# On the Rollout-Training Mismatch In Modern RL Systems

Feng Yao

**UCSD** 

# Efficient RL systems are rising

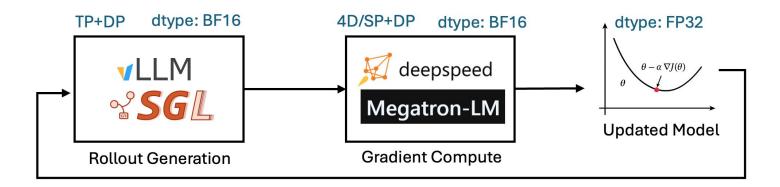
VeRL/OpenRLHF adopts hybrid engines

# Efficient RL systems are rising

- VeRL/OpenRLHF adopts hybrid engines
  - Rollout: Advanced LLM inference engines (vLLM, SGLang)
  - Training: Modern LLM training backends (FSDP, Megatron)

## Efficient RL systems are rising

- VeRL/OpenRLHF adopts hybrid engines
  - Rollout: Advanced LLM inference engines (vLLM, SGLang)
  - Training: Modern LLM training backends (FSDP, Megatron)



- Rollout-Training Mismatch
  - Expected

$$heta \leftarrow heta + \mu \cdot \mathbb{E}_{\underbrace{a \sim \pi( heta)}_{rollout}}[R(a) \cdot \underbrace{
abla_{ heta} \log \pi(a, heta)}_{training}]$$

- Rollout-Training Mismatch
  - Expected

$$heta \leftarrow heta + \mu \cdot \mathbb{E}_{\underbrace{a \sim \pi( heta)}_{rollout}}[R(a) \cdot \underbrace{
abla_{ heta} \log \pi(a, heta)}_{training}]$$

Implementation: Rollout engine (vLLM) + Training backends (FSDP)

$$heta \leftarrow heta + \mu \cdot \mathbb{E}_{a \sim \pi_{ ext{vilm}}( heta)}[R(a) \cdot 
abla_{ heta} \log \pi_{ ext{fsdp}}(a, heta)]$$

- Rollout-Training Mismatch
  - Expected

$$heta \leftarrow heta + \mu \cdot \mathbb{E}_{\underbrace{a \sim \pi( heta)}_{rollout}}[R(a) \cdot \underbrace{
abla_{ heta} \log \pi(a, heta)}_{training}]$$

Implementation: Rollout engine (vLLM) + Training backends (FSDP)

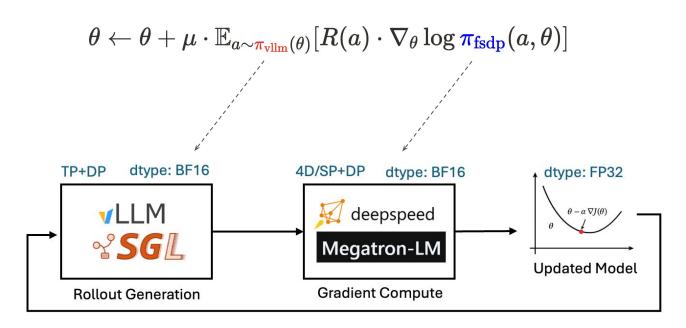
$$heta \leftarrow heta + \mu \cdot \mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta)}[R(a) \cdot 
abla_{ heta} \log \pi_{ ext{fsdp}}(a, heta)]$$

Mismatch!

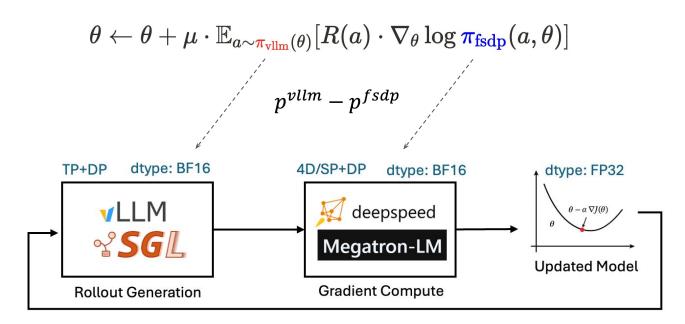
- Rollout-Training Mismatch
  - For the **same** rollout & model parameter

$$heta \leftarrow heta + \mu \cdot \mathbb{E}_{a \sim oldsymbol{\pi_{ ext{slm}}}( heta)}[R(a) \cdot 
abla_{ heta} \log oldsymbol{\pi_{ ext{fsdp}}}(a, heta)]$$

- Rollout-Training Mismatch
  - For the same rollout & model parameter



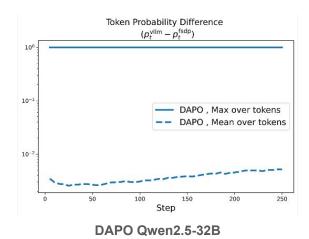
- Rollout-Training Mismatch
  - For the same rollout & model parameter

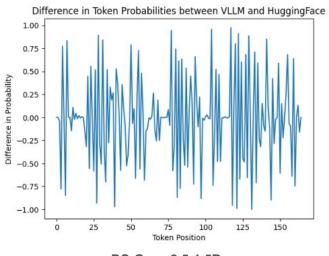


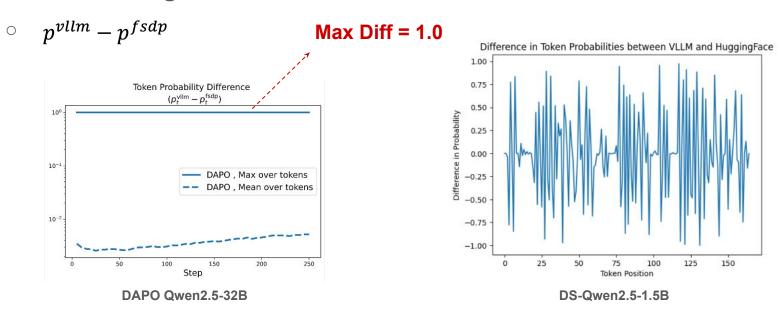
$$\circ \quad p^{vllm} - p^{fsdp}$$

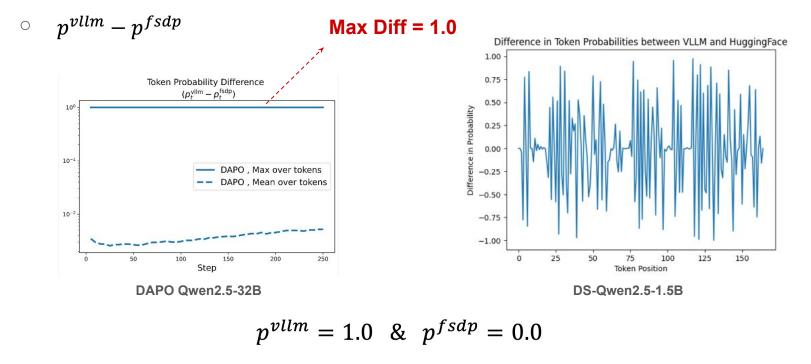
#### Rollout-Training Mismatch

 $\circ \quad p^{vllm} - p^{fsdp}$ 

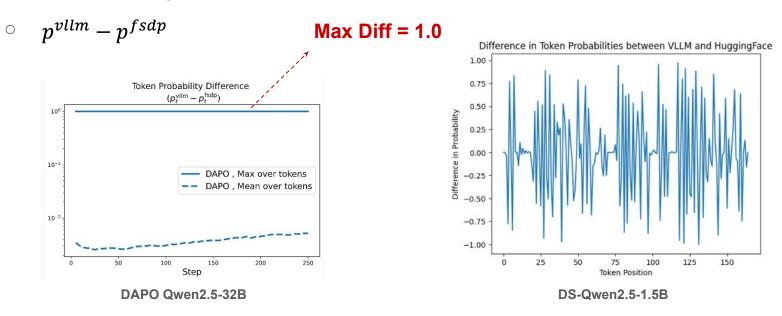








#### Implicitly makes RL "Off-Policy"!



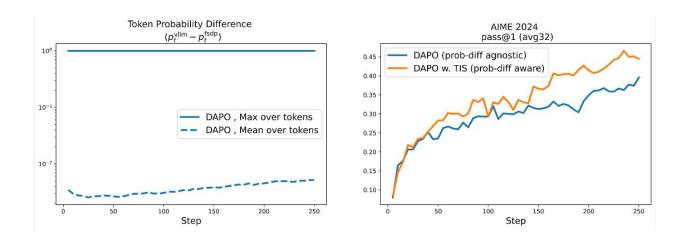
$$p^{vllm} = 1.0 \& p^{fsdp} = 0.0$$

## But it can be fixed effectively

Using the classic Truncated Importance Sampling (TIS) technique

# But it can be fixed effectively

- Using the classic *Truncated Importance Sampling (TIS)* technique
  - We show that fix it with TIS can improve training effectiveness



# Harvesting the Off-Policyness via Quantization

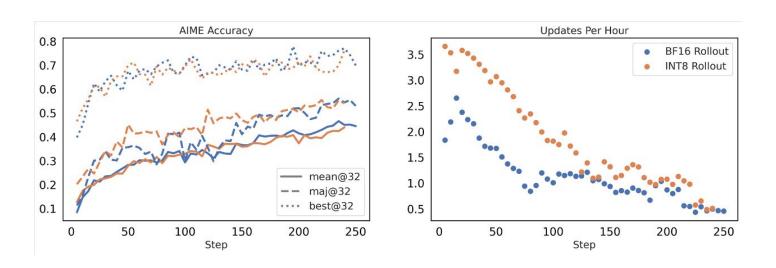
Since TIS is able to handle the mismatch.

# Harvesting the *Off-Policyness* via Quantization

- Since TIS is able to handle the mismatch.
  - Can we go even more "off-policy" and thus faster?

# Harvesting the Off-Policyness via Quantization

- Since TIS is able to handle the mismatch.
  - Can we go even more "off-policy" and thus faster?



#### **Outline**

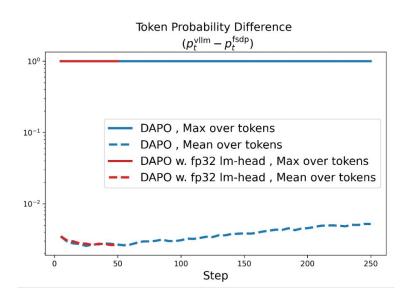
- Why Rollout-Training Mismatch Occurs
- How to Fix the Off-Policy Issue It Brings
- Harvesting Rollout-Training Mismatch via Quantization
- Analyzing the Effectiveness of Different Fixes

#### **Outline**

- Why Rollout-Training Mismatch Occurs
- How to Fix the Off-Policy Issue It Brings
- Harvesting Rollout-Training Mismatch via Quantization
- Analyzing the Effectiveness of Different Fixes

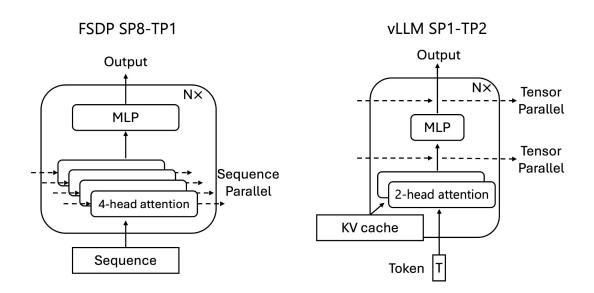
Two common believes

- Two common believes
  - Inaccessible true sampling probabilities
    - Add additional gap
  - Backend numerical differences
    - Hard to fix



- Hybrid Engine & Error Propagation
  - Different compute patterns via different backends & parallelism

- Hybrid Engine & Error Propagation
  - Different compute patterns via different backends & parallelism



#### **Outline**

- Why Rollout-Training Mismatch Occurs
- How to Fix the Off-Policy Issue It Brings
- Harvesting Rollout-Training Mismatch via Quantization
- Analyzing the Effectiveness of Different Fixes

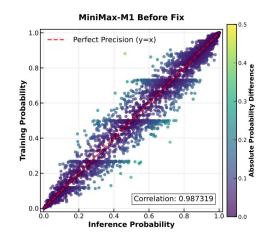
- Trial 1 Mitigate the system-level mismatch
  - vLLM seems to be the root cause

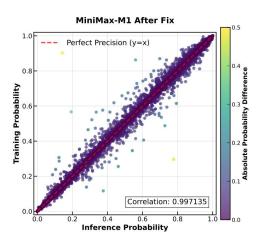
- Trial 1 Mitigate the system-level mismatch
  - vLLM seems to be the root cause → Patch vLLM to:

- Trial 1 Mitigate the system-level mismatch
  - vLLM seems to be the root cause → Patch vLLM to:
    - Return the actual sampling probabilities for vLLM V1 engine
    - Improve the numerical precision by using FP32 LM\_Head

- Trial 1 Mitigate the system-level mismatch
  - vLLM seems to be the root cause → Patch vLLM to:
    - Return the actual sampling probabilities for vLLM V1 engine
    - Improve the numerical precision by using FP32 LM\_Head

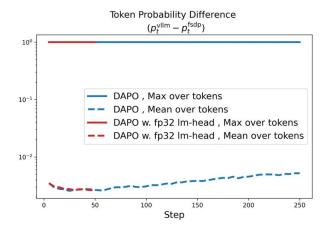
It helps, but ...

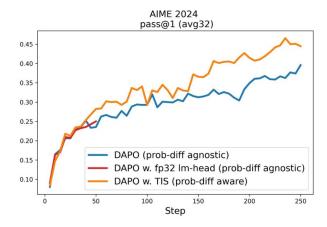




- Trial 1 Mitigate the system-level mismatch
  - vLLM seems to be the root cause → Patch vLLM to:
    - Return the actual sampling probabilities for vLLM V1 engine
    - Improve the numerical precision by using FP32 LM\_Head

It helps, but the gap still exists





- Trial 2 Apply algorithm-level fix
  - Be aware of the mismatch → **Importance sampling correction**:

- Trial 2 Apply algorithm-level fix
  - Be aware of the mismatch → Importance sampling correction:
    - Recall: Vanilla Importance Sampling

$$\mathbb{E}_{a \sim \pi_{ ext{fsdp}}( heta)}[R(a)] = \mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta)} \left[ \underbrace{rac{\pi_{ ext{fsdp}}(a, heta)}{\pi_{ ext{vllm}}(a, heta)}}_{ ext{importance ratio}} \cdot R(a) 
ight]$$

- Trial 2 Apply algorithm-level fix
  - Be aware of the mismatch → **Importance sampling correction**:
    - Expected gradient

$$\mathbb{E}_{a \sim \pi_{ ext{fsdp}}( heta)}[R(a) \cdot 
abla_{ heta} \log \pi_{ ext{fsdp}}(a, heta)]$$

- Trial 2 Apply algorithm-level fix
  - Be aware of the mismatch → Importance sampling correction:
    - Expected gradient

$$\mathbb{E}_{a \sim \pi_{\mathrm{fsdp}}( heta)}[R(a) \cdot 
abla_{ heta} \log \pi_{\mathrm{fsdp}}(a, heta)]$$

But currently we have

$$\mathbb{E}_{a \sim oldsymbol{\pi}_{ ext{vllm}}( heta)}[R(a) \cdot 
abla_{ heta} \log oldsymbol{\pi}_{ ext{fsdp}}(a, heta)]$$

- Trial 2 Apply algorithm-level fix
  - Be aware of the mismatch → Importance sampling correction:
    - Expected gradient

$$\mathbb{E}_{a \sim \pi_{\mathrm{fsdp}}( heta)}[R(a) \cdot 
abla_{ heta} \log \pi_{\mathrm{fsdp}}(a, heta)]$$

But currently we have

$$\mathbb{E}_{a \sim oldsymbol{\pi}_{ ext{vllm}}( heta)}[R(a) \cdot 
abla_{ heta} \log oldsymbol{\pi}_{ ext{fsdp}}(a, heta)]$$

So we should fix the gradient as:

$$\mathbb{E}_{a \sim oldsymbol{\pi_{ ext{vilm}}}( heta)} \Big[ rac{oldsymbol{\pi_{ ext{fsdp}}}(a, heta)}{oldsymbol{\pi_{ ext{vilm}}}(a, heta)} \cdot R(a) \cdot 
abla_{ heta} \log oldsymbol{\pi_{ ext{fsdp}}}(a, heta) \Big]$$

- Trial 2 Apply algorithm-level fix
  - Be aware of the mismatch → Importance sampling correction:
    - Expected gradient

$$\mathbb{E}_{a \sim \pi_{\mathrm{fsdp}}( heta)}[R(a) \cdot 
abla_{ heta} \log \pi_{\mathrm{fsdp}}(a, heta)]$$

But currently we have

$$\mathbb{E}_{a \sim oldsymbol{\pi}_{ ext{vllm}}( heta)}[R(a) \cdot 
abla_{ heta} \log oldsymbol{\pi}_{ ext{fsdp}}(a, heta)]$$

■ In practice, we use **Truncated Importance Sampling (TIS)**:

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta)} \Big[ \underbrace{\min \Big( rac{\pi_{ ext{fsdp}}(a, heta)}{\pi_{ ext{vllm}}(a, heta)}, C \Big)}_{ ext{truncated importance ratio}} \cdot R(a) \cdot 
abla_{ heta} \log \pi_{ ext{fsdp}}(a, heta) \Big]$$

Extend to General Case

- Extend to General Case
  - Expected Policy Gradient (PPO)

$$\mathbb{E}_{a \sim \pi_{ ext{fsdp}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, \operatorname{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

- Extend to General Case
  - Expected Policy Gradient (PPO)

$$\mathbb{E}_{a \sim \pi_{ ext{fsdp}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, \operatorname{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

VeRL/OpenRLHF's Implementation (recompute)

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \min \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})}, 1 - \epsilon, \; 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

#### Extend to General Case

Expected Policy Gradient (PPO)

$$\mathbb{E}_{a \sim \pi_{ ext{fsdp}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, \operatorname{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

VeRL/OpenRLHF's Implementation (recompute)

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \min \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})}, 1 - \epsilon, \; 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

Truncated Importance Sampling (TIS)

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

Variants of TIS

- Variants of TIS
  - PPO Importance Sampling (PPO-IS)

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{vllm}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{vllm}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

- Variants of TIS
  - PPO Importance Sampling (PPO-IS)



$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{vllm}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{vllm}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

- Variants of TIS
  - PPO Importance Sampling (PPO-IS)



$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{vllm}}(a, \, heta_{ ext{old}})} \, \hat{A}, \operatorname{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{vllm}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$
Can break out of the trust region

- Variants of TIS
  - PPO Importance Sampling (PPO-IS)



$$\mathbb{E}_{a \sim \pi_{\text{vllm}}(\theta_{\text{old}})} \Big[ \nabla_{\theta} \, \min \Big( \frac{\pi_{\text{fsdp}}(a, \, \theta)}{\pi_{\text{vllm}}(a, \, \theta_{\text{old}})} \, \hat{A}, \text{clip} \Big( \frac{\pi_{\text{fsdp}}(a, \, \theta)}{\pi_{\text{vllm}}(a, \, \theta_{\text{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$
Can break out of the trust region

Vanilla Importance Sampling (Vanilla-IS)

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

- Variants of TIS
  - PPO Importance Sampling (PPO-IS)



$$\mathbb{E}_{a \sim \pi_{\text{vllm}}(\theta_{\text{old}})} \left[ \nabla_{\theta} \ \min \left( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})} \ \hat{A}, \text{clip} \left( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \right) \hat{A} \right) \right]$$
Can break out of the trust region

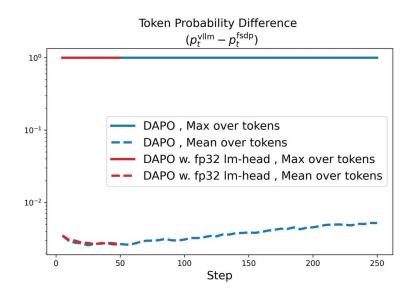
Vanilla Importance Sampling (Vanilla-IS)

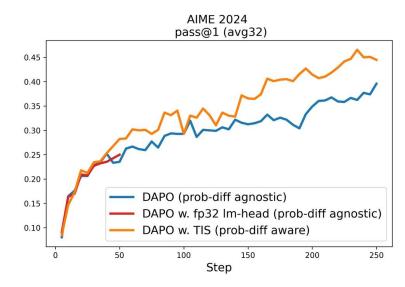
$$\mathbb{E}_{\frac{\pi_{\text{vllm}}(\theta_{\text{old}})}{\text{importance ratio}}} \frac{\pi_{\text{fsdp}}(a, \theta_{\text{old}})}{\pi_{\text{vllm}}(a, \theta_{\text{old}})} \cdot \nabla_{\theta} \ \min \left( \frac{\pi_{\text{fsdp}}(a, \theta)}{\pi_{\text{fsdp}}(a, \theta_{\text{old}})} \ \hat{A}, \text{clip} \left( \frac{\pi_{\text{fsdp}}(a, \theta)}{\pi_{\text{fsdp}}(a, \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \right) \hat{A} \right) \right]$$

Can be too large and makes training crash

#### How well can TIS fix it?

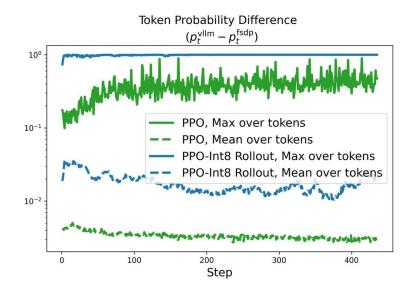
#### DAPO 32B Setting

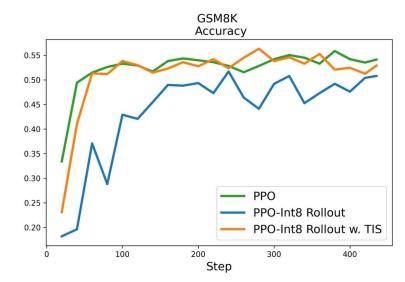




#### How well can TIS fix it?

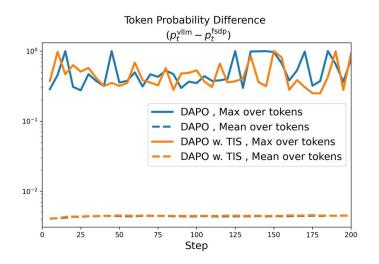
- GSM8K 0.5B Setting
  - Normal RL: Max Diff is smaller (~0.4) than 1.0 (in DAPO-32B setting)
  - INT8 Rollout: Max Diff is larger (~1.0) than normal RL setting

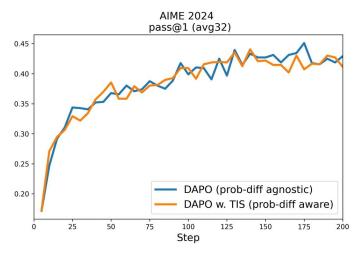




# Does TIS always help?

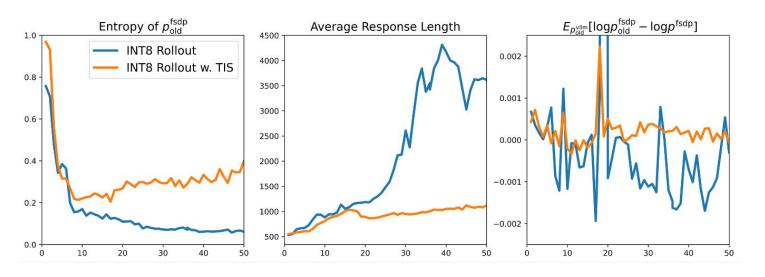
- DAPO 1.5B Setting
  - In settings where prob diff is relatively small
    - TIS does not always help, but doesn't hurt





# **Does the Mismatch really matter?**

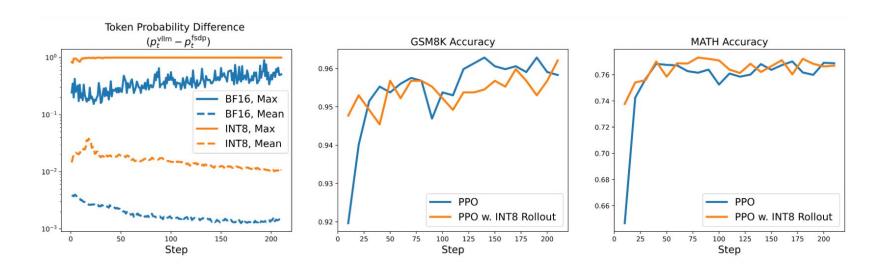
Unexpected training instability on challenging tasks



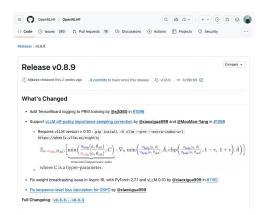
DAPO Qwen2.5-32B

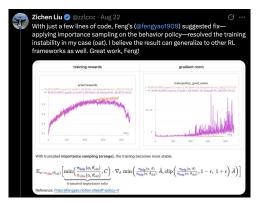
### **Does the Mismatch really matter?**

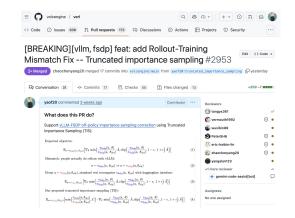
Possible negligible on simple tasks

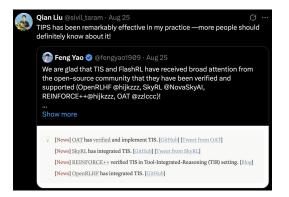


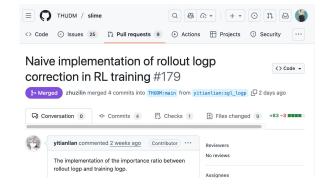
#### **Community Verification**











#### REINFORCE++-baseline is all you need in RLVR



#### **Tool-Integrated Reasoning and Agent Experiments**

We have thoroughly validated the effectiveness of global standard deviation and global advantage normalization in complex multi-turn tool call scenarios. Our experiments utilize the framework established by <a href="arXiv:2505.07773">arXiv:2505.07773</a>, which features a zero-shot agent environment designed for large language models to tackle mathematical problems using Qwen 2.5 Base 7B.

	300-step/avg@32/iternum-2(inference)	aime24	aime25	hmmt_feb_ 2025	hmmt_feb_ 2024	dimmc	avg
reinforce++	with-baseline	30.83333	27.1875	17.91667	18.95833	25.625	24.10416667
reinforce++	with-baseline,vllm-correction	31.5625	23.9583	20.625	20.1042	29.375	25.125
grpo		31.66667	21.875	16.97917	17.70833	24.6875	22.58333333
рро		30.20833	21.66667	15	18.4375	23.95833	21.85416667

REINFORCE++-baseline achieves the best performance in the multi-turn tool call tasks.

The REINFORCE++ baseline can be combined with dynamic sampling, cliphigher, and <u>truncated importance sampling</u> (vLLM correction in the figure) to continuously improve performance.

#### **Outline**

- Why Rollout-Training Mismatch Occurs
- How to Fix the Off-Policy Issue It Brings
- Harvesting Rollout-Training Mismatch via Quantization
- Analyzing the Effectiveness of Different Fixes

#### Harvesting Off-Policy in Quantization

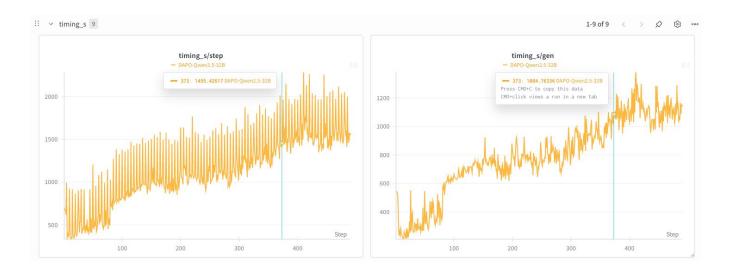
As TIS handles the gap, can we go even **further off-policy** for **speedup**?

### Harvesting Off-Policy in Quantization

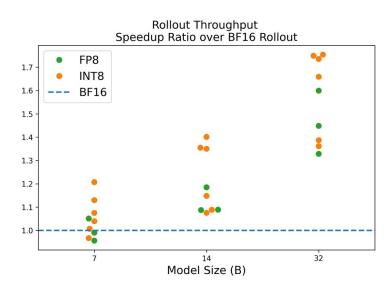
As TIS handles the gap, can we go even further off-policy for speedup?

**Rollout generation** is a bottleneck in RL training efficiency:

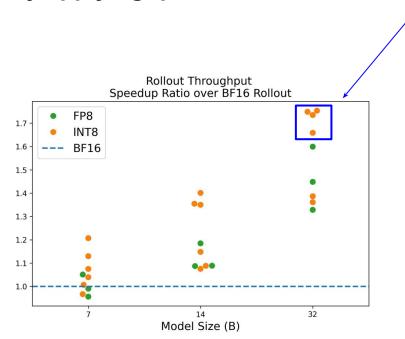
In DAPO-32B setting, rollout takes up ~70% of the training time



#### Naively applying quantization can accelerate rollout speed

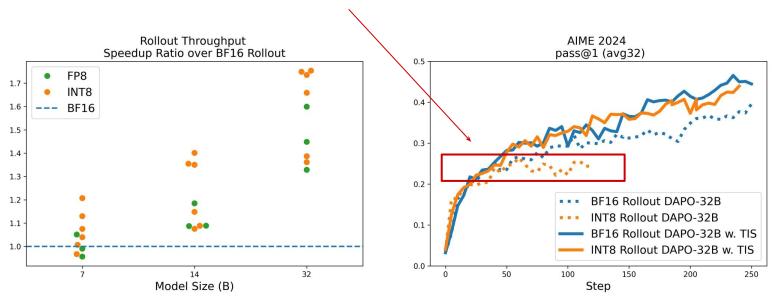


Naively applying quantization can accelerate rollout speed



#### Naively applying quantization can accelerate rollout speed

But the performance is also degraded!



Naively applying quantization can accelerate rollout speed

But the performance is also degraded!

#### This can be expected, as quantization introduces more mismatch

$$\underbrace{\mathbb{E}_{a \sim \pi_{\text{bf16}}(\theta_{\text{old}})}}_{\text{int8 Rollout: } \pi_{\text{bf16}} \rightarrow \pi_{\text{int8}}} \left[ \nabla_{\theta} \min \left( \frac{\pi_{\text{bf16}}(a, \theta)}{\pi_{\text{bf16}}(a, \theta_{\text{old}})} \, \hat{A}, \, \operatorname{clip} \left( \frac{\pi_{\text{bf16}}(a, \theta)}{\pi_{\text{bf16}}(a, \theta_{\text{old}})}, \, 1 - \epsilon, \, 1 + \epsilon \right) \hat{A} \right) \right]$$

63

This can be expected, as quantization introduces more mismatch

$$\underbrace{\mathbb{E}_{a \sim \pi_{\text{bf16}}(\theta_{\text{old}})}}_{\text{int8 Rollout: } \pi_{\text{bf16}} \rightarrow \pi_{\text{int8}}} \left[ \nabla_{\theta} \min \left( \frac{\pi_{\text{bf16}}(a, \theta)}{\pi_{\text{bf16}}(a, \theta_{\text{old}})} \, \hat{A}, \text{ clip} \left( \frac{\pi_{\text{bf16}}(a, \theta)}{\pi_{\text{bf16}}(a, \theta_{\text{old}})}, \, 1 - \epsilon, \, 1 + \epsilon \right) \hat{A} \right) \right]$$

#### FlashRL fixes it with TIS

$$\mathbb{E}_{a \sim \pi_{\text{int8}}(\theta_{\text{old}})} \left[ \underbrace{\min \left( \frac{\pi_{\text{bf16}}(a, \theta_{\text{old}})}{\pi_{\text{int8}}(a, \theta_{\text{old}})}, C \right)}_{\text{truncated importance ratio}} \cdot \nabla_{\theta} \ \min \left( \frac{\pi_{\text{bf16}}(a, \theta)}{\pi_{\text{bf16}}(a, \theta_{\text{old}})} \, \hat{A}, \text{clip} \left( \frac{\pi_{\text{bf16}}(a, \theta)}{\pi_{\text{bf16}}(a, \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \right) \hat{A} \right) \right]$$

#### FlashRL fixes it with TIS

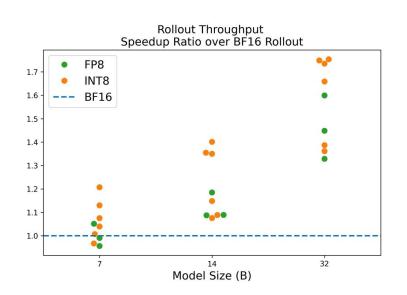
$$\mathbb{E}_{a \sim \pi_{ ext{int8}}( heta_{ ext{old}})} \left[ \underbrace{\min\left(rac{\pi_{ ext{bf16}}(a, heta_{ ext{old}})}{\pi_{ ext{int8}}(a, heta_{ ext{old}})}, C
ight)}_{ ext{truncated importance ratio}} \cdot 
abla_{ heta} \, \min\left(rac{\pi_{ ext{bf16}}(a, heta)}{\pi_{ ext{bf16}}(a, heta_{ ext{old}})} \, \hat{A}, ext{clip}\left(rac{\pi_{ ext{bf16}}(a, heta)}{\pi_{ ext{bf16}}(a, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon
ight) \hat{A} 
ight) 
brace$$

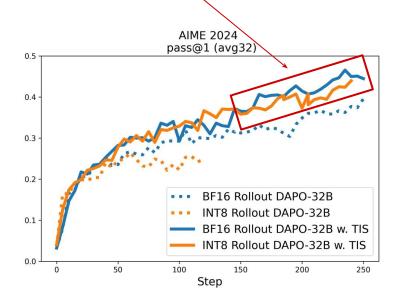
#### FlashRL is implemented as a PyPI package to patch vLLM

```
pip install flash-llm-rl  # install with pip
export FLASHRL_CONFIG='fp8' # turn on env variable
bash your-rl-training-script # no code change needed!
```

#### **DAPO 32B Setting**

Matches the performance of BF16 rollout with TIS

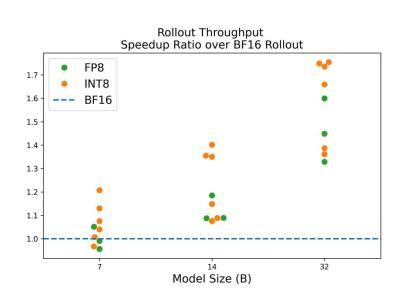


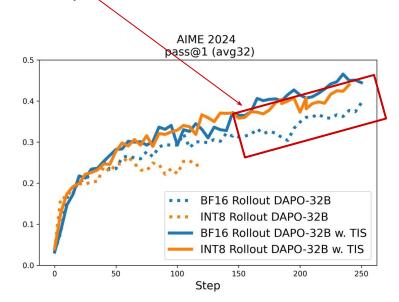


#### **DAPO 32B Setting**

Matches the performance of BF16 rollout with TIS

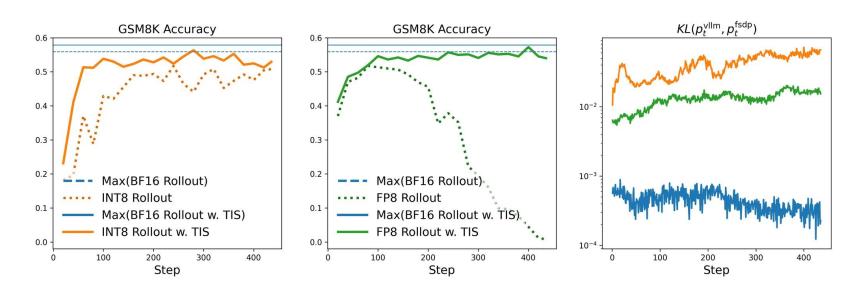
Outperforms naive BF16 rollout (without TIS)





#### **GSM8K 0.5B Setting**

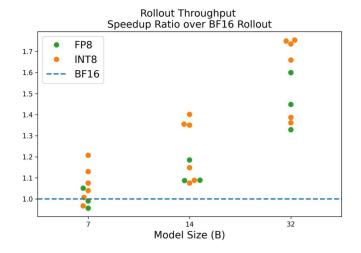
TIS works both in INT8 and FP8 setting



**Rollout Speedup** 

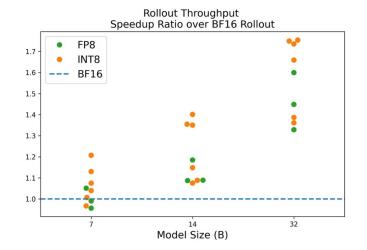
**Rollout Speedup** 

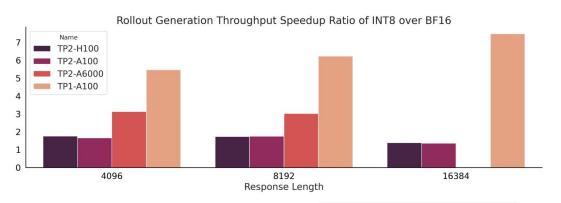
Regular RL Setting



#### **Rollout Speedup**

Regular RL Setting
Standard Inference Setting





**Figure 3.** Throughput speedup ratio of INT8-quantized Deepseek-R1-Distill-Qwen-32B relative to BF16 in 4 inference-only configurations, measured across varying response lengths

**End-to-End Speedup & Effectiveness** 

## More detailed analysis

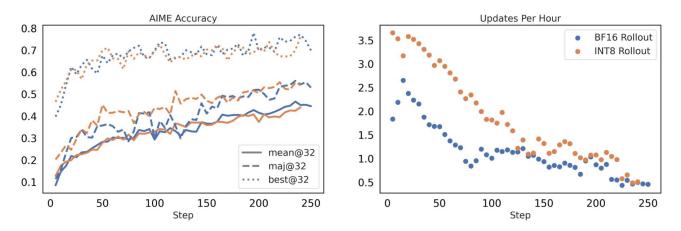
**End-to-End Speedup & Effectiveness** 

INT8 as a pressure test

## More detailed analysis

#### **End-to-End Speedup & Effectiveness**

INT8 as a pressure test



**Figure 4. Left:** Downstream performance of RL training with BF16 vs. INT8 rollout precision. **Right:** Updates per hour achieved with BF16 and INT8 rollout. All experiments use the DAPO recipe on Qwen2.5-32B, trained for 250 steps on 4 nodes with 8×H100 GPUs. [wandb]

FP8 quantization can be naturally conducted in an online manner

FP8 quantization can be naturally conducted in an online manner

**INT8 quantization** requires complicated calibration process

FP8 quantization can be naturally conducted in an online manner

**INT8 quantization** requires complicated calibration process

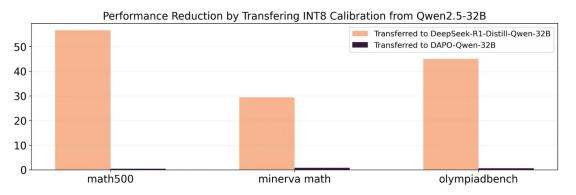
Our solution: Online INT8 Quantization via Calibration Transfer

calculate the calibration result once at the beginning of training and reuse it at every online step

#### Online INT8 Quantization via Calibration Transfer

calculate the calibration result once at the beginning of training and reuse it at every online step

Observation: RL changes model weights less aggressively comparing to SFT



**Figure 6.** We experiment on models finetuned via SFT / RL from Qwen2.5-32B base model. We find that compared to SFT, reusing the base model calibration result by applying it to RL funetuned model rarely change the performance. This indicates that INT8 online quantization is practically possible in RL by reusing previous calibration result.

#### **Outline**

- Why Rollout-Training Mismatch Occurs
- How to Fix the Off-Policy Issue It Brings
- Harvesting Rollout-Training Mismatch via Quantization
- Analyzing the Effectiveness of Different Fixes

PPO

$$\mathbb{E}_{a \sim \pi_{ ext{fsdp}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

Recompute

$$\mathbb{E}_{a \sim \pi_{\text{vllm}}(\theta_{\text{old}})} \Big[ \nabla_{\theta} \ \min \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{fsdp}}(a, \ \theta_{\text{old}})} \, \hat{A}, \text{clip} \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{fsdp}}(a, \ \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

PPO

$$\mathbb{E}_{a \sim \pi_{ ext{fsdp}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

Recompute

$$\mathbb{E}_{a \sim m{\pi}_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \min \Big( rac{m{\pi}_{ ext{fsdp}}(a, \; heta)}{m{\pi}_{ ext{fsdp}}(a, \; heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{m{\pi}_{ ext{fsdp}}(a, \; heta)}{m{\pi}_{ ext{fsdp}}(a, \; heta_{ ext{old}})}, 1 - \epsilon, \; 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

PPO-IS

$$\mathbb{E}_{a \sim \pi_{\text{vllm}}(\theta_{\text{old}})} \Big[ \nabla_{\theta} \ \min \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})} \ \hat{A}, \text{clip} \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \Big) \ \hat{A} \Big) \Big]$$

PPO

$$\mathbb{E}_{a \sim \pi_{ ext{fsdp}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

Recompute

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

PPO-IS

$$\mathbb{E}_{a \sim \pi_{\text{vllm}}(\theta_{\text{old}})} \Big[ \nabla_{\theta} \ \min \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})} \ \hat{A}, \text{clip} \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \Big) \ \hat{A} \Big) \Big]$$

Vanilla-IS

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

PPO

$$\mathbb{E}_{a \sim \pi_{ ext{fsdp}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \min \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})}, 1 - \epsilon, \; 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

Recompute

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \min \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})}, 1 - \epsilon, \; 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

PPO-IS

$$\mathbb{E}_{a \sim \pi_{\text{vllm}}(\theta_{\text{old}})} \left[ \nabla_{\theta} \ \min \left( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})} \, \hat{A}, \text{clip} \left( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \right) \hat{A} \right) \right]$$

Vanilla-IS

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

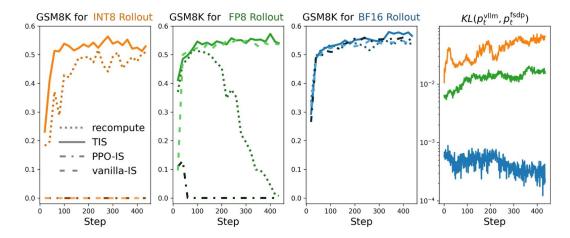
TIS

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

#### **Comparison with TIS-Variants**

#### GSM8K, PPO, Qwen2.5-0.5B-Instruct

Only TIS works consistently



**Figure 5.** We ablate different rollout-training mismatch mitigation strategies on Qwen2.5-0.5B with GSM8k. Note PPO-IS and Vanilla-IS achieves near 0 accuracy for INT8 rollouts thus being highly overlapped. We also plot the KL divergence between vLLM sampled distribution and the FSDP distribution on the right. [int8 wandb][fp8 wandb]

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

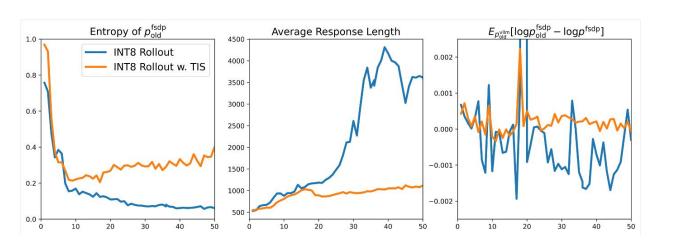
#### Recompute

- The mismatch can lead to entropy collapse
  - Gradient computation vs. rollout generation

$$\mathbb{E}_{a \sim \pi_{\text{vllm}}(\theta_{\text{old}})} \Big[ \nabla_{\theta} \ \min \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{fsdp}}(a, \ \theta_{\text{old}})} \, \hat{A}, \text{clip} \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{fsdp}}(a, \ \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

#### Recompute

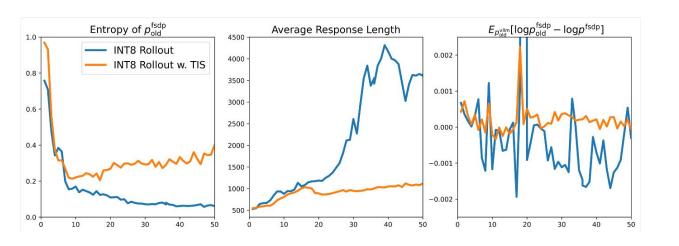
- The mismatch can lead to entropy collapse
  - Gradient computation vs. rollout generation



$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

#### Recompute

- The mismatch can lead to entropy collapse
  - Gradient computation vs. rollout generation

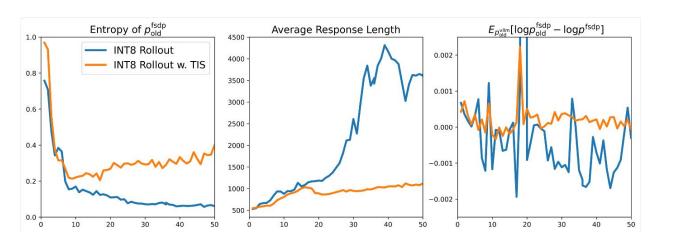


For a with A < 0

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

#### Recompute

- The mismatch can lead to entropy collapse
  - Gradient computation vs. rollout generation



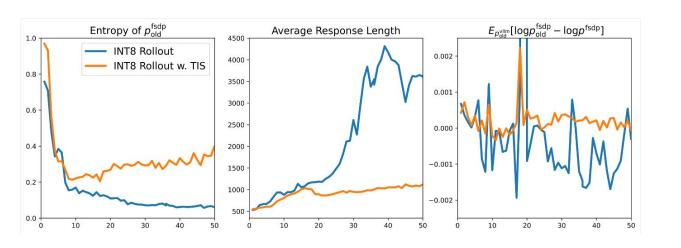
For a with A < 0

 $\pi_{ ext{fsdp}}(a)$  becomes smaller

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \min \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})}, 1 - \epsilon, \; 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

#### Recompute

- The mismatch can lead to entropy collapse
  - Gradient computation vs. rollout generation



For a with A < 0

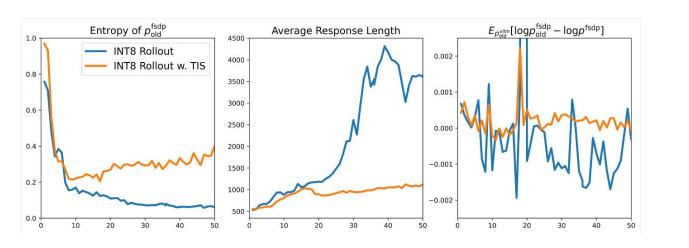
 $\pi_{ ext{fsdp}}(a)$  becomes smaller

With a large gap w. INT8,  $\pi_{\text{vllm}}(a)$  stays the same

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \min \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \; heta)}{\pi_{ ext{fsdp}}(a, \; heta_{ ext{old}})}, 1 - \epsilon, \; 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

#### Recompute

- The mismatch can lead to entropy collapse
  - Gradient computation vs. rollout generation



For a with A < 0

 $\pi_{ ext{fsdp}}(a)$  becomes smaller

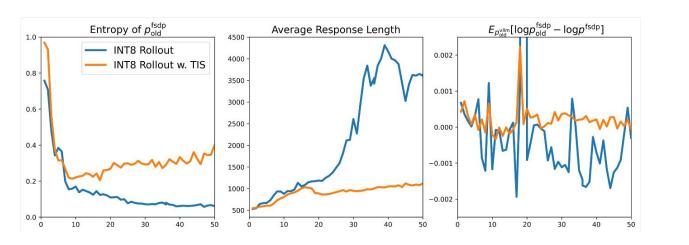
With a large gap w. INT8,  $\pi_{\text{vilm}}(a)$  stays the same

 $\pi_{\mathrm{fsdp}}(a)$  is over-penalized

$$\mathbb{E}_{a \sim \pi_{ ext{vllm}}( heta_{ ext{old}})} \Big[ 
abla_{ heta} \, \min \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} \Big( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon \Big) \, \hat{A} \Big) \Big]$$

#### Recompute

- The mismatch can lead to entropy collapse
  - Gradient computation vs. rollout generation



For a with A < 0

 $\pi_{
m fsdp}(a)$  becomes smaller

With a large gap w. INT8,  $\pi_{\text{vilm}}(a)$  stays the same

 $\pi_{
m fsdp}(a)$  is over-penalized

Small entropy

### Why PPO-IS fails

$$\mathbb{E}_{a \sim \pi_{\text{vllm}}(\theta_{\text{old}})} \Big[ \nabla_{\theta} \ \min \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})} \ \hat{A}, \text{clip} \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \Big) \ \hat{A} \Big) \Big]$$

#### PPO-IS

PPO-IS is still "biased" from the PPO gradient

$$\mathbb{E}_{a \sim \pi_{ ext{fsdp}}( heta_{ ext{old}})} igg[ 
abla_{ heta} \, \min igg( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})} \, \hat{A}, ext{clip} igg( rac{\pi_{ ext{fsdp}}(a, \, heta)}{\pi_{ ext{fsdp}}(a, \, heta_{ ext{old}})}, 1 - \epsilon, \, 1 + \epsilon igg) \, \hat{A} igg) igg]$$

### Why PPO-IS fails

$$\mathbb{E}_{a \sim \pi_{\text{vllm}}(\theta_{\text{old}})} \left[ \nabla_{\theta} \ \min \left( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})} \ \hat{A}, \text{clip} \left( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \right) \hat{A} \right) \right]$$

#### PPO-IS

PPO-IS is still "biased" from the PPO gradient

$$\mathbb{E}_{a \sim \pi_{\mathrm{fsdp}}(\theta_{\mathrm{old}})} \Big[ \nabla_{\theta} \ \min \Big( \frac{\pi_{\mathrm{fsdp}}(a, \ \theta)}{\pi_{\mathrm{fsdp}}(a, \ \theta_{\mathrm{old}})} \ \hat{A}, \mathrm{clip} \Big( \frac{\pi_{\mathrm{fsdp}}(a, \ \theta)}{\pi_{\mathrm{fsdp}}(a, \ \theta_{\mathrm{old}})}, 1 - \epsilon, \ 1 + \epsilon \Big) \ \hat{A} \Big) \Big]$$

- The clip in PPO is designed for "trust region"
  - lacktriangle At time step 0,  $heta= heta_{
    m old}$  , we don't want to clip but PPO-IS may clip
  - PPO-clip works differently than TIS

### Why PPO-IS fails

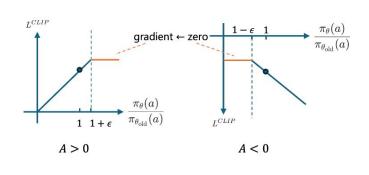
$$\mathbb{E}_{a \sim \pi_{\text{vllm}}(\theta_{\text{old}})} \Big[ \nabla_{\theta} \ \min \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})} \ \hat{A}, \text{clip} \Big( \frac{\pi_{\text{fsdp}}(a, \ \theta)}{\pi_{\text{vllm}}(a, \ \theta_{\text{old}})}, 1 - \epsilon, \ 1 + \epsilon \Big) \ \hat{A} \Big) \Big]$$

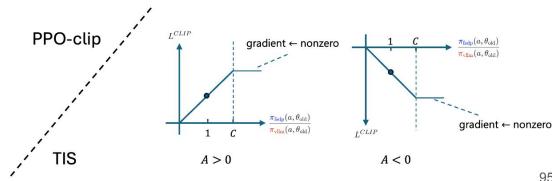
#### PPO-IS

PPO-IS is still "biased" from the PPO gradient

$$\mathbb{E}_{a \sim \pi_{\mathrm{fsdp}}(\theta_{\mathrm{old}})} \bigg[ \nabla_{\theta} \ \min \bigg( \frac{\pi_{\mathrm{fsdp}}(a, \ \theta)}{\pi_{\mathrm{fsdp}}(a, \ \theta_{\mathrm{old}})} \, \hat{A}, \mathrm{clip} \bigg( \frac{\pi_{\mathrm{fsdp}}(a, \ \theta)}{\pi_{\mathrm{fsdp}}(a, \ \theta_{\mathrm{old}})}, 1 - \epsilon, \ 1 + \epsilon \bigg) \, \hat{A} \bigg) \bigg]$$

- The clip in PPO is designed for "trust region" 0
  - At time step 0,  $\theta = \theta_{\rm old}$ , we don't want to clip but PPO-IS may clip
  - PPO-clip works differently than TIS





#### Vanilla-IS

- Uncapped importance ratio amplifies the gradient noise
  - Leading to unstable training

#### Vanilla-IS

- Uncapped importance ratio amplifies the gradient noise
  - Leading to unstable training

$$\frac{1}{N^2} \cdot \sum_{a_1}^{a_N} \left( \frac{\pi_{\mathrm{fsdp}}(a, \theta_{\mathrm{old}})}{\pi_{\mathrm{vllm}}(a, \theta_{\mathrm{old}})} \right)^2 \cdot Var[\nabla_{\theta} \cdots] \qquad = \qquad Var[ \qquad ]$$
importance ratio

$$\frac{1}{N}E\big[Var[\nabla_{\theta}(x,y)]\big] \approx \frac{1}{N^2} \sum_{x,y} Var[\nabla_{\theta}(x,y)] \qquad = \qquad Var\Big[\frac{1}{N} \sum_{x,y} \nabla_{\theta}(x,y)\Big]$$

Gradient noise

# What's beyond?

## What's beyond?

- The gap can be amplified in MoE RL
  - Dynamic Routing
  - Specially Optimized Kernels

## What's beyond?

- The gap can be amplified in MoE RL
  - Dynamic Routing
  - Specially Optimized Kernels

- TIS is orthogonal and compatible with existing GxPOs
  - GxPOs adjust the computation of advantage / importance ratio
  - TIS addresses the system-level mismatch problem

## **Takeaways**

 Mixing inference backend with training backends brings off-policy RL training, even if they share the same weights

Truncated Importance Sampling (TIS) is effective mitigating the gap

 With TIS integrated, rollout generation can be accelerated via quantization without sacrificing the performance

# **Thanks for Listening!**

Feng Yao

https://yaof20.github.io/