# AReaL: A Large-Scale Asynchronous Reinforcement Learning System for Language Reasoning

**Wei Fu (me)**, Jiaxuan Gao , Xujie Shen , Chen Zhu , Zhiyu Mei , Chuyi He , Shusheng Xu, Guo Wei , Jun Mei , Jiashu Wang, Tongkai Yang, Binhang Yuan, Yi Wu

**Contact & Collaboration**: fuwth17@gmail.com & jxwuyi@gmail.com

香港科技大學 THE HONG KONG UNIVERSITY OF SCIENCE AND TECHNOLOGY

清華大學 交叉信息研究院 Institute for Interdisciplinary Information Sciences, Tsinghua University
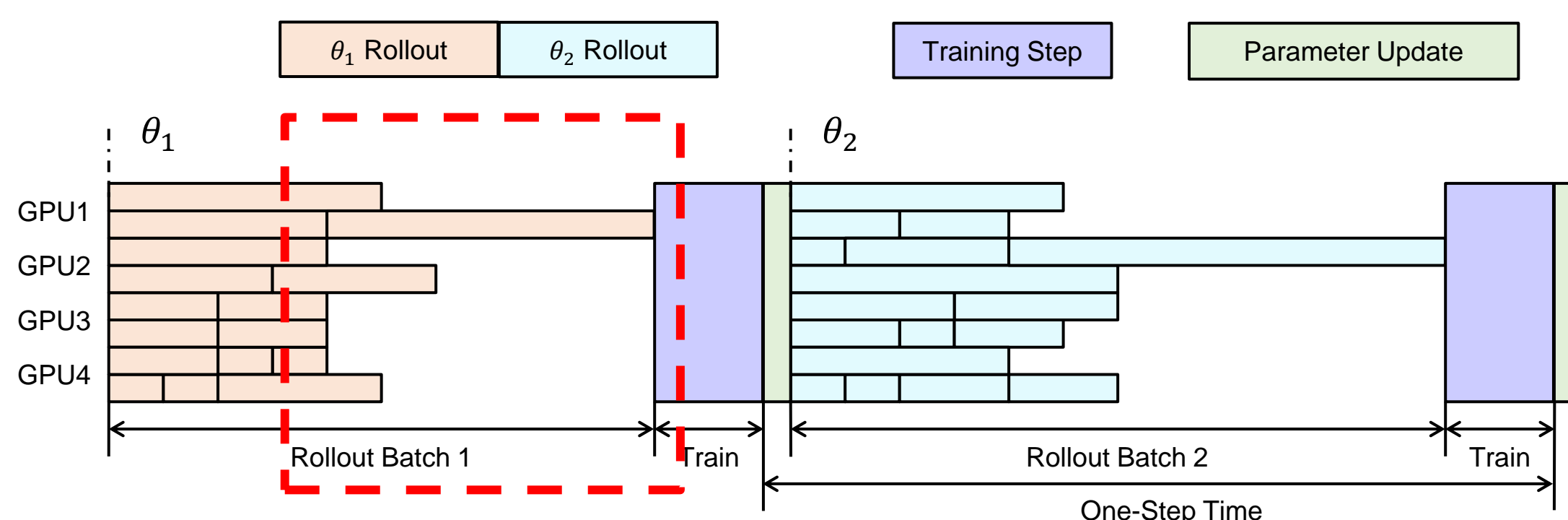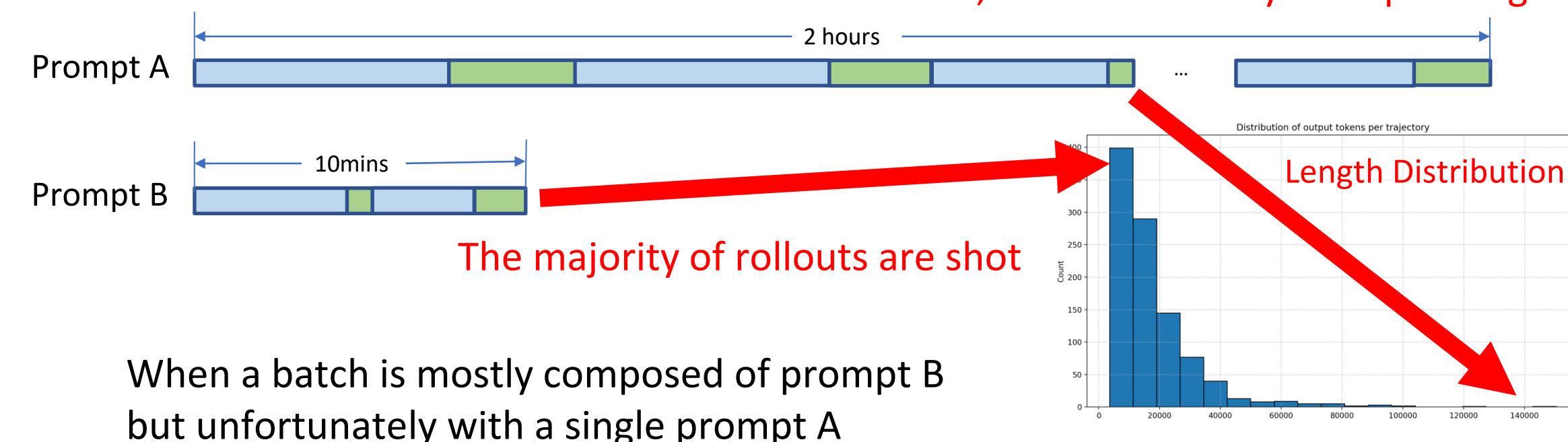
蚂蚁集团 ANT GROUP

NEURAL INFORMATION PROCESSING SYSTEMS

## CHALLENGES OF SYNC. RL SYSTEMS

**Challenge 1**: Significant *GPU idle times during* inference/generation when lengths vary

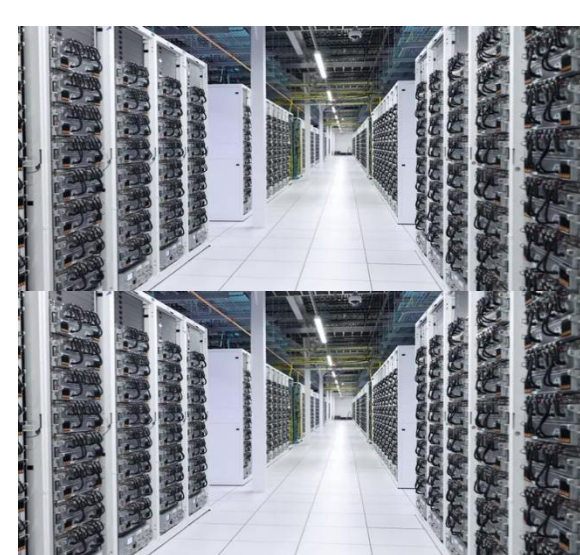### Multi-Turn (>128 turns) Agentic RL with Reasoning
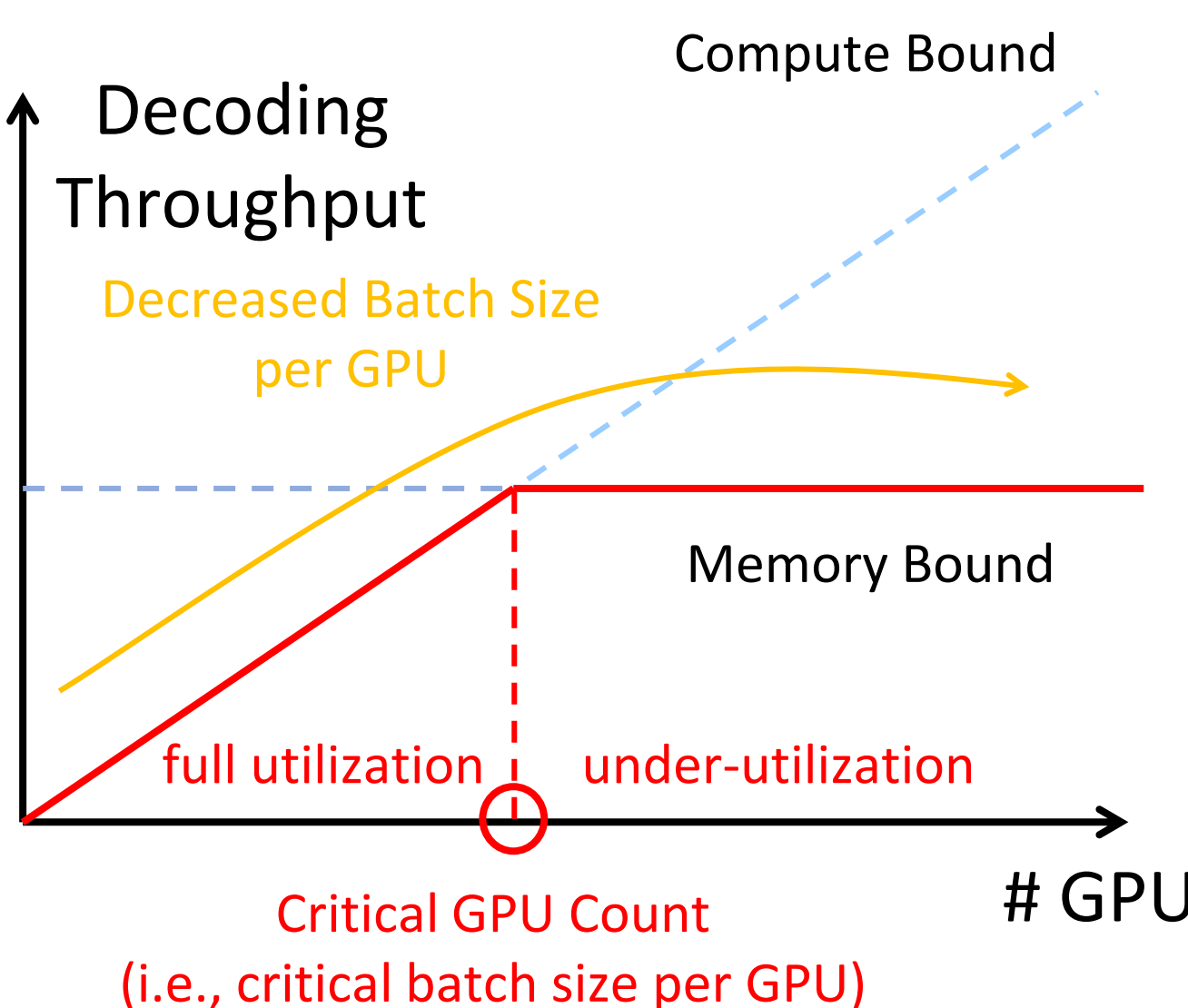


..., while some may be super long

Prompt A — 2 hours

Prompt B — 10mins

Length Distribution

The majority of rollouts are shot

When a batch is mostly composed of prompt B but unfortunately with a single prompt A

2 hours

**GPU: Just relax**

$\theta_1$ Rollout | $\theta_2$ Rollout | Training Step | Parameter Update



GPU1 GPU2 GPU3 GPU4

Rollout Batch 1 | Train | Rollout Batch 2 | Train

One-Step Time

**Challenge 2**: *Hard to scale up*

= 100 hours RL training

≠ 50 hours RL training



Compute Bound

Decoding Throughput

Decreased Batch Size per GPU

Memory Bound

full utilization | under-utilization

Critical GPU Count (i.e., critical batch size per GPU)

# GPU

## FULLY ASYNCHRONOUS RL IN AReaL

**Both training and inference achieve full GPU utilization**

Inference Servers

$\pi_t$  $\pi_{t+1}$  $\pi_{t+2}$

PPO Actor

$\pi_t$ → $\pi_{t+1}$

**Updating new parameters immediately**

AReaL -lite
Lightning-Fast RL

**GitHub**

$\pi_{t+1}$ → $\pi_{t+2}$

**Starting training as long as a new batch arrives**

$\pi_{t+2}$

**Keys to the challenges:**

1. **Continuous Rollout Batching with Interruption**

2. **Disaggregated and Overlapped Training & Inference**

Rollout Interruption & Parameter Update
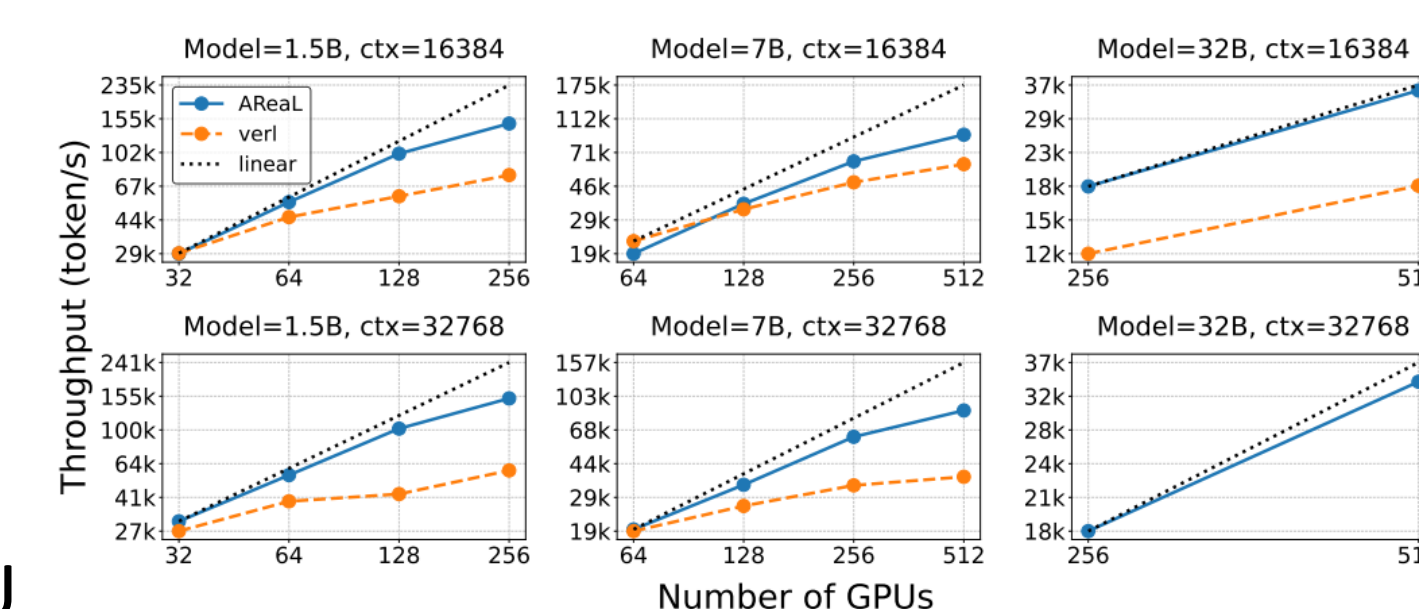
Full Trajectory

Cross-Version Trajectories

### Algorithmic Challenges

• **Data Staleness**

Data generated from old policies impedes efficient learning.

• **Inconsistent Policy Versions**

Interrupted trajectories involve segments generated from different policy versions, which violates the PPO objective.

• **Staleness-Aware Training**

Control the maximum staleness with algorithm-system co-design.

• **Decoupled PPO objective**
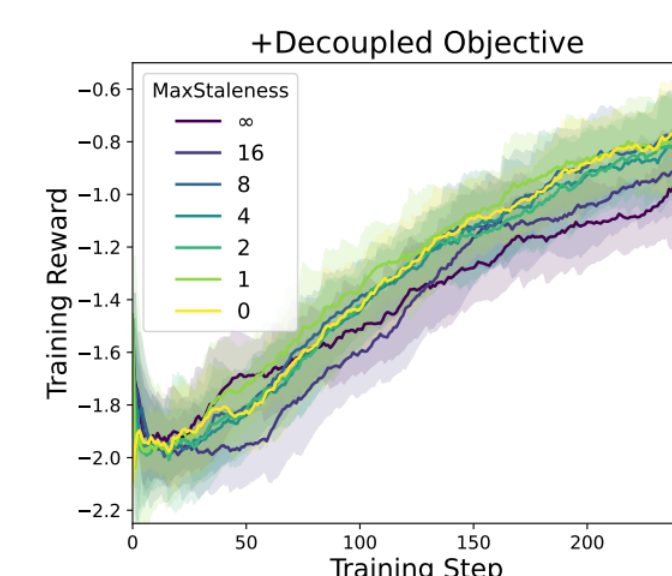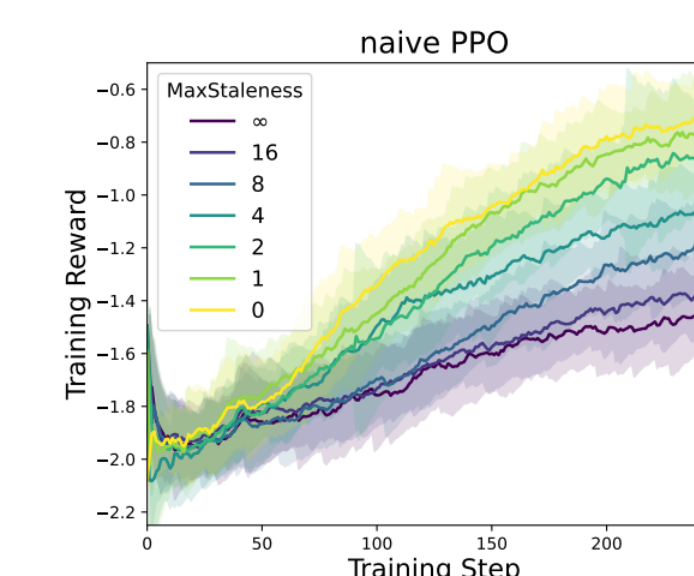
$$J(\theta) = \mathbb{E}_{q\sim\mathcal{D}, a_t\sim\pi_{\text{behav}}}\left[\sum_{t=1}^{H}\min\left(\frac{\pi_\theta}{\pi_{\text{behav}}}\hat{A}_t, \frac{\pi_{\text{prox}}}{\pi_{\text{behav}}}\text{clip}\left(\frac{\pi_\theta}{\pi_{\text{prox}}}, 1-\epsilon, 1+\epsilon\right)\hat{A}_t\right)\right]$$

Importance Ratio | Importance Ratio | Trust Region Center | Importance Ratio



Model=1.5B, ctx=16384 | Model=7B, ctx=16384 | Model=32B, ctx=16384
Model=1.5B, ctx=32768 | Model=7B, ctx=32768 | Model=32B, ctx=32768

AReaL / veri / linear

Throughput (token/s) — Number of GPUs

*2x Higher* Effective Training Throughput Compared with veRL

| Model | LiveCodeBench ↑ | Training Hours ↓ |
|---|---|---|
| 14B basemodel | 53.4 | - |
| w/ VeRL | 57.9* | 44.4 |
| w/ Sync.AReaL | 56.7 | 48.8 |
| w/ AReaL (ours) | **58.1** | **21.9** |
| 32B basemodel | 57.4 | - |
| w/ VeRL | - | 46.4 |
| w/ Sync.AReaL | **61.2** | 51.1 |
| w/ AReaL (ours) | 61.0 | **31.1** |

*2x Convergence Speed* in Competitive Programming Tasks



naive PPO | +Decoupled Objective

MaxStaleness: ∞ / 16 / 8 / 4 / 2 / 1 / 0

Training Reward — Training Step

*Both algorithm modifications are imperative to async. RL performance*