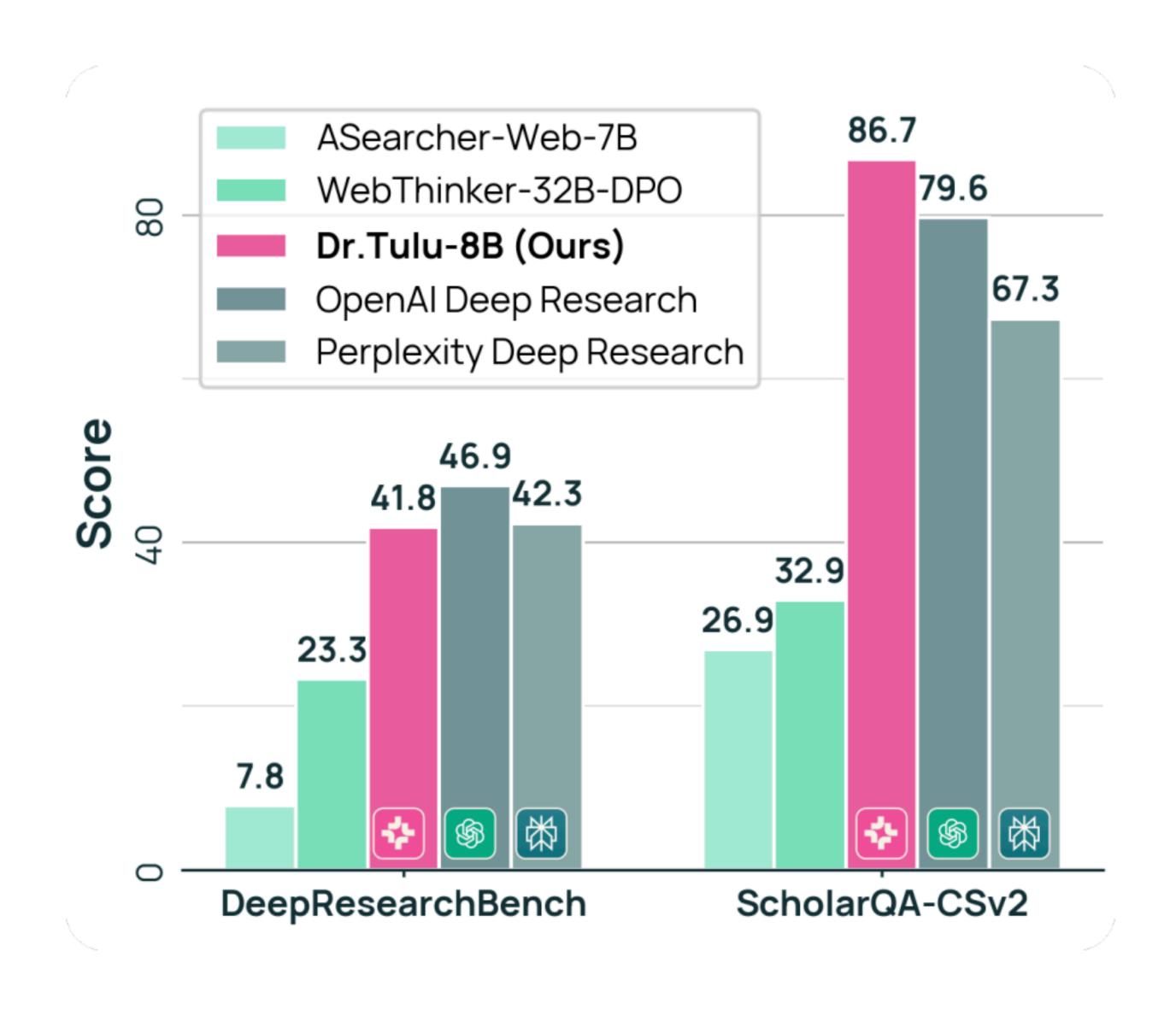
- Cites "IL-10-engineered T cells **reduced colitis severity**"
- States "engineering T cells to overexpress IL-10 suppresses macrophage TNF-a"



The rubric buffer is also used for generating rubrics New rubrics per instance r_1 The response **states** that IL-10 from engineered T cells suppresses Rubrics macrophage TNF-α via **STAT3** generation activation. Rubric-based scoring by contrasting r_2 The response states that systemic rollouts **IL-6 did not increase** and no $S(x,y) = \sum w_k \cdot \text{Judge}(r_k, y)$ cytokine-release syndrome was "positive rubric" observed after infusion. r_3 The response contains wrong hegative rubric" **claims** that an anti-inflammatory cytokine upregulates...

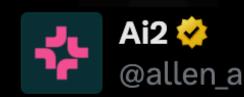


Practical example: DR Tulu



Practical example:





Announcing Olmo 3, a leading fully open LM suite built for reasoning, chat, & tool use, and an open model flow—not just the final weights, but the entire training journey.

Best fully open 32B reasoning model & best 32B base model.



6:04 AM · Nov 20, 2025 · 274.5K Views

Practical example: 401mo3

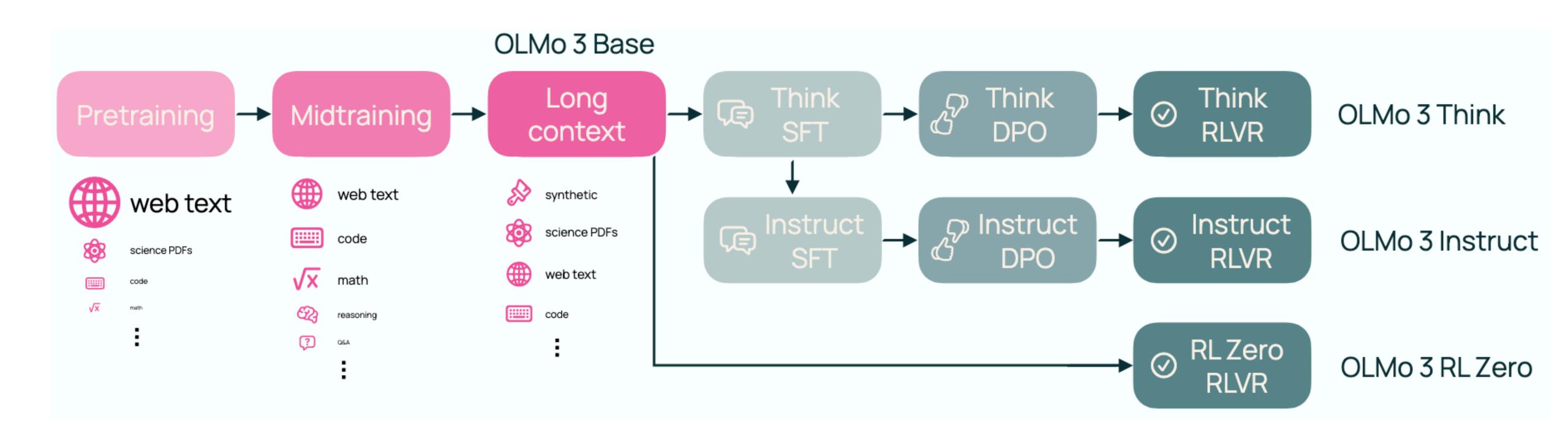


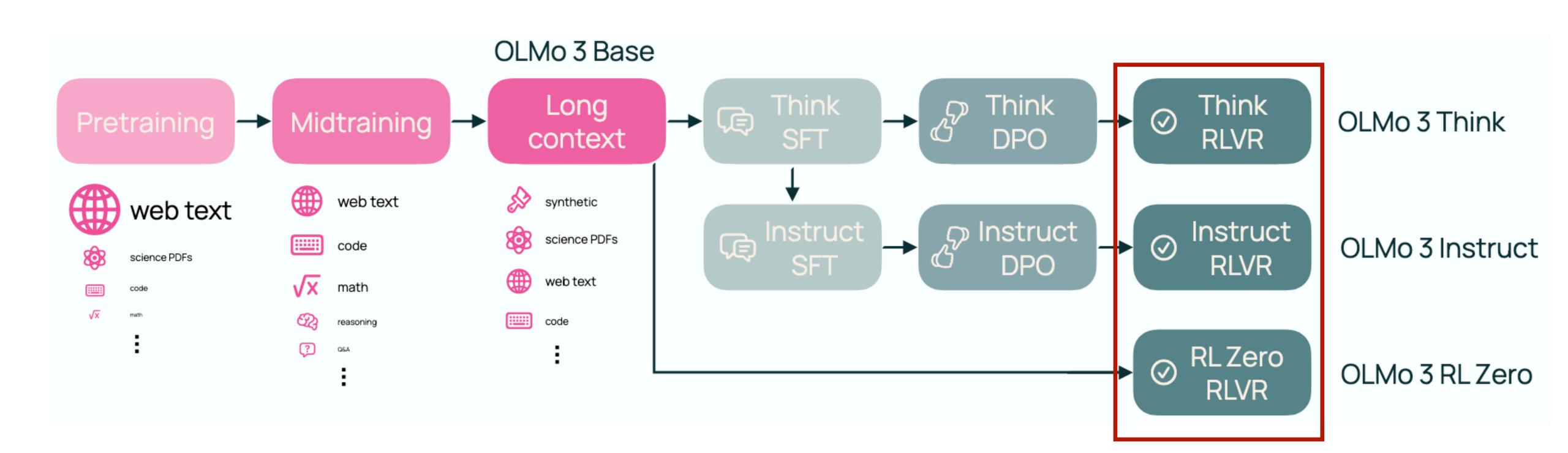
Announcing Olmo 3, a leading fully open LM suite built for reasoning, chat, & tool use, and an open model flow—not just the final weights, but the entire training journey.

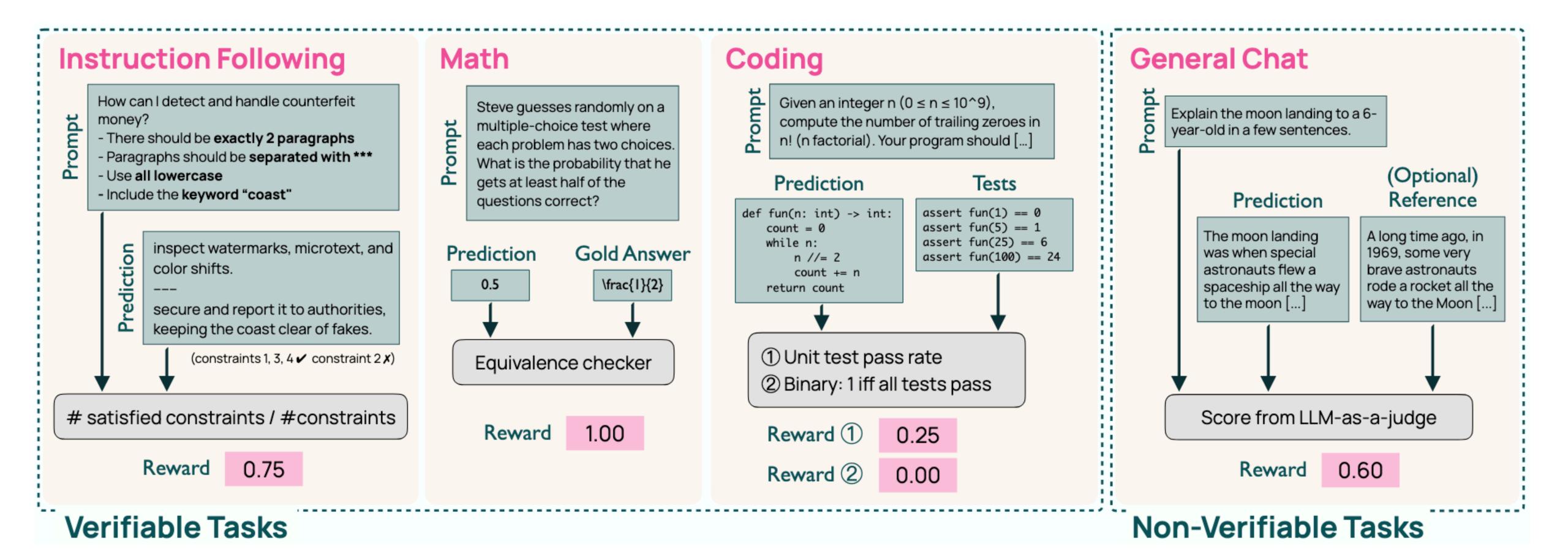
Best fully open 32B reasoning model & best 32B base model.



6:04 AM · Nov 20, 2025 **274.5K** Views





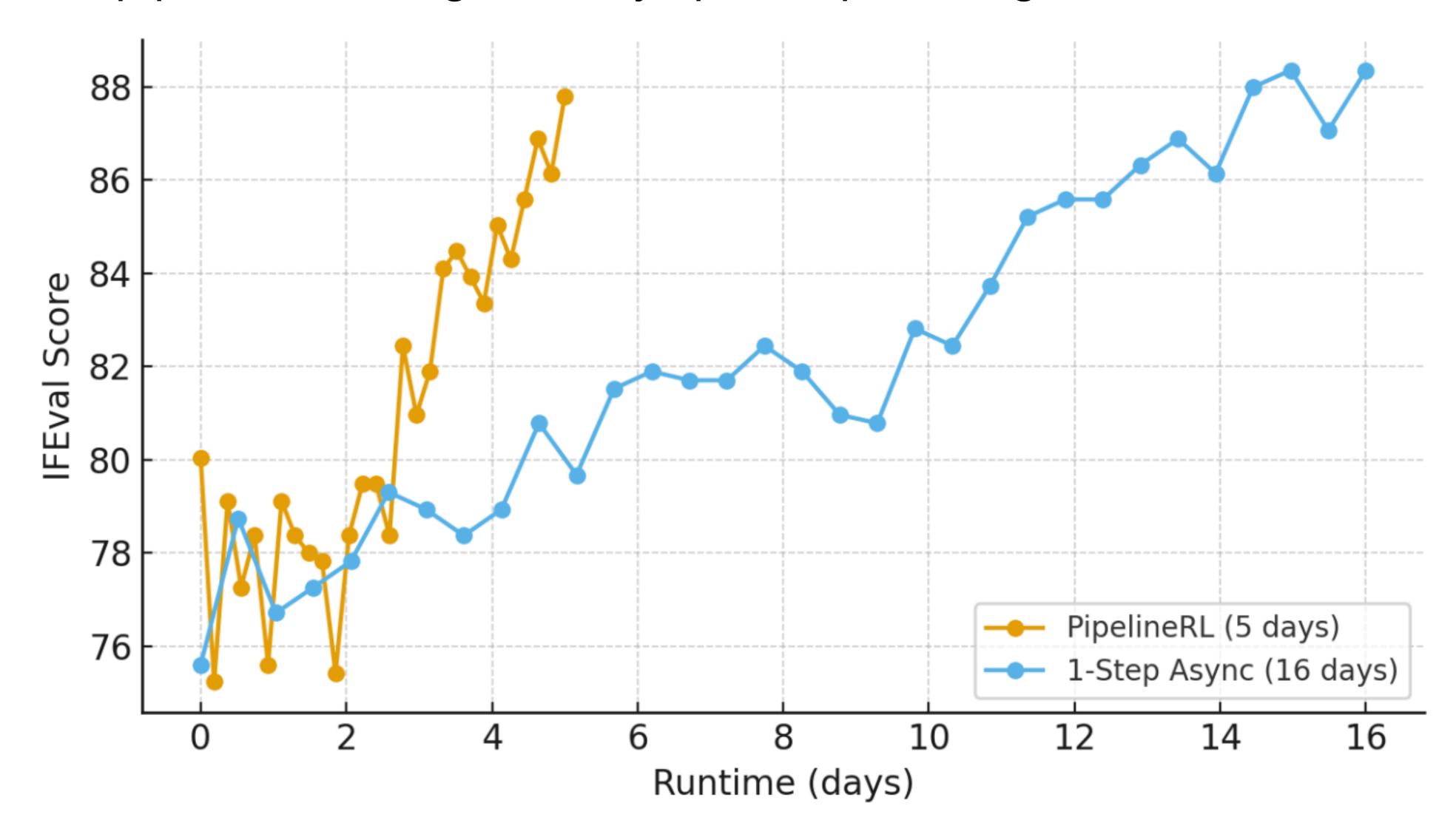


We make the following improvements over vanilla GRPO:²²

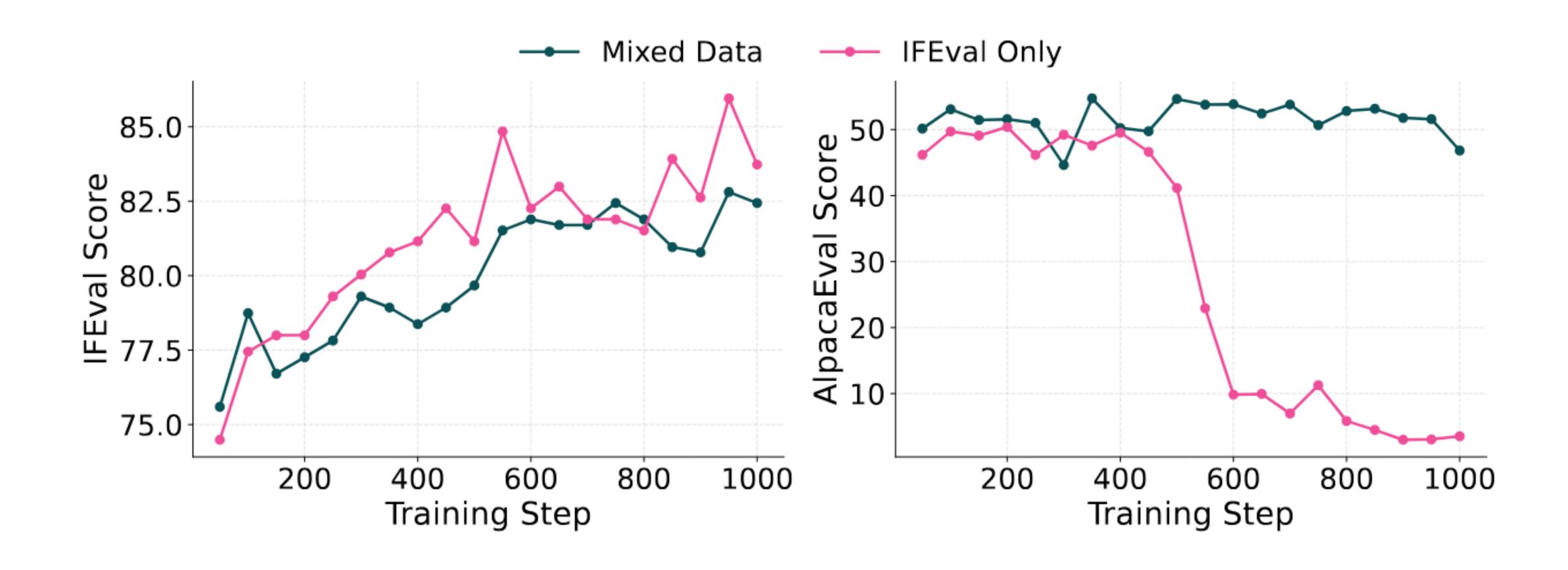
- **Zero Gradient Signal Filtering**: We remove groups of instances whose rewards are all identical (i.e., a batch with zero standard deviation in their advantage) to avoid training on samples that provide zero gradient, similar to DAPO (Yu et al., 2025).
- Active Sampling: We maintain a consistent batch size in spite of zero gradient filtering with a nove, more efficient version of dynamic sampling (Yu et al., 2025), see Olmore Infra for details.
- Token-level loss: We use a token-level loss to normalize the loss by the total number of tokens across the batch (Yu et al., 2025), rather than per-sample to avoid a length bias.
- No KL Loss We remove the KL loss as a common practice (GLM-4.5 Team et al., 2025; Yu et al., 2025; Liu et al., 2025b) as it allows less-restricted policy updates, and removing it does not lead to over-optimization or destabilized training.
- Clip Higher: We set the upper-bound clipping term in the loss to a slightly higher value than the lower bound to enable larger updates on tokens, as proposed by Yu et al. (2025).
- Truncated Importance Sampling: To adjust for differences between log probabilities from the inference and training engines, we multiply the loss by the truncated importance sampling ratio, following Yao et al. (2025).
- No standard deviation normalization: When calculating advantage, we do not normalize by the standard deviation of the group, following Liu et al. (2025b). This removes a difficulty bias, where questions with low standard deviation in their rewards (e.g. too hard or too easy) have their advantages significantly increased by the normalization term.

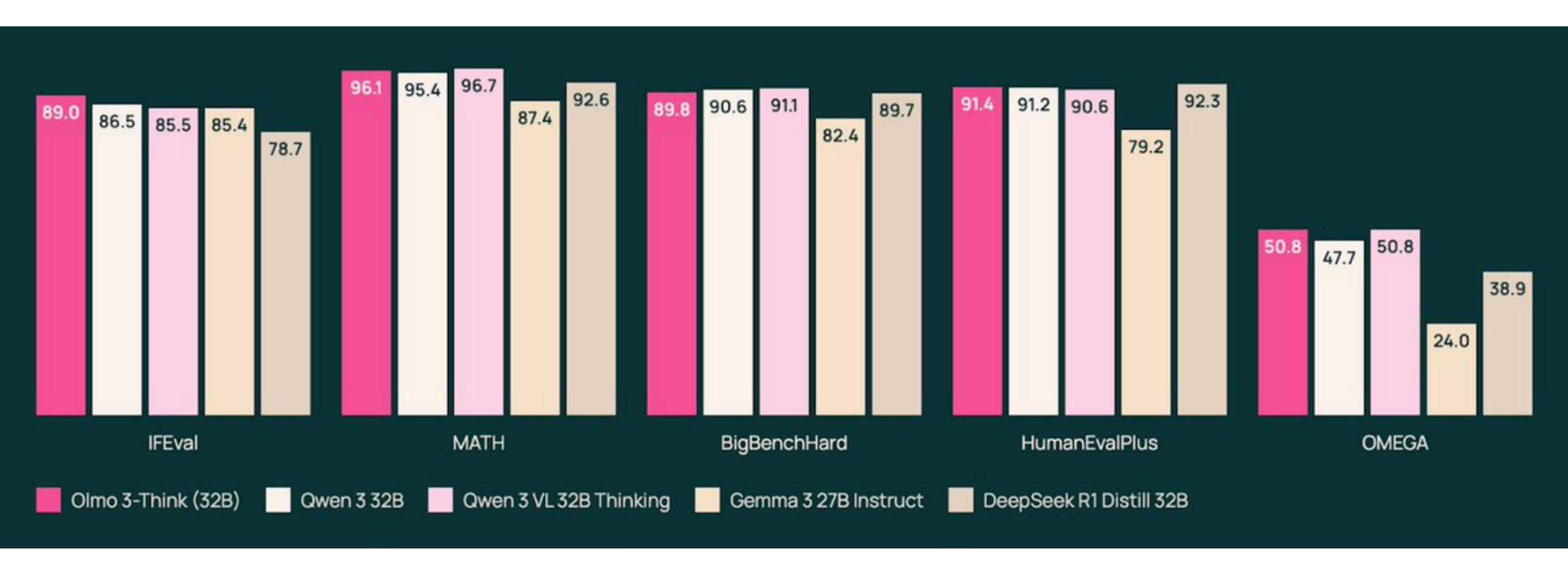
OlmoRL Formulation Our final objective function includes a token-level loss, truncated importance sampling, clip-higher, and no standard deviation in the advantage calculation:

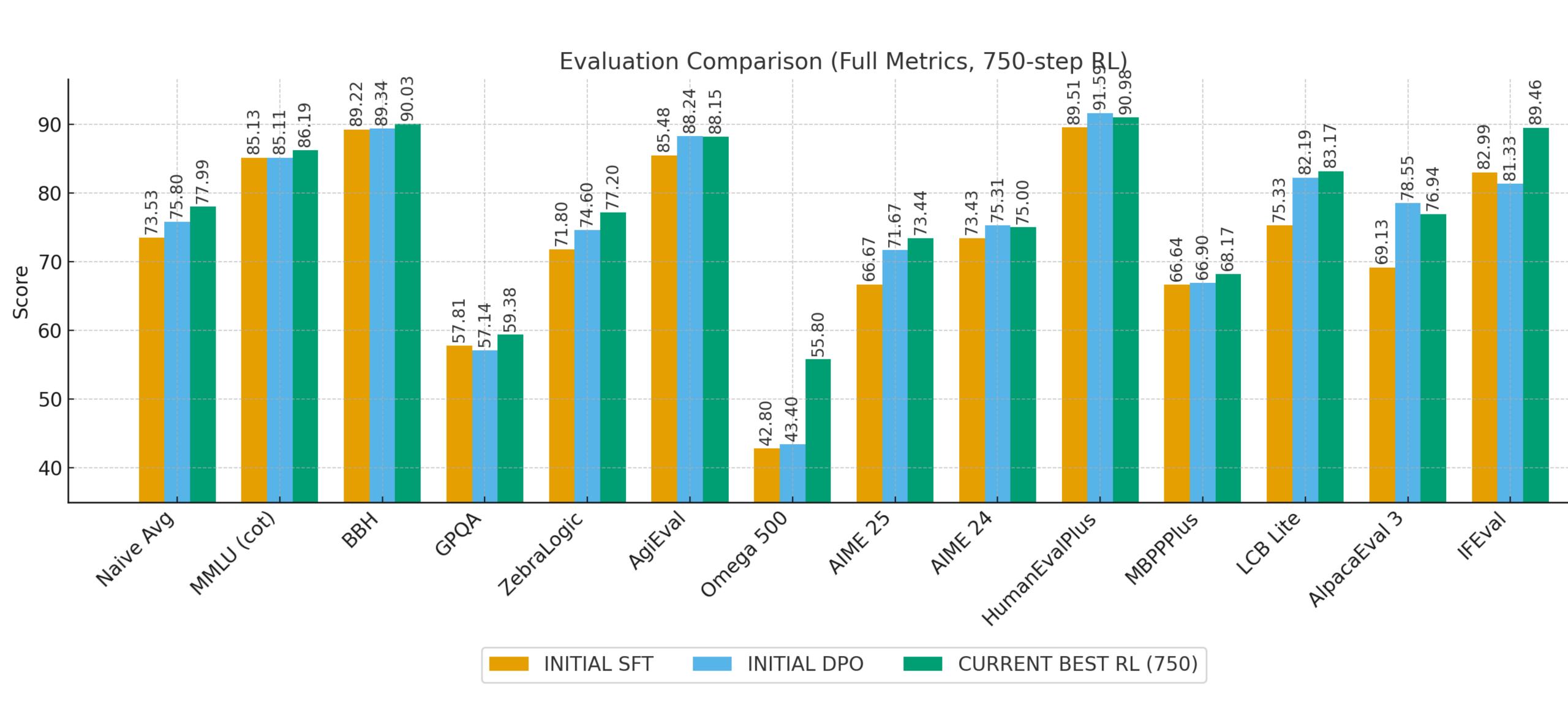
Also used pipelineRL to significantly speed up training

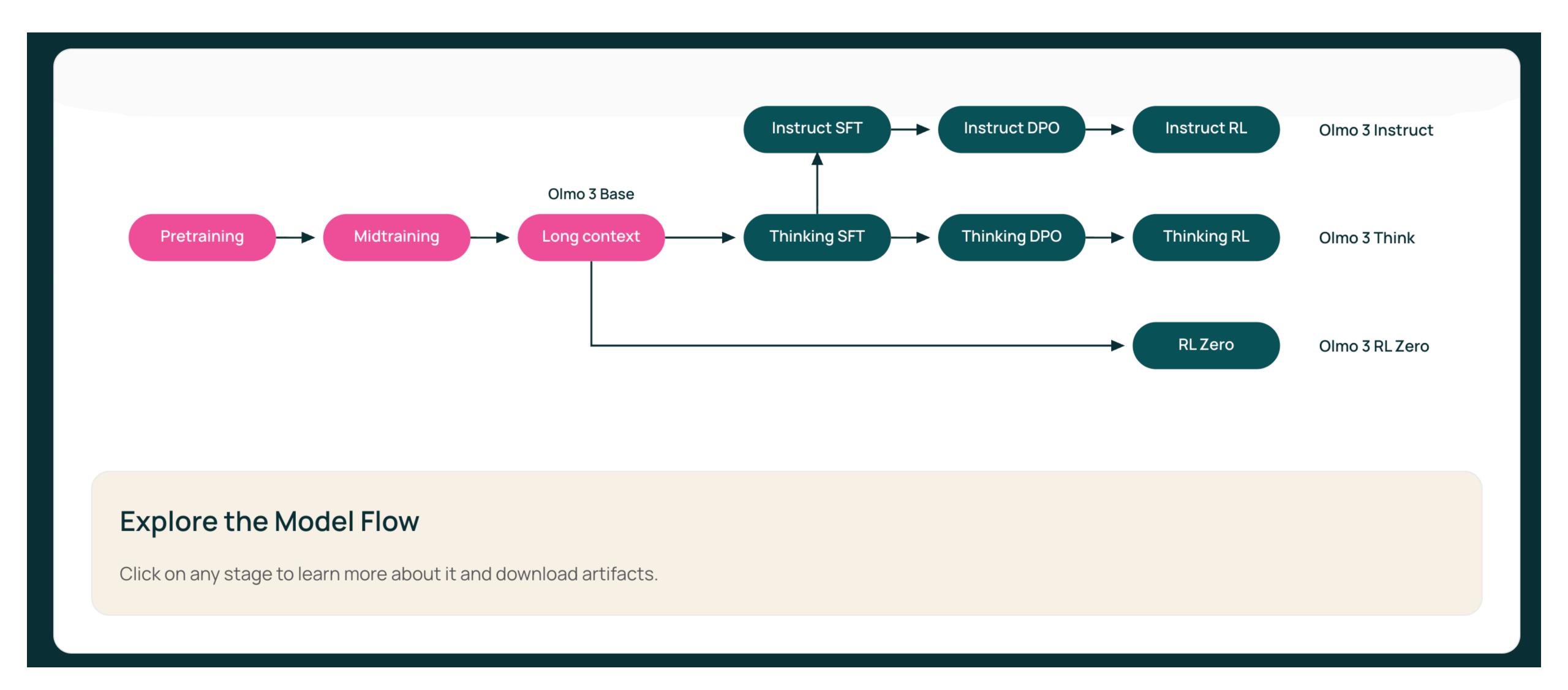


Mixing data prevents overfitting and collapse









https://allenai.org/blog/olmo3

Go forth and train SOTA!

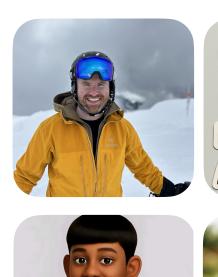
Further reading

- Always start with Sutton & Barto (chapter 13)
- Then OpenAl's spinning up RL course is great (https://spinningup.openai.com/en/latest/ index.html)
- Costa's 37 implementation details of PPO is amazing (https://iclr-blog-track.github.io/2022/03/25/ppoimplementation-details/)
- ScaleRL (Khatri et al., 2025) is a great overview of recent approaches.



Thanks to the many contributors for this work and slides:

Finbarr Timbers, Nathan Lambert, Teng Xiao, Michael Noukhovitch, Saumya Malik, Scott Geng, Faeze Brahman, Valentina Pyatkin, Saurabh Shah, Costa Huang, Hannaneh Hajishirzi, H2Lab as a whole.



























Thanks to the many contributors for this work and slides:

All my code for RL training can be found on github @ open-instruct

Happy to answer further questions etc:

hamishiv@cs.washington.edu

@hamishivi (and most other things)



Citations for the work discussed:

- 1. Levine S. Supervised Learning of Behaviors (Lecture 2). Published online August 2021.
- 2. Liu J, Li Y, Fu Y, Wang J, Liu Q, Shen Y. When Speed Kills Stability: Demystifying RL Collapse from the Inference-Training Mismatch. Published online September 2025. https://yingru.notion.site/When-Speed-Kills-Stability-Demystifying-RL-Collapse-from-the-Inference-Training-Mismatch-271211a558b7808d8b12d403fd15edda
- 3. Silver D, Schrittwieser J, Simonyan K, et al. Mastering the game of Go without human knowledge. Nature. 2017;550(7676):354-359. doi:10.1038/nature24270
- 4. Berseth G. Is Exploration or Optimization the Problem for Deep Reinforcement Learning? arXiv [csLG]. Published online 2025. http://arxiv.org/abs/2508.01329
- 5. Sutton RS, Barto AG. Reinforcement Learning: An Introduction. Second. The MIT Press; 2018. http://incompleteideas.net/book/the-book-2nd.html
- 6. Schulman J, Moritz P, Levine S, Jordan M, Abbeel P. High-Dimensional Continuous Control Using Generalized Advantage Estimation. In: *Proceedings of the International Conference on Learning Representations (ICLR)*.; 2016.
- 7. Shao Z, Wang P, Zhu Q, et al. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. arXiv [csCL]. Published online 2024. http://arxiv.org/abs/2402.03300
- 8. Ahmadian A, Cremer C, Gallé M, et al. Back to Basics: Revisiting REINFORCE Style Optimization for Learning from Human Feedback in LLMs. arXiv [csLG]. Published online 2024. http://arxiv.org/abs/2402.14740
- 9. Hu J, Liu JK, Xu H, Shen W. REINFORCE++: An Efficient RLHF Algorithm with Robustness to Both Prompt and Reward Models. arXiv [csCL]. Published online 2025. http://arxiv.org/abs/2501.03262
- 10. Kazemnejad A, Aghajohari M, Portelance E, et al. VinePPO: Unlocking RL Potential For LLM Reasoning Through Refined Credit Assignment. arXiv [csLG]. Published online 2024. http://arxiv.org/abs/2410.01679
- 11. Schulman J, Wolski F, Dhariwal P, Radford A, Klimov O. Proximal Policy Optimization Algorithms. arXiv [csLG]. Published online 2017. http://arxiv.org/abs/1707.06347
- 12. Yu Q, Zhang Z, Zhu R, et al. DAPO: An Open-Source LLM Reinforcement Learning System at Scale. arXiv [csLG]. Published online 2025. http://arxiv.org/abs/2503.14476
- 13. MiniMax, :, Chen A, et al. MiniMax-M1: Scaling Test-Time Compute Efficiently with Lightning Attention. arXiv [csCL]. Published online 2025. http://arxiv.org/abs/2506.13585
- 14. Zheng C, Liu S, Li M, et al. Group Sequence Policy Optimization. arXiv [csLG]. Published online 2025. http://arxiv.org/abs/2507.18071
- 15. Yao F, Liu L, Zhang D, Dong C, Shang J, Gao J. Your Efficient RL Framework Secretly Brings You Off-Policy RL Training. Feng Yao's Notion. Published online August 2025. https://fengyao.notion.site/off-policy-rl
- 16.Qi P, Liu Z, Zhou X, et al. Defeating the Training-Inference Mismatch via FP16. arXiv [csLG]. Published online 2025. http://arxiv.org/abs/2510.26788
- 17.Fu W, Gao J, Shen X, et al. AReaL: A Large-Scale Asynchronous Reinforcement Learning System for Language Reasoning. arXiv [csLG]. Published online 2025. http://arxiv.org/abs/2505.24298
- 18. Piché A, Kamalloo E, Pardinas R, Chen X, Bahdanau D. PipelineRL: Faster On-policy Reinforcement Learning for Long Sequence Generation. arXiv [csLG]. Published online 2025. http://arxiv.org/abs/2509.19128
- 19. Khatri D, Madaan L, Tiwari R, et al. The Art of Scaling Reinforcement Learning Compute for LLMs. arXiv [csLG]. Published online 2025. http://arxiv.org/abs/2510.13786
- 20.P. Aggarwal and S. Welleck. L1: Controlling how long a reasoning model thinks with reinforcement learning, 2025.