

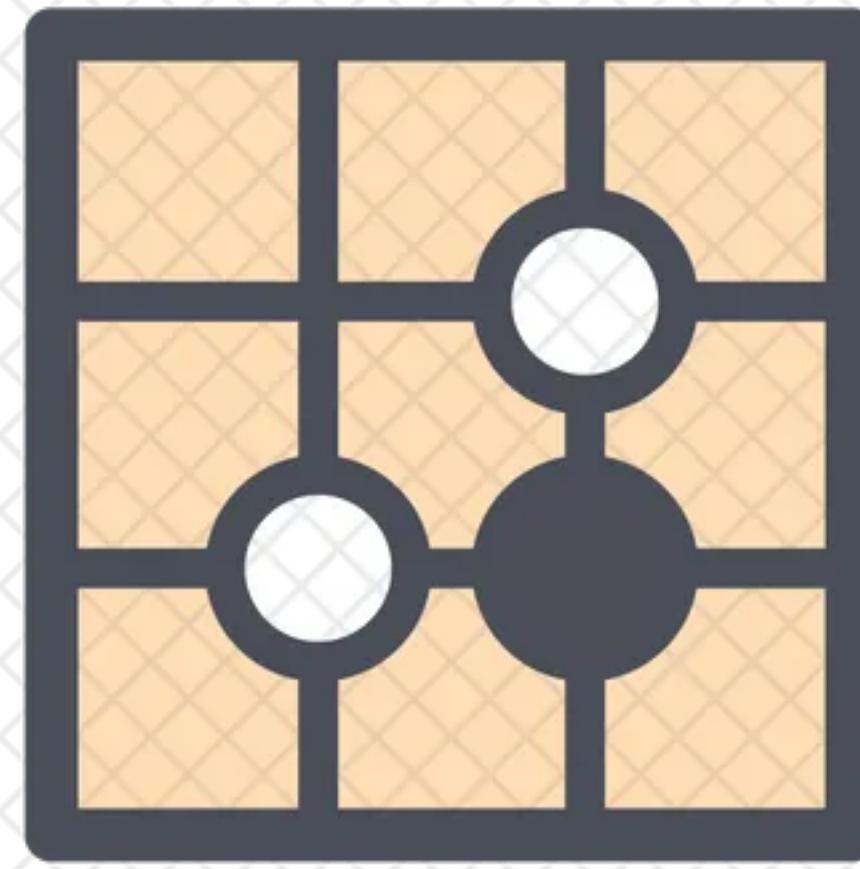
★ Best Student Paper Finalist ★

Stellaris: Staleness-Aware Distributed Reinforcement Learning with Serverless Computing

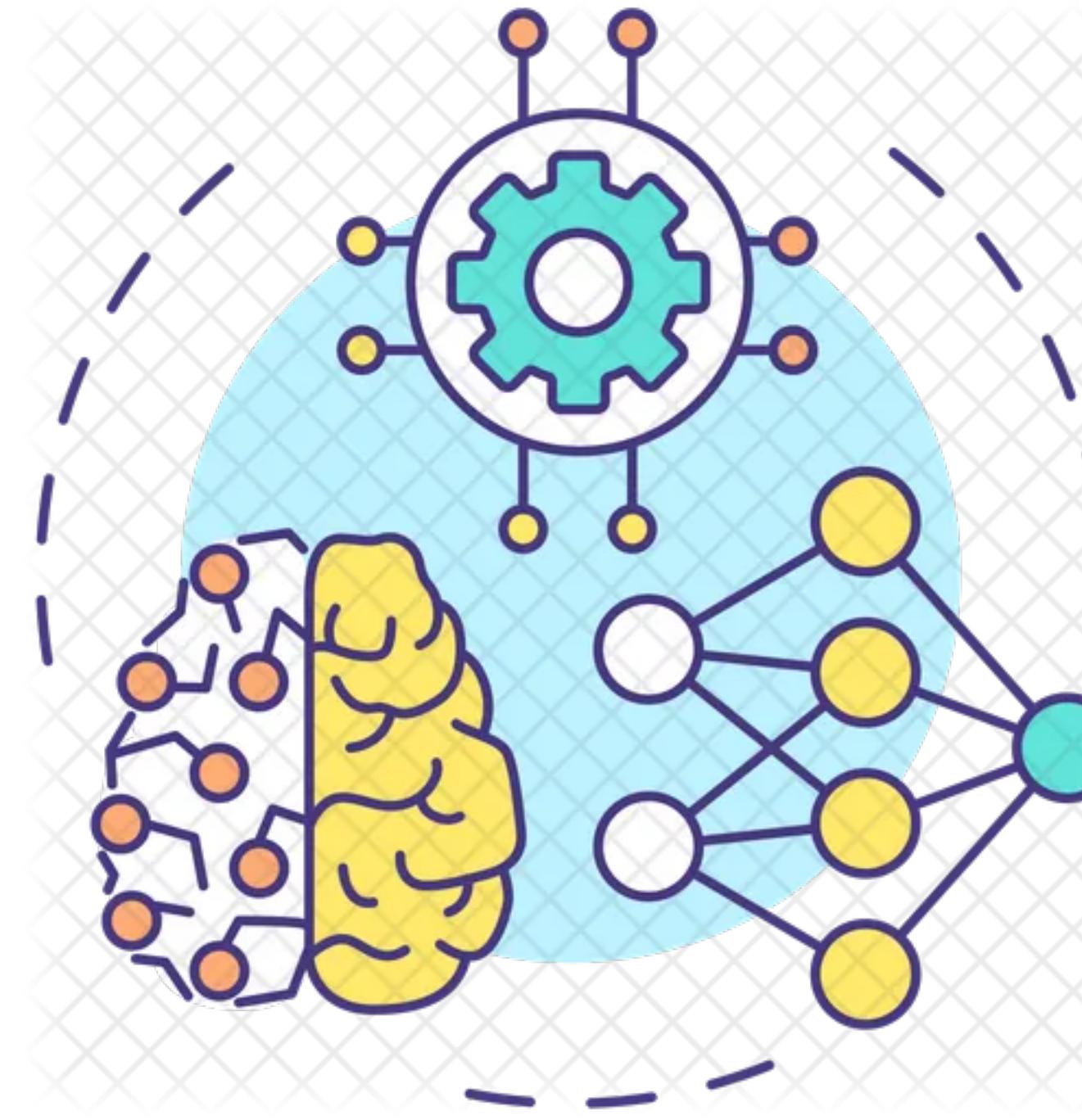
Hanfei Yu¹, Hao Wang¹, Devesh Tiwari², Jian Li³, Seung-Jong Park⁴

Stevens Institute of Technology¹, Northeastern University², Stony Brook University³, Missouri University of Science & Technology⁴

Deep Reinforcement Learning (DRL)



AlphaGo
Gaming AIs

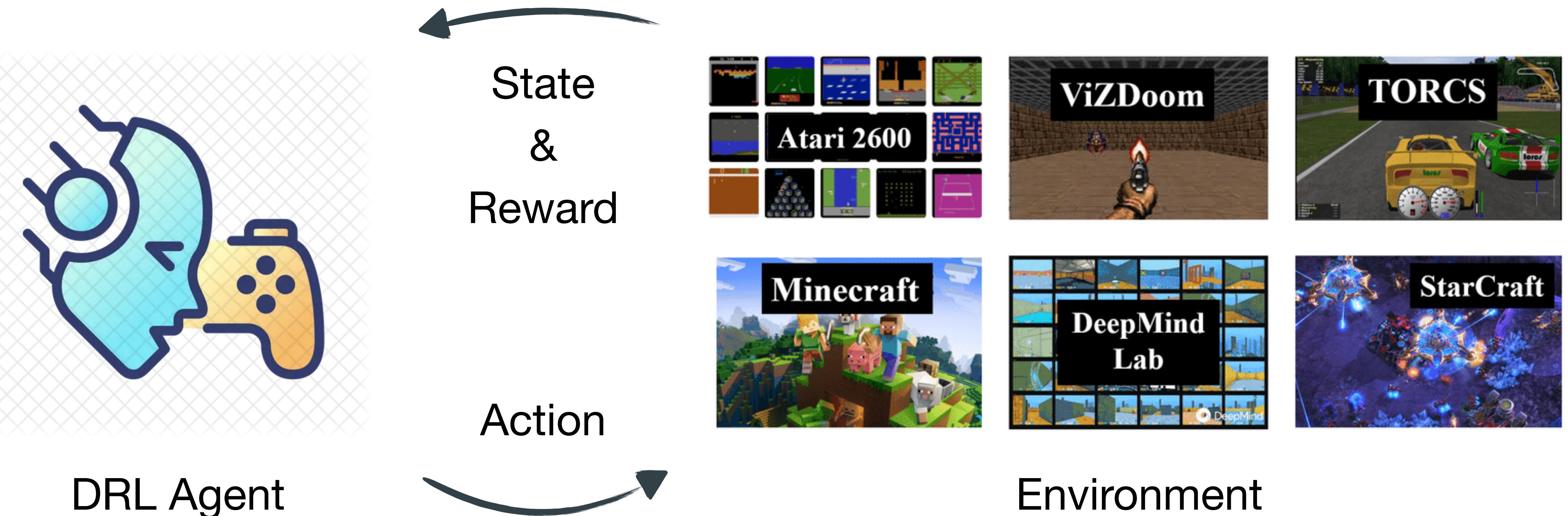


Large Language Models
RL from Human Feedback

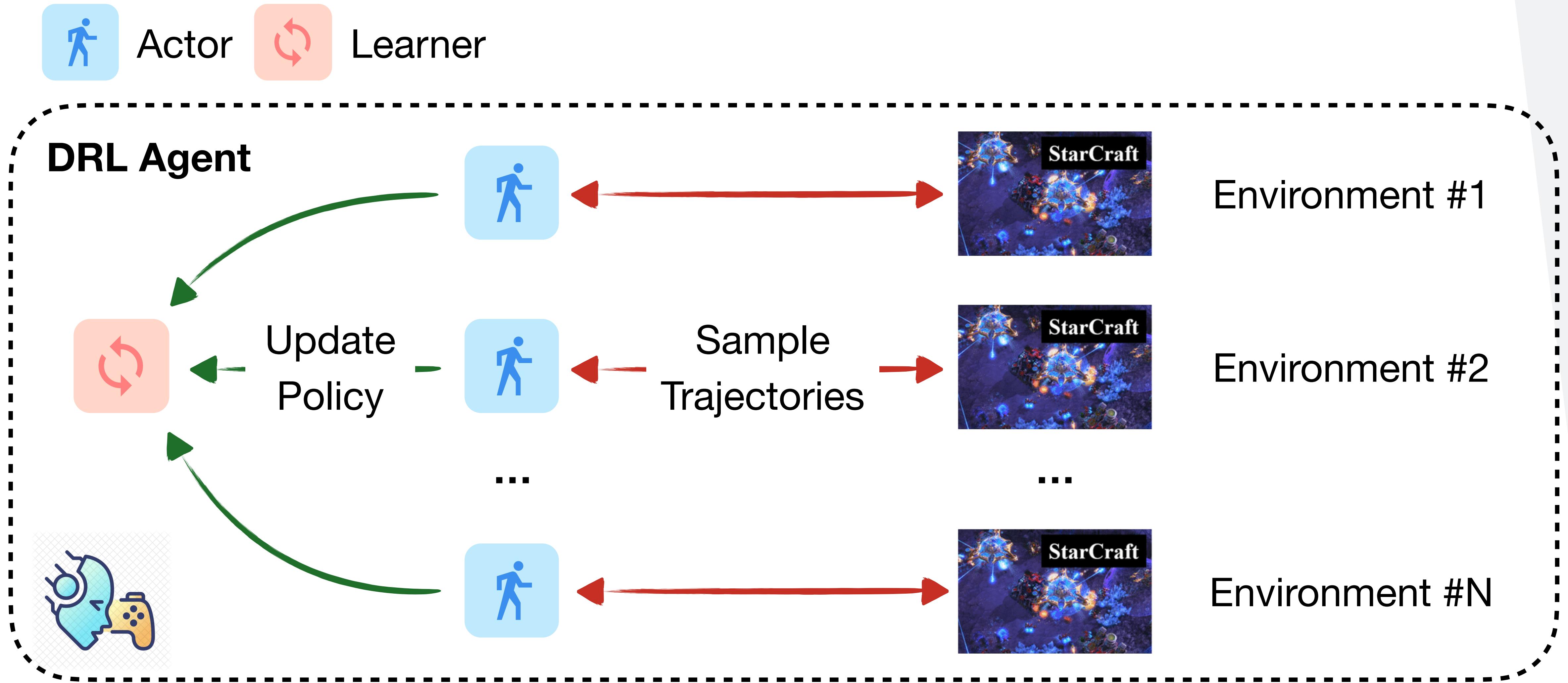


HPC System Scheduling

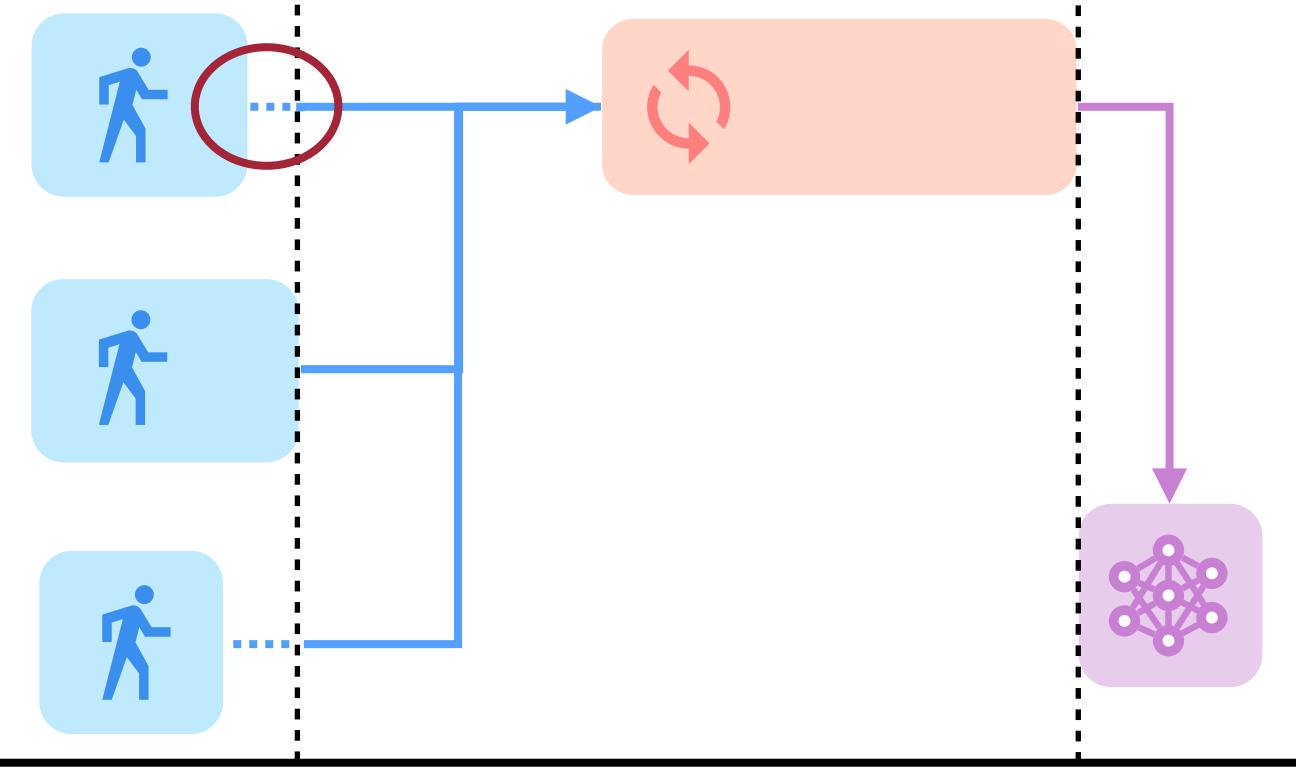
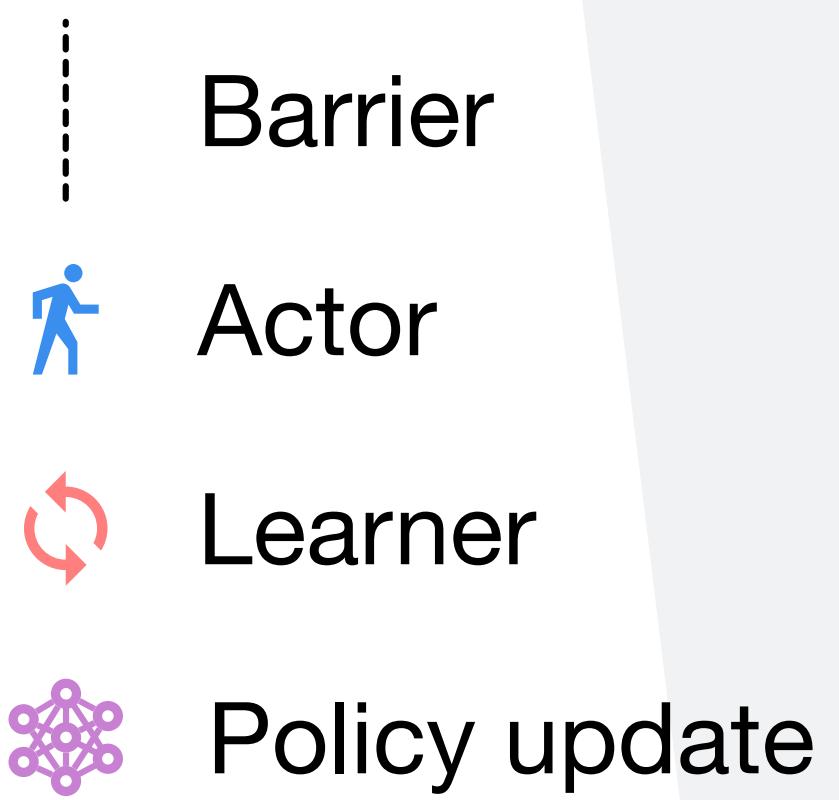
Training DRL Agent in Environments



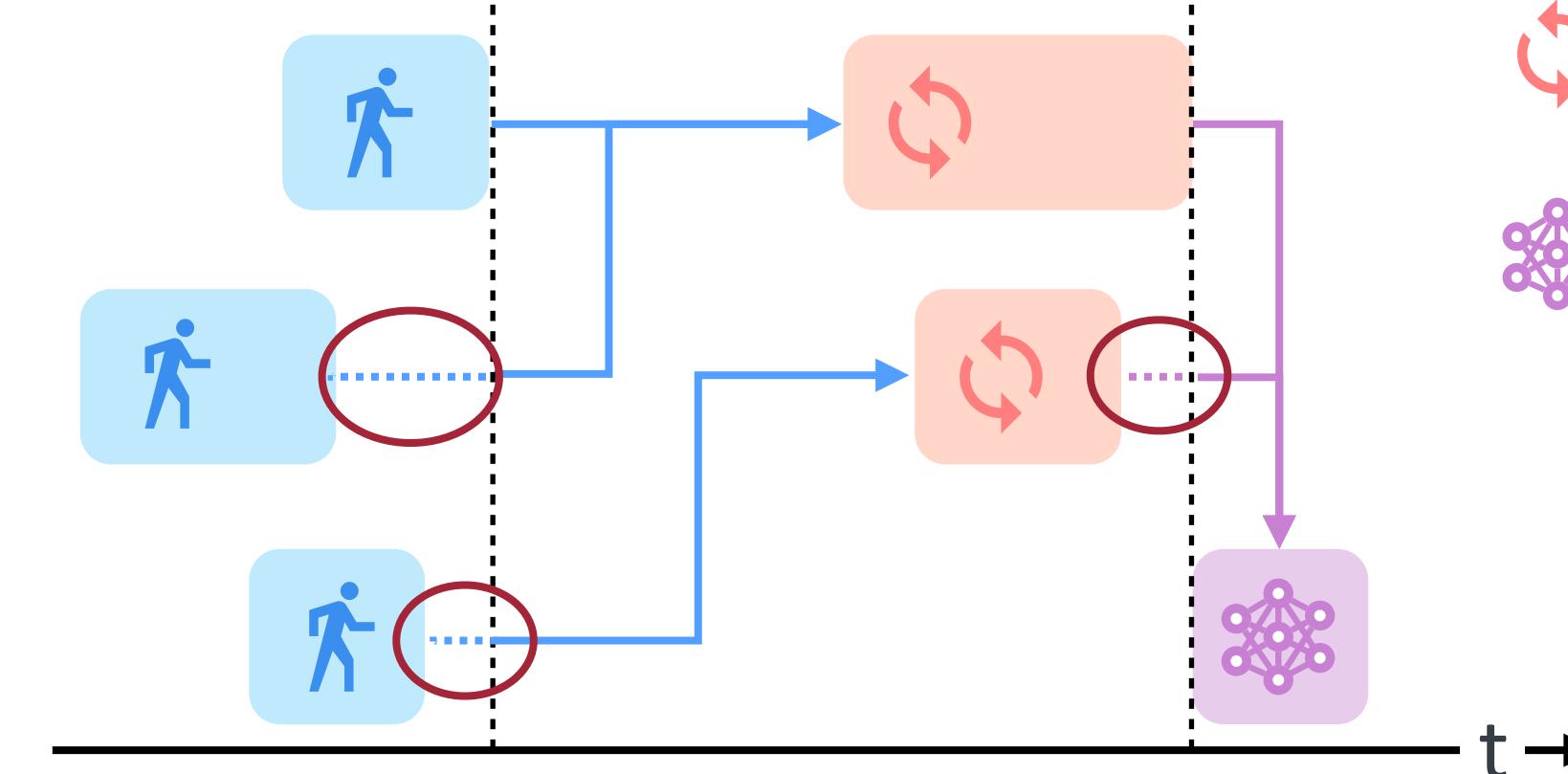
Scaling DRL Training: Actor-Learner



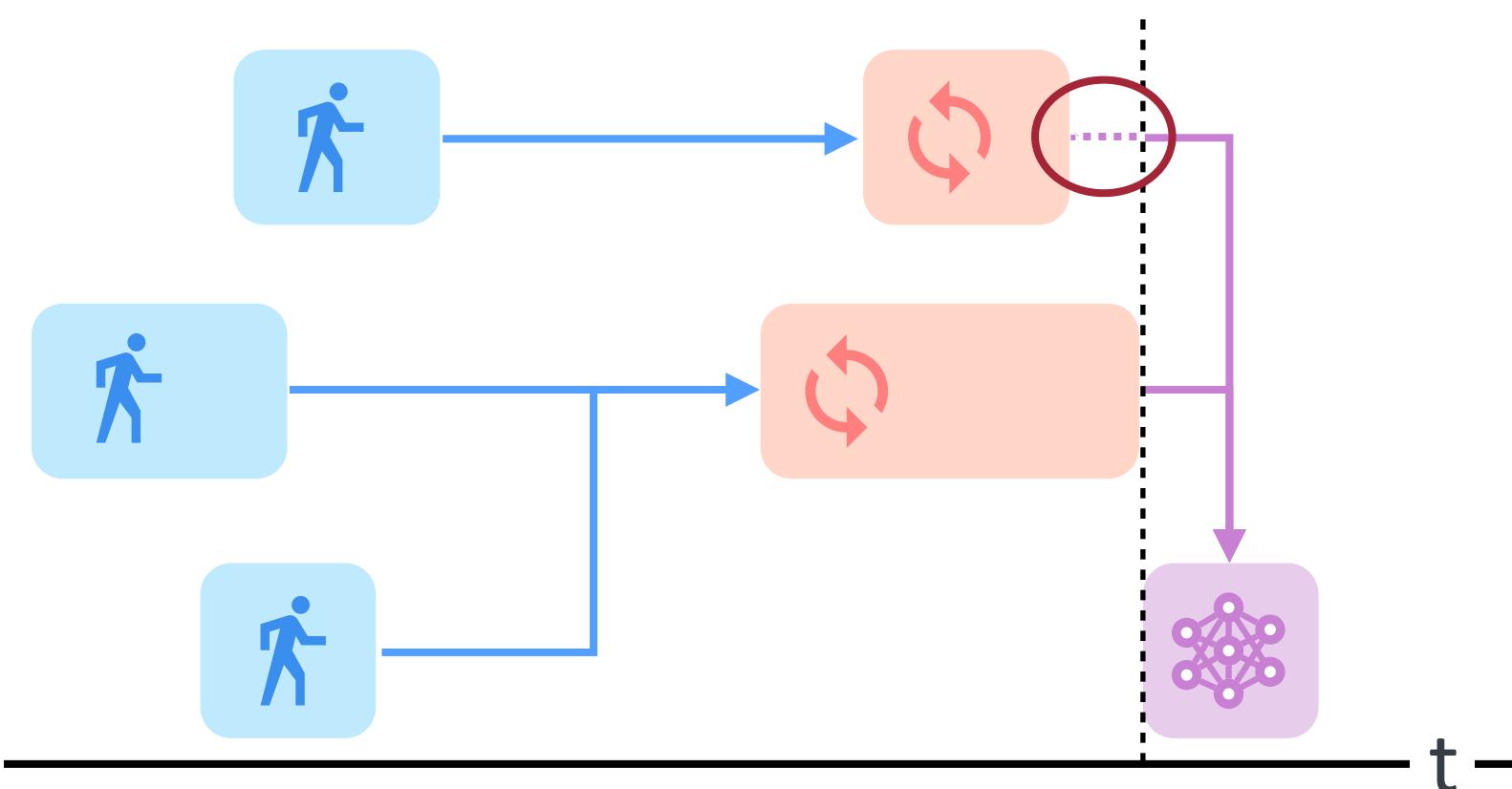
Existing Actor-Learner Architectures



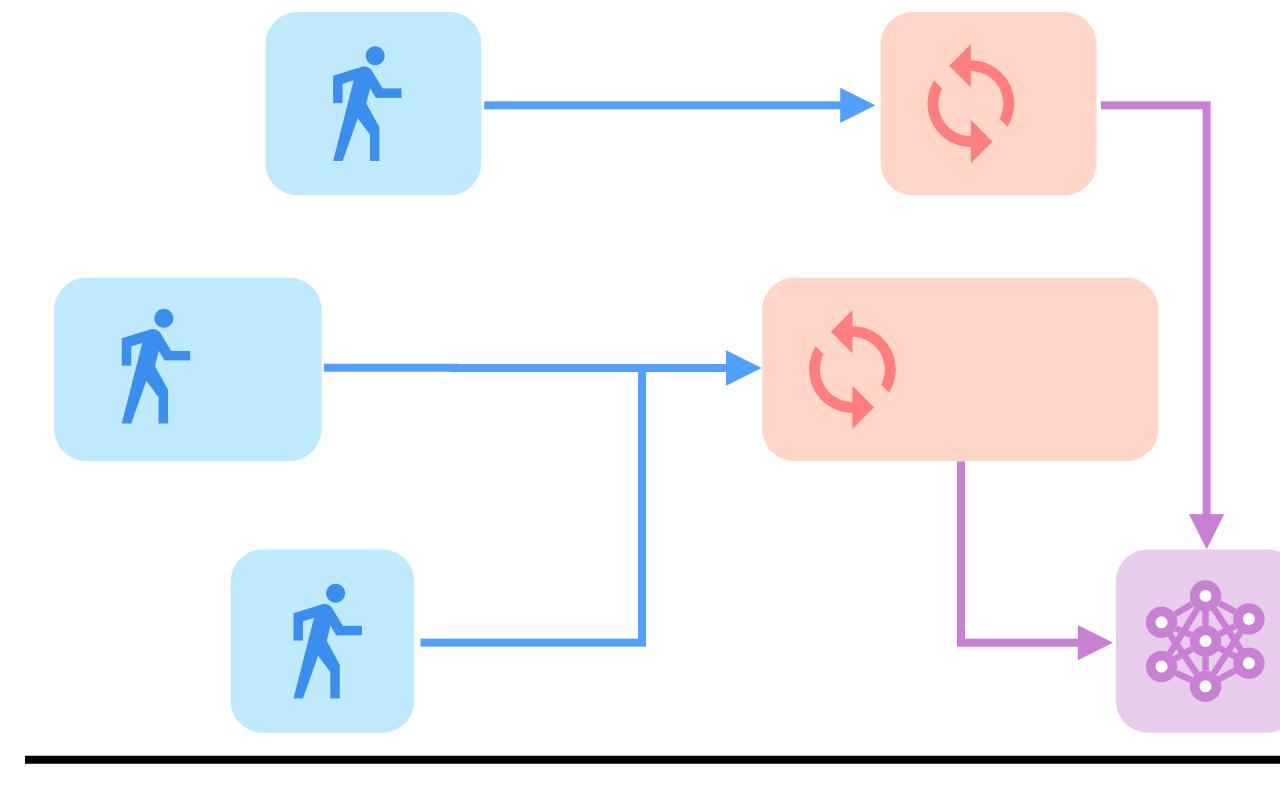
Synchronous Actors + Centralized Sync. Learner



Synchronous Actors + Synchronous Learners



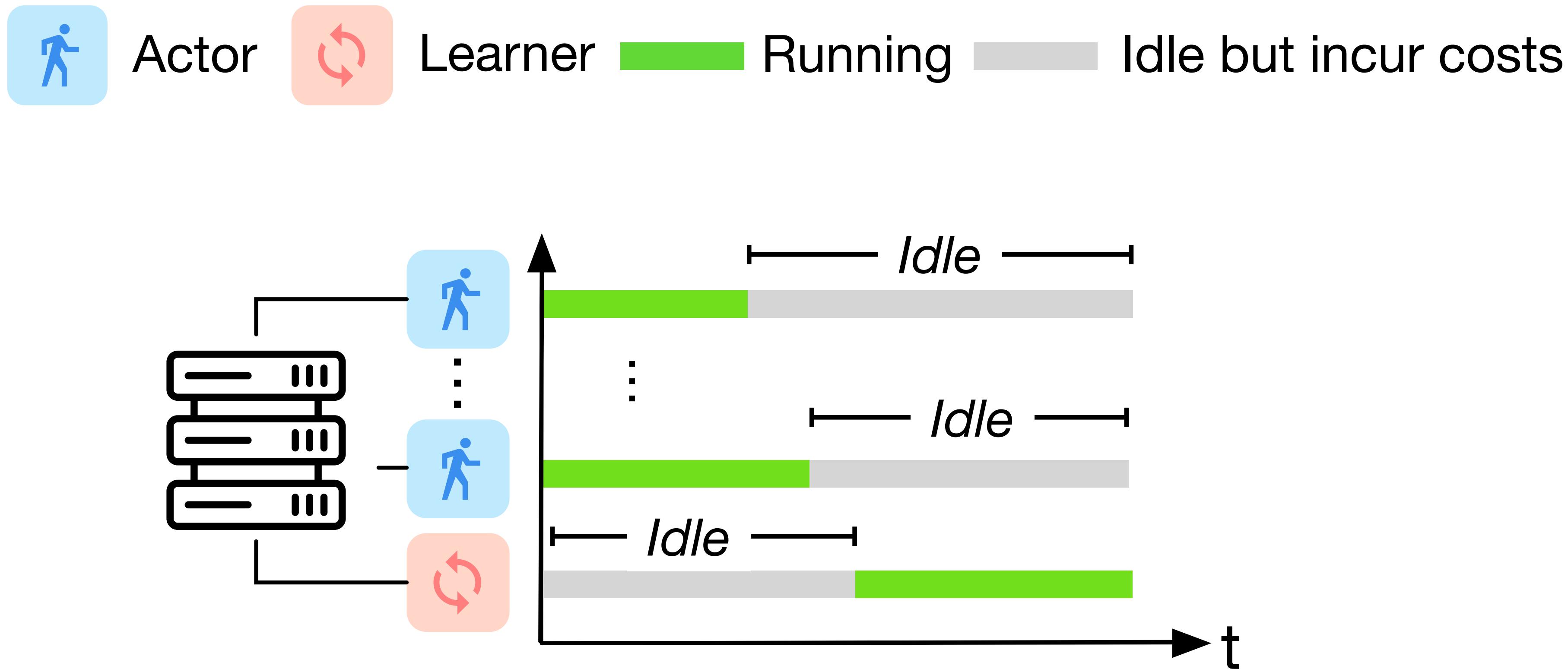
Asynchronous Actors + Synchronous Learners



(A)synchronous Actors + Asynchronous Learners

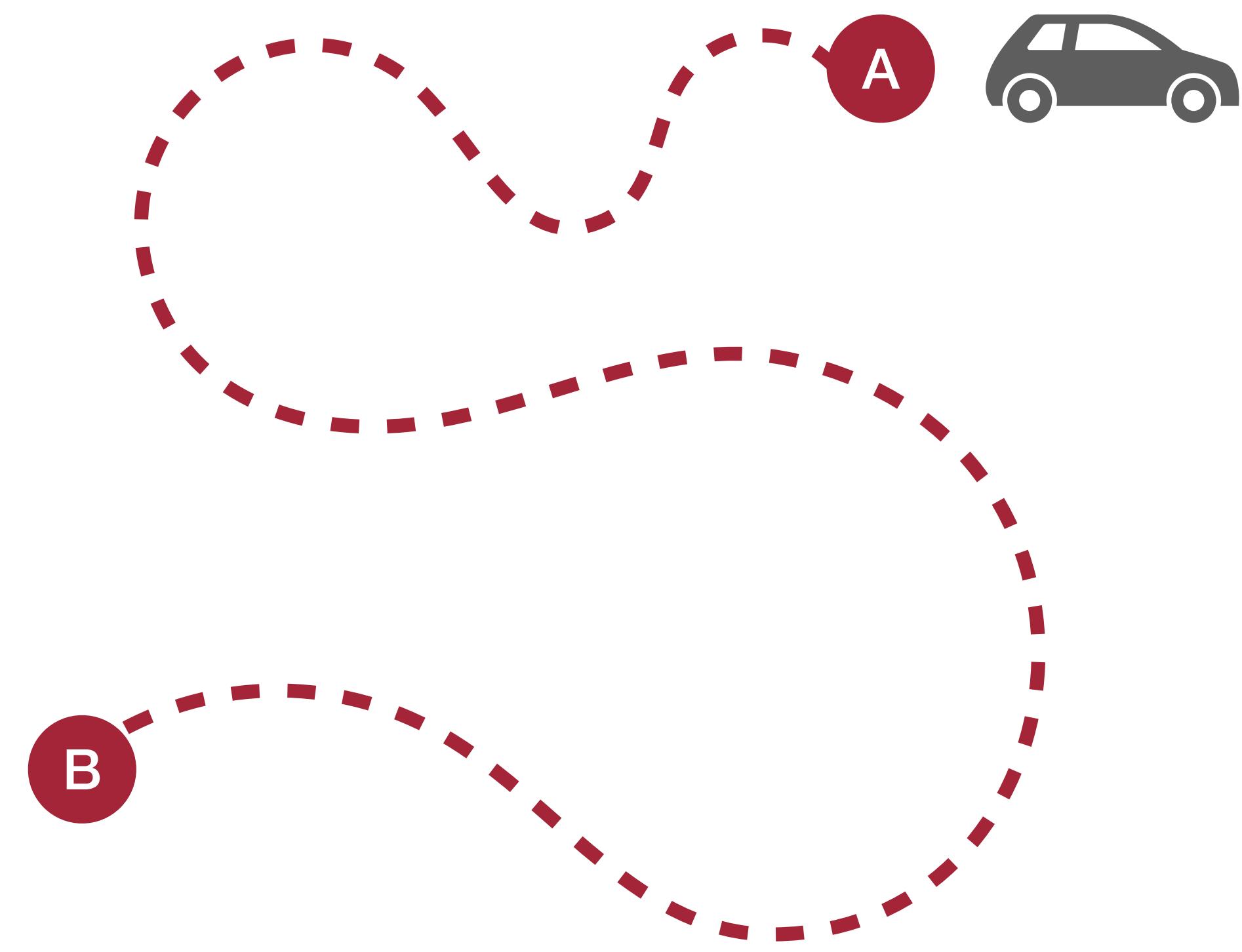
Stellaris
Fully asynchronous
Eliminate blocking
and idle waste

Dynamic Usage in Serverful DRL Training



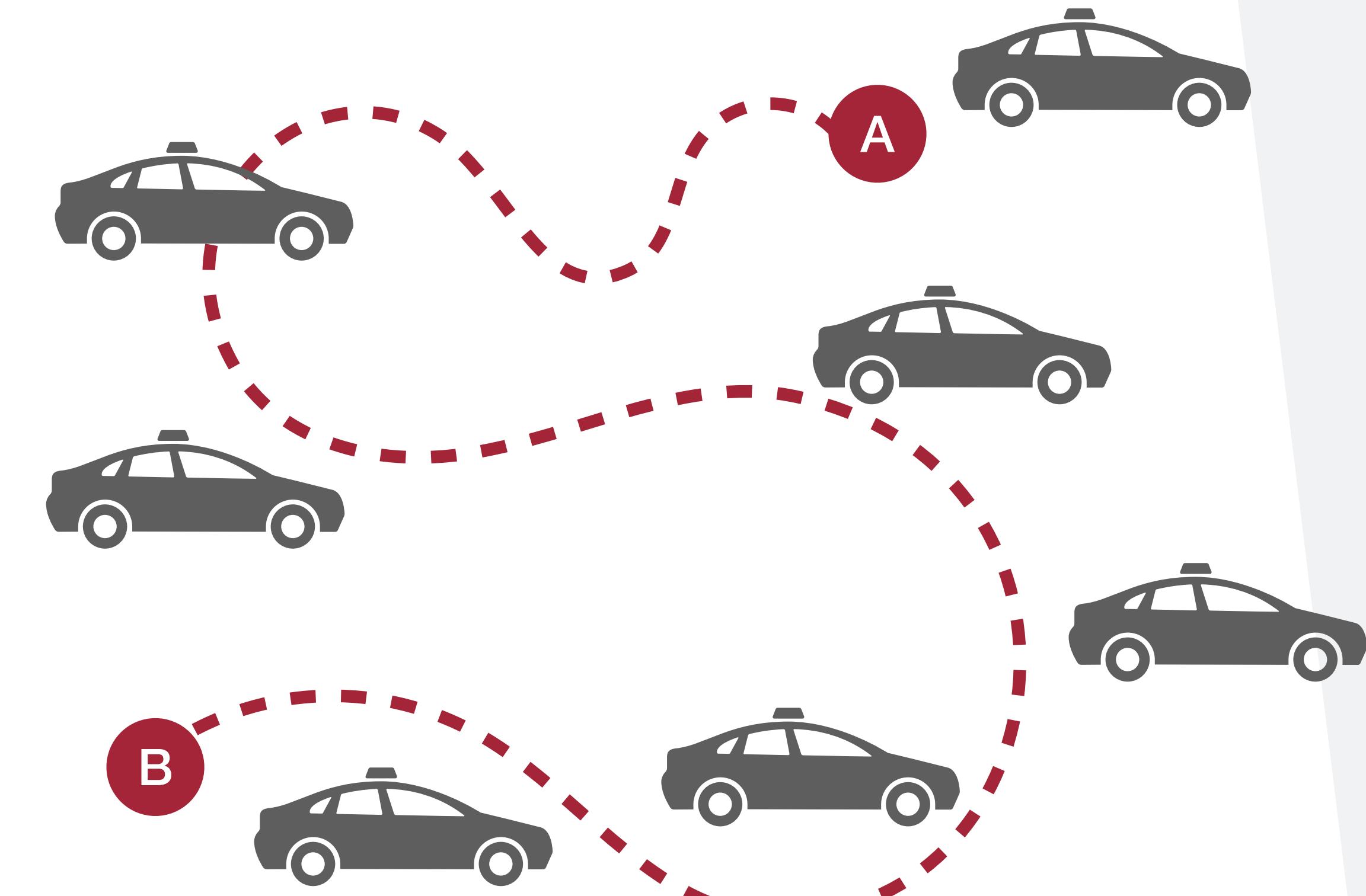
Can DRL training go **serverless**? Yes!

Cloud / HPC

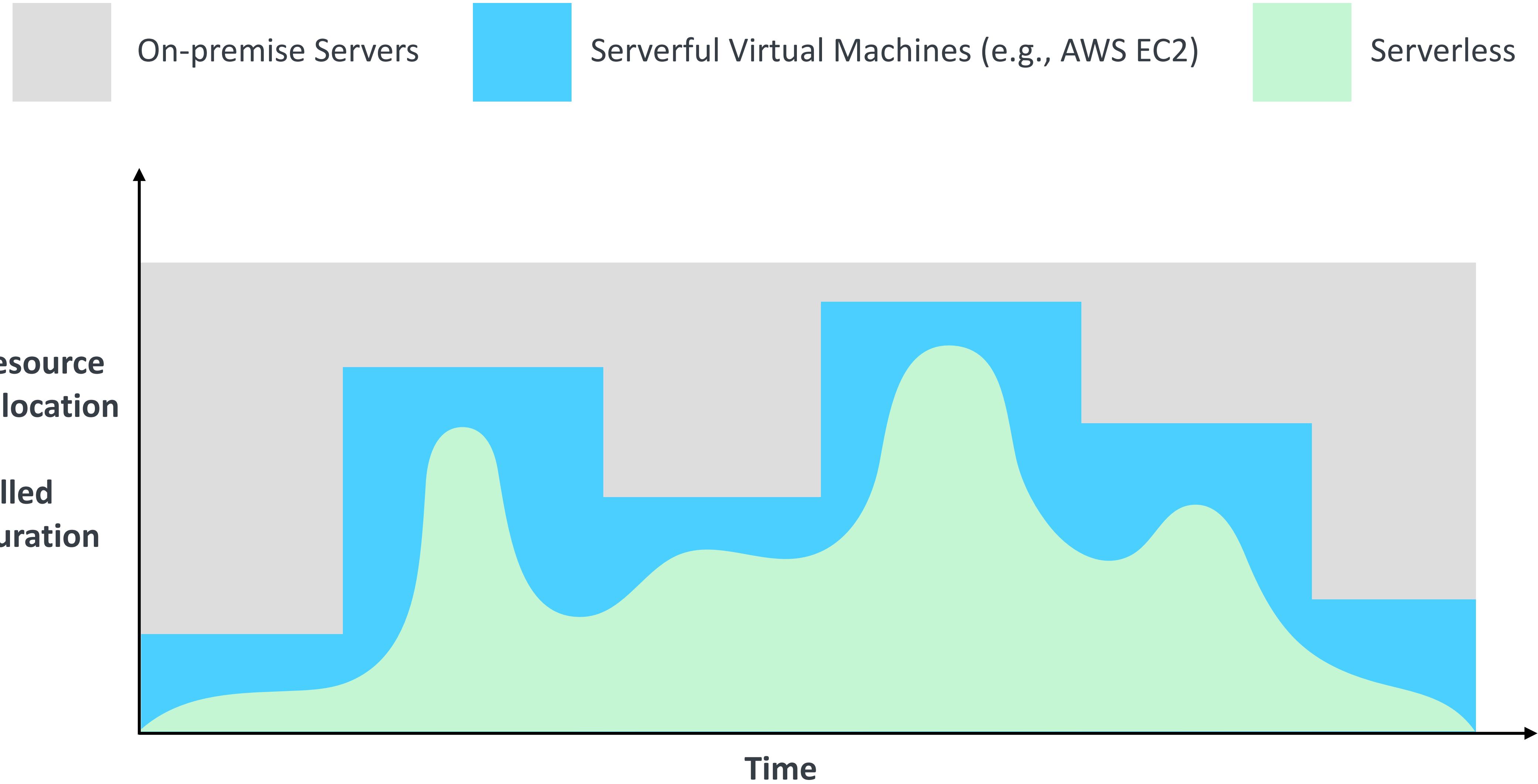


Car rental

Serverless

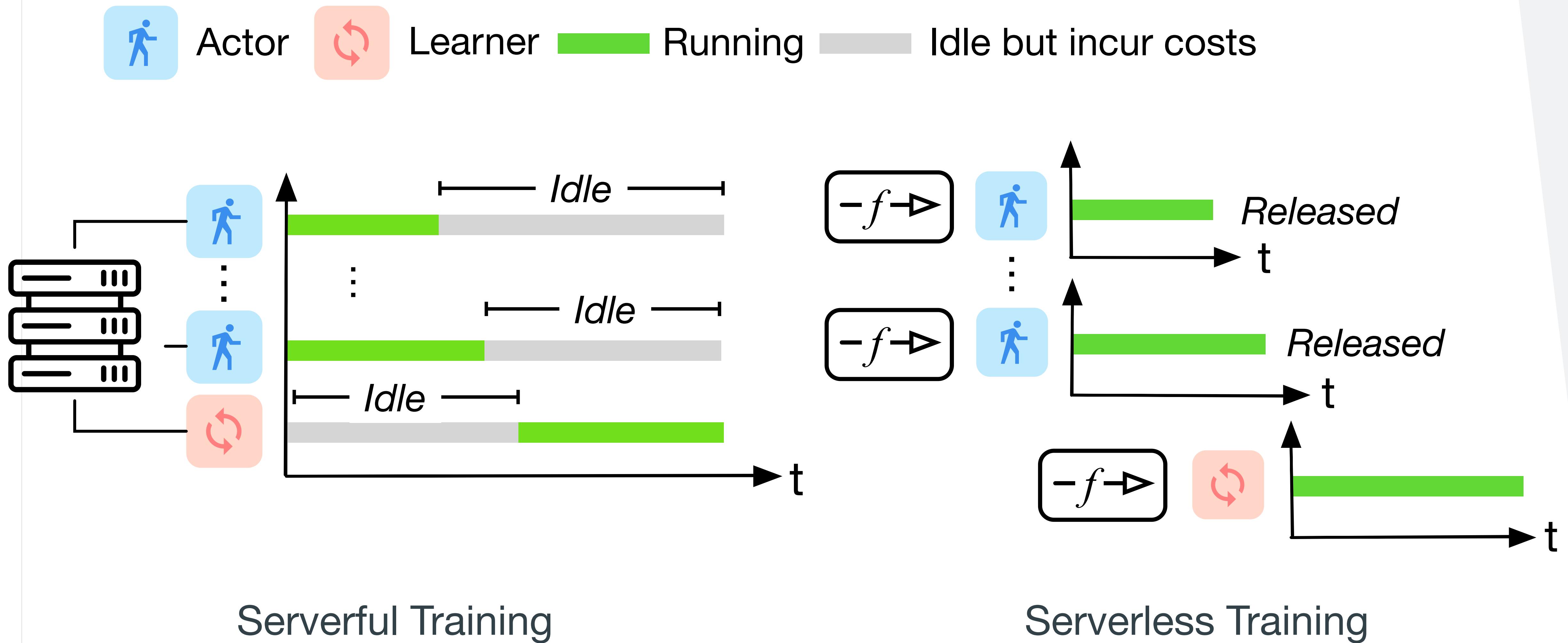


Cruise (Self-driving Taxi)



Schleier-Smith, Johann, et al. "What serverless computing is and should become: The next phase of cloud computing." Communications of the ACM 64.5 (2021): 76-84.

Serverful vs. Serverless DRL Training

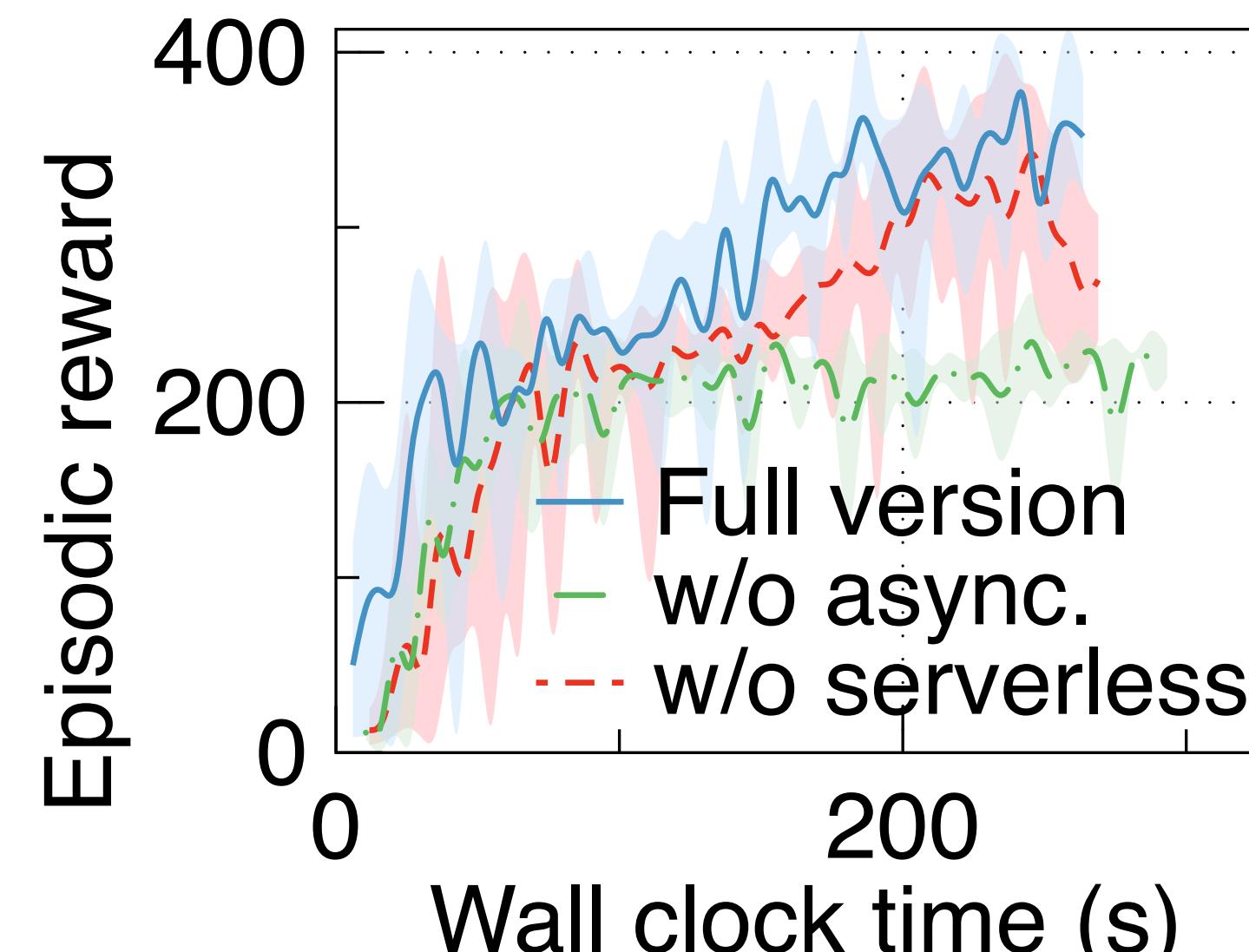


Serverful Training

Serverless Training

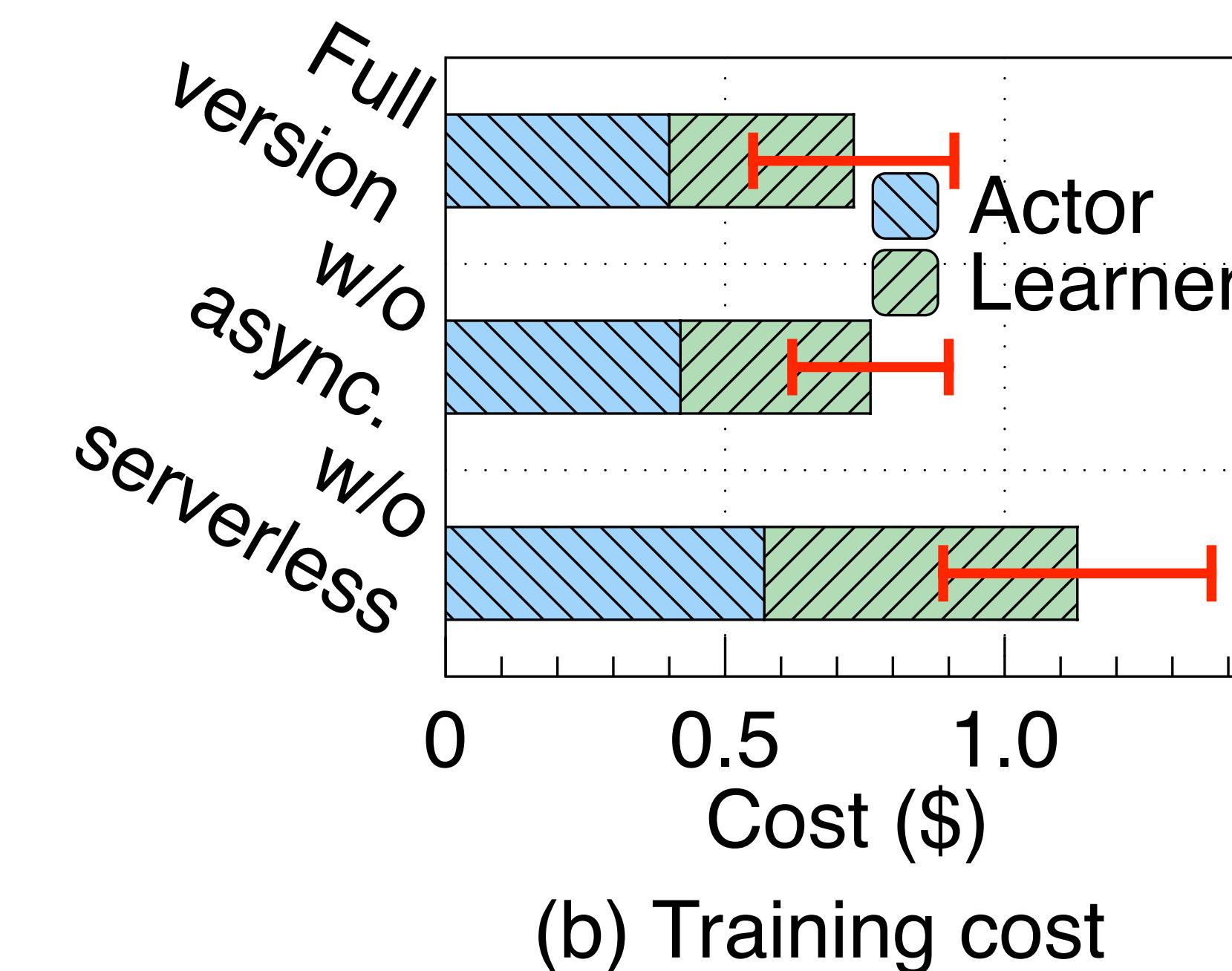
Asynchronous and Serverless DRL Training

High training efficiency



(a) Training performance

Low training cost



(b) Training cost

Proximal Policy Optimization (PPO) on MuJoCo Hopper-v4

Existing Works

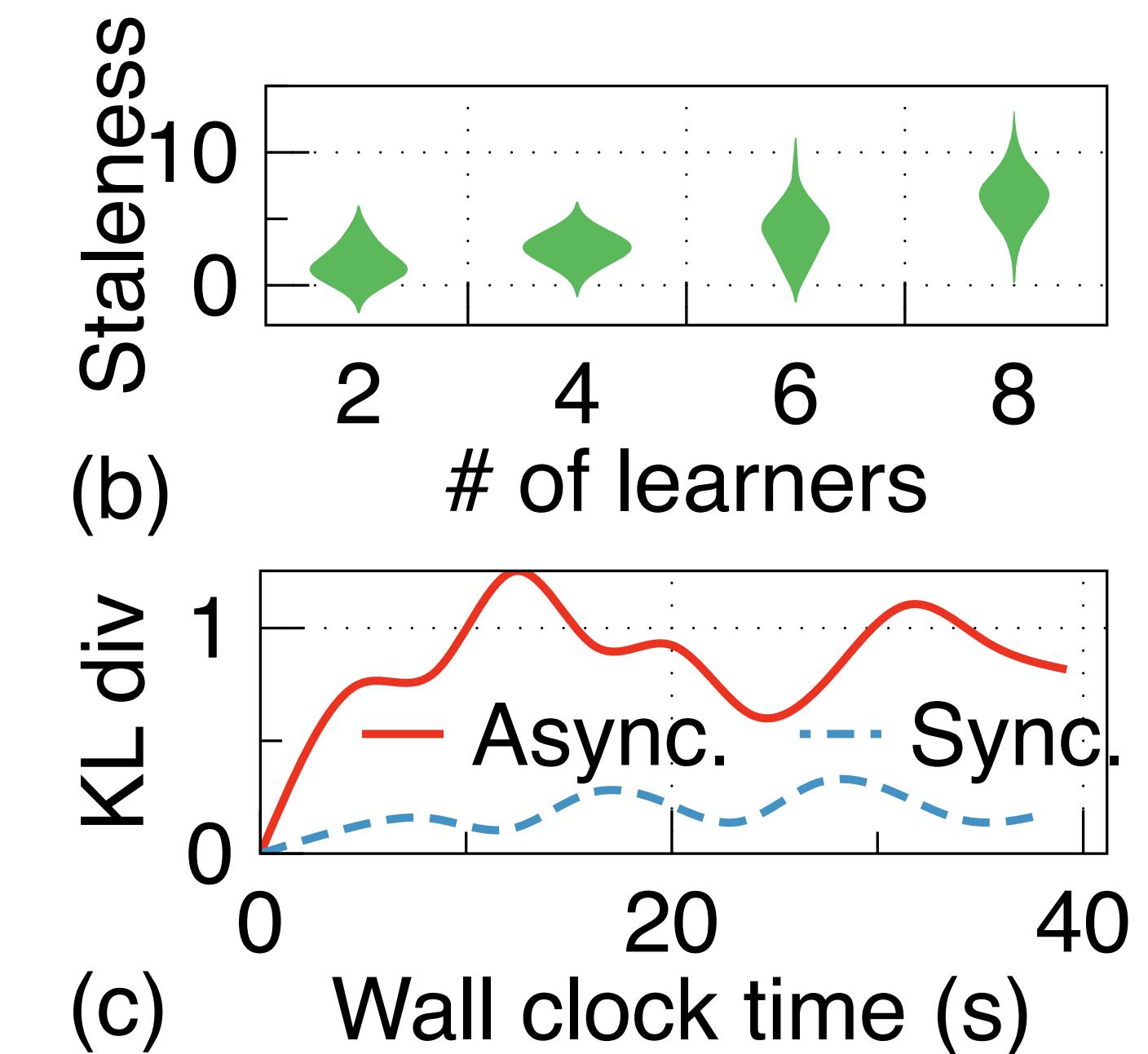
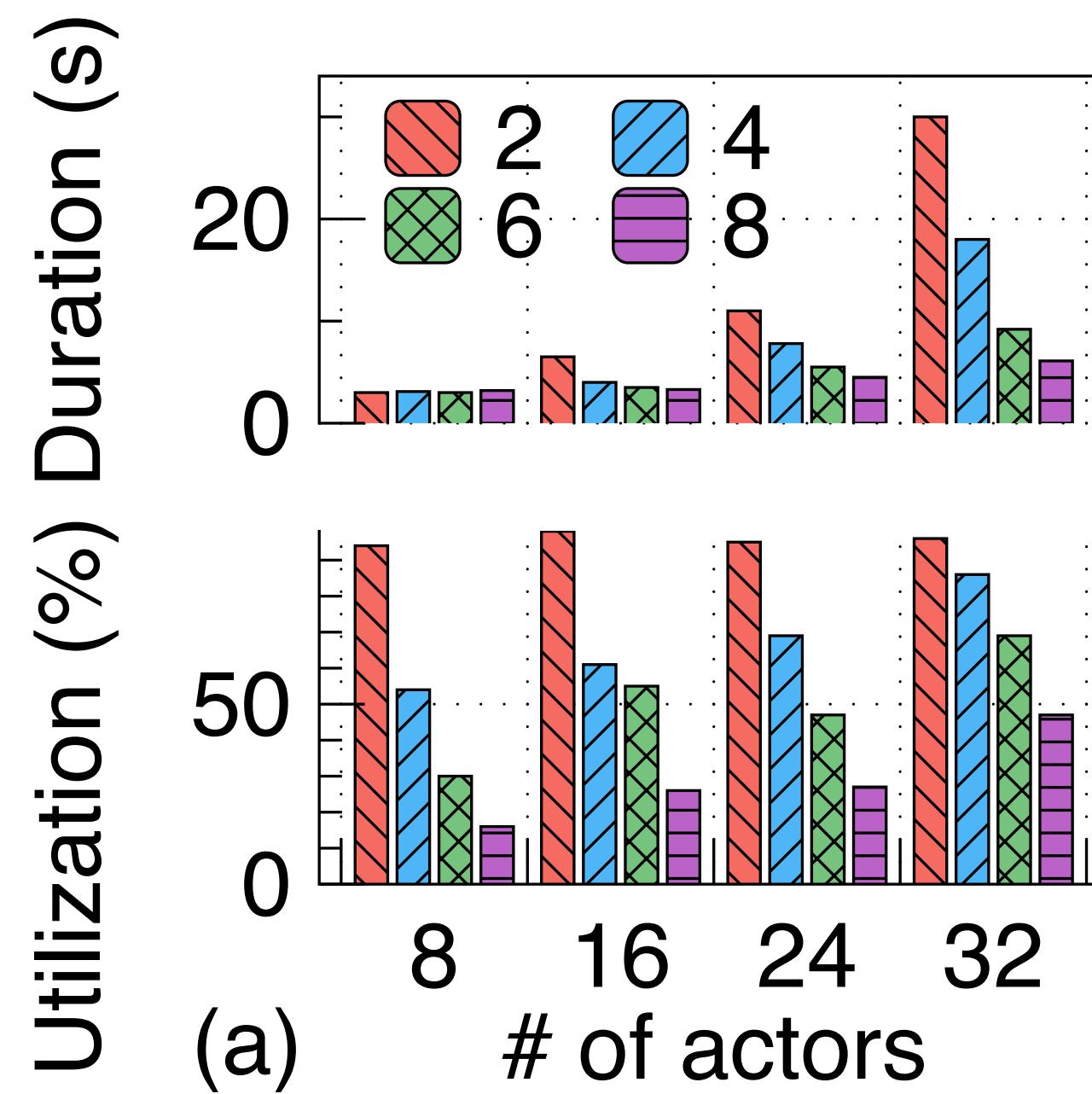
Framework	Asynchronous Learners	Scalable Actors	On-policy and Off-policy	Serverless
Ray RLlib (ICML 2018)	✗	✗	✓	✗
MSRL (ATC 2023)	✗	✗	✓	✗
SEED RL (ICLR 2020)	✗	✗	✓	✗
SRL (ICLR 2024)	✗	✗	✗	✗
MinionsRL (AAAI 2024)	✗	✓	✗	✓
Stellaris (SC 24)	✓	✓	✓	✓

Challenges

Dynamic learner orchestration

Dynamic staleness

Unstable policy updates



KL div: Kullback–Leibler divergence

Design Goals

Dynamic learner
orchestration



On-Demand Serverless Learners

Dynamic staleness



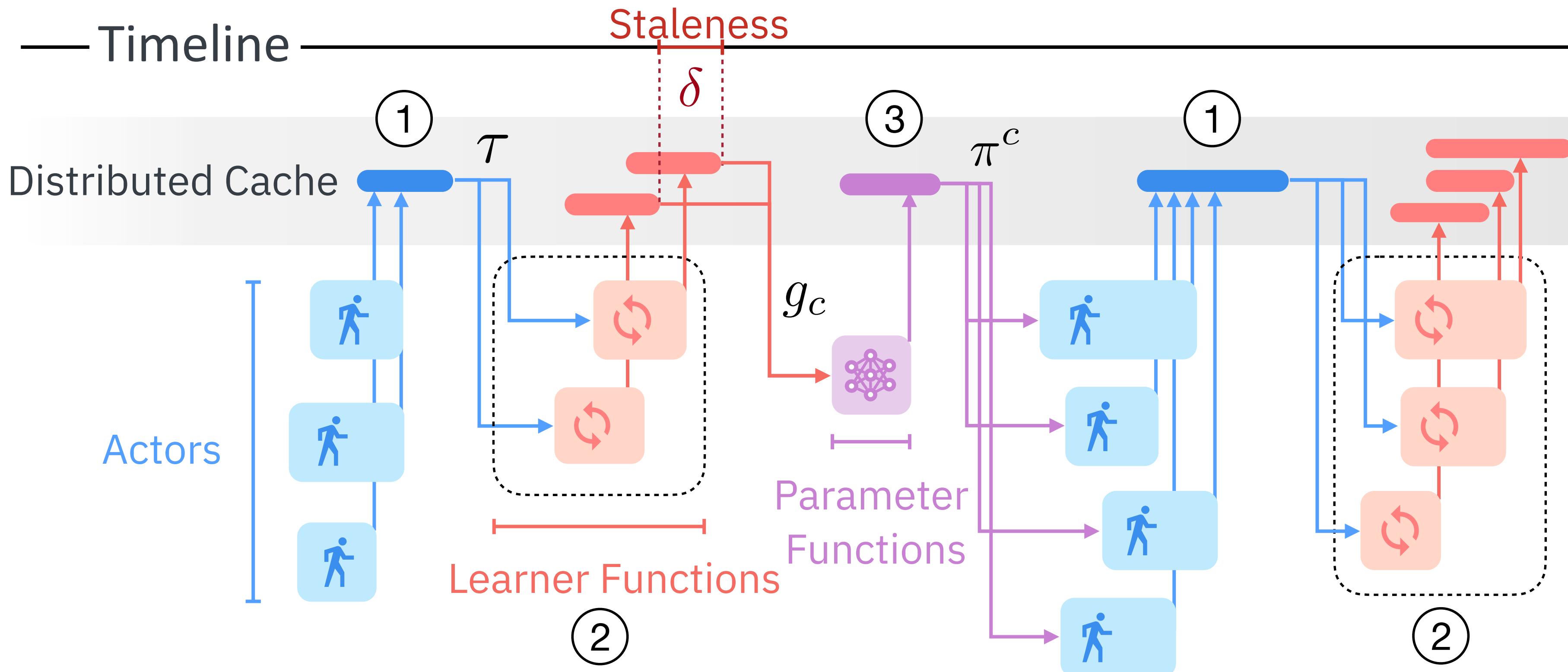
Staleness-Aware Gradient
Aggregation

Unstable policy updates



Global Importance Sampling
Truncation

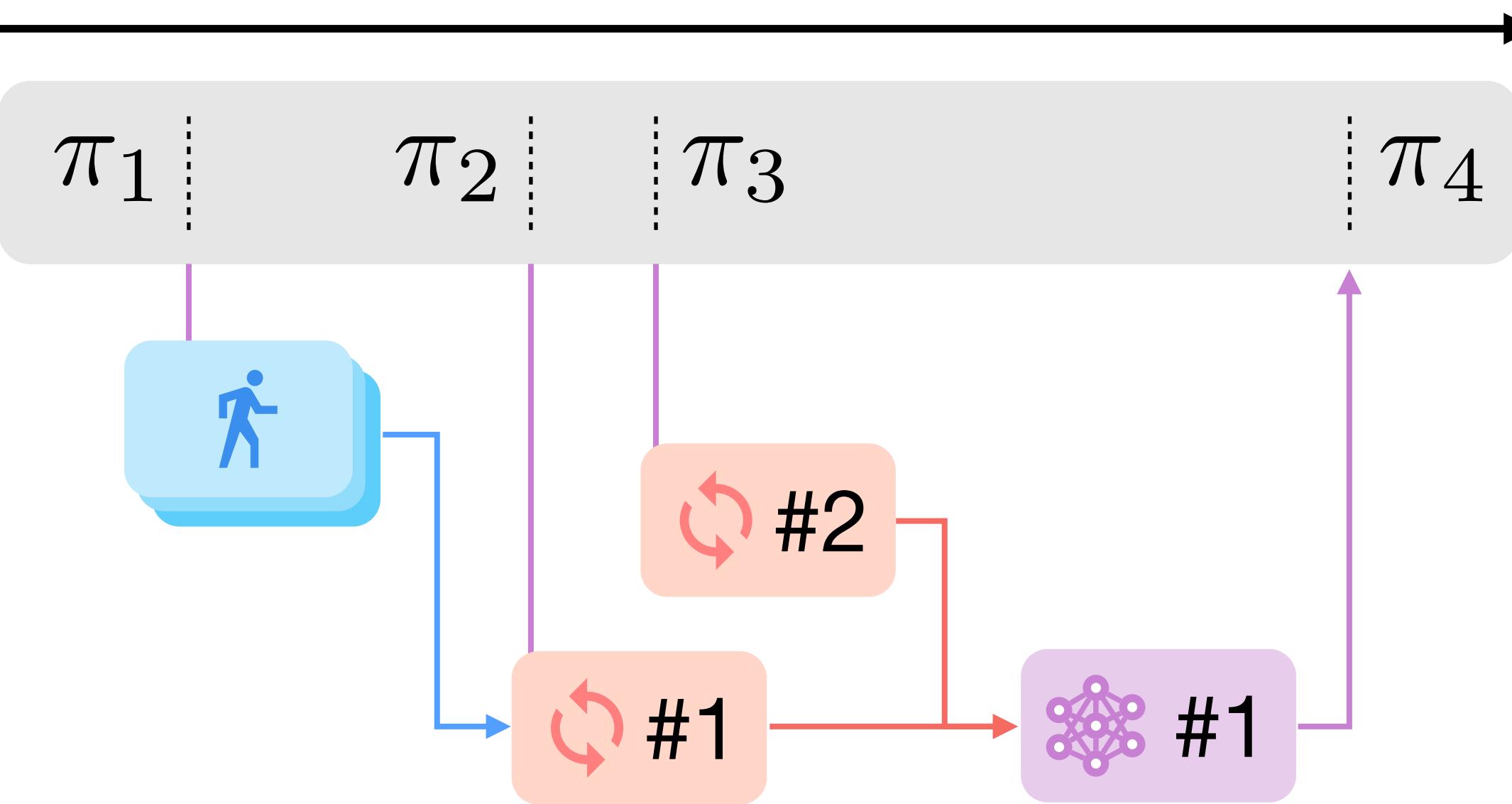
Stellaris



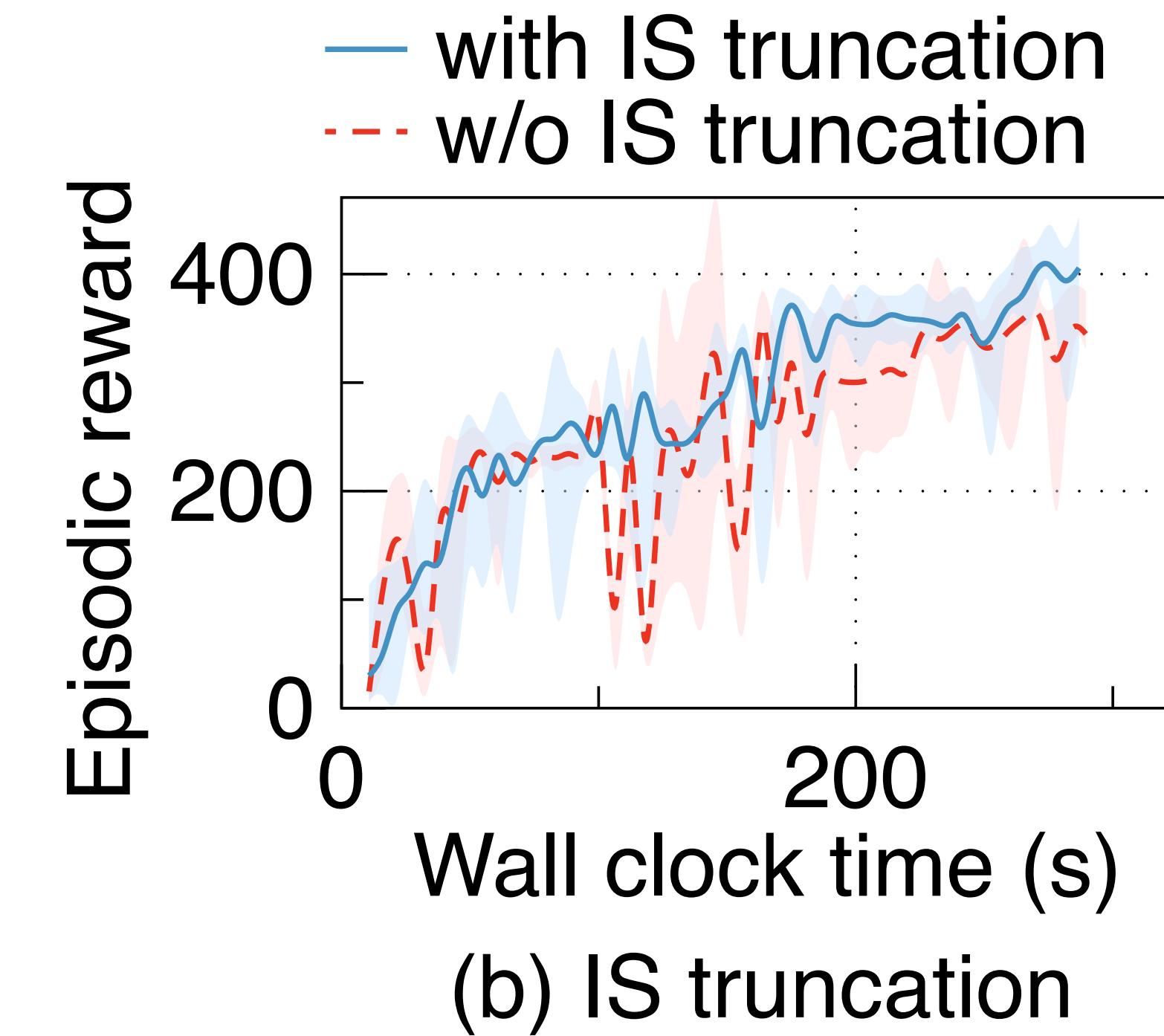
- ① Importance sampling driven trajectory collection
- ② On-demand gradient calculation
- ③ Staleness-aware gradient aggregation

Global Importance Sampling Truncation

Policy Version Timeline



$$P \text{ Truncate } \left(\frac{\pi_2}{\pi_1}, \frac{\pi_3}{\pi_1} \right)$$



Actor

Learner Function



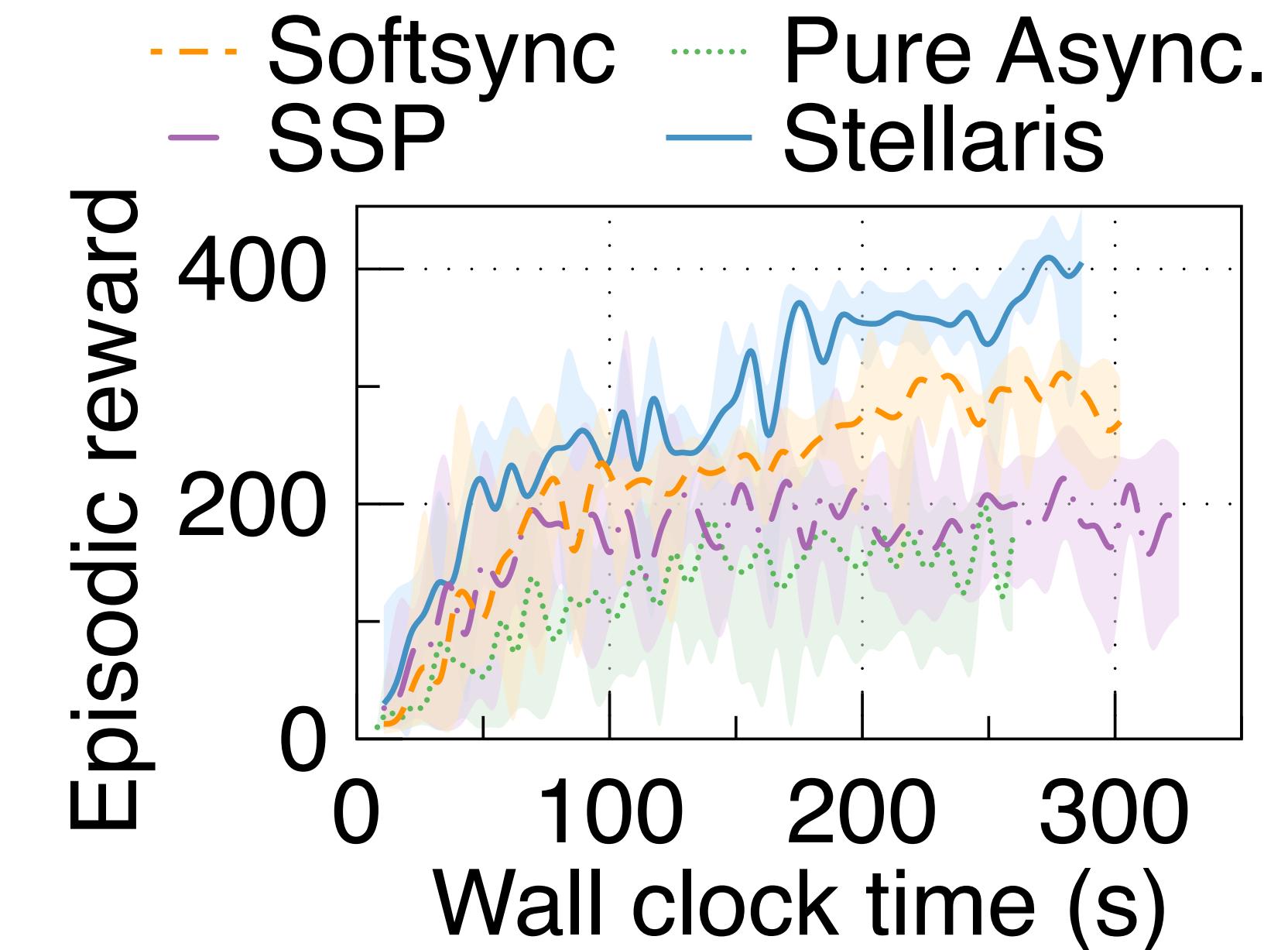
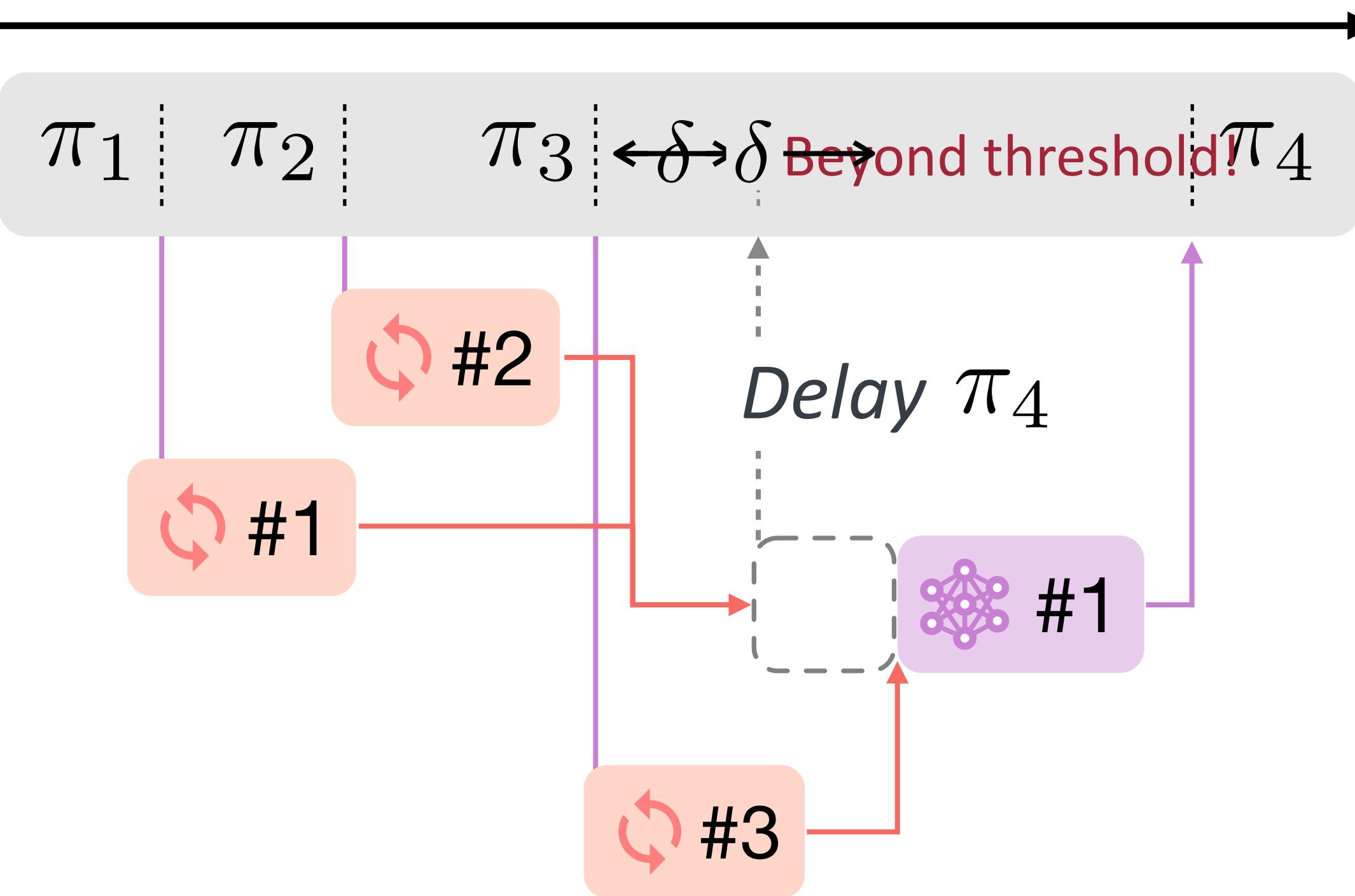
Parameter Function



Delay

Staleness-Aware Gradient Aggregation

Policy Version Timeline



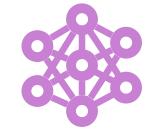
(a) Gradient aggregation



Actor



Learner Function



Parameter Function



Delay

Theoretical Guarantees

Importance Sampling Truncation

Lower bound on
Monotonic reward improvement

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}(\|\nabla J(\theta_t)\|^2) \leq 2\sqrt{\frac{2C_1C_2}{Tb}}$$

Gradient Aggregation

$\mathcal{O}\left(\frac{1}{\sqrt{Tb}}\right)$ near-linear convergence rate

$$J(\pi_i) - J(\mu) \geq -\frac{\gamma\epsilon^{\pi_i}\sqrt{2\log\rho}}{(1-\gamma)^2}$$

Please refer details to our paper

Implementation

Ray RLlib
Docker Containers
AWS EC2

Metrics

Episodic reward
Training cost

Baselines

Ray RLlib [1]
MinionsRL [2]

MuJoCo [3]

Hopper
Humanoid
Walker2d

Evaluation

Benchmarks

Atari [4]

Gravitar
SpaceInvaders
Qbert

[1] Liang, E.; et al. RLlib: Abstractions for Distributed Reinforcement Learning. ICML 2018

[2] H. Yu; et al. Cheaper and Faster: Distributed Deep Reinforcement Learning with Serverless Computing. AAAI 2024

[3] Todorov, E.; et al. Mujoco: A Physics Engine for Model-based Control. IROS 2012

[4] Marc, B.; et al. The Arcade Learning Environment: An Evaluation Platform for General Agents. JAIR 2013

Testbed Clusters

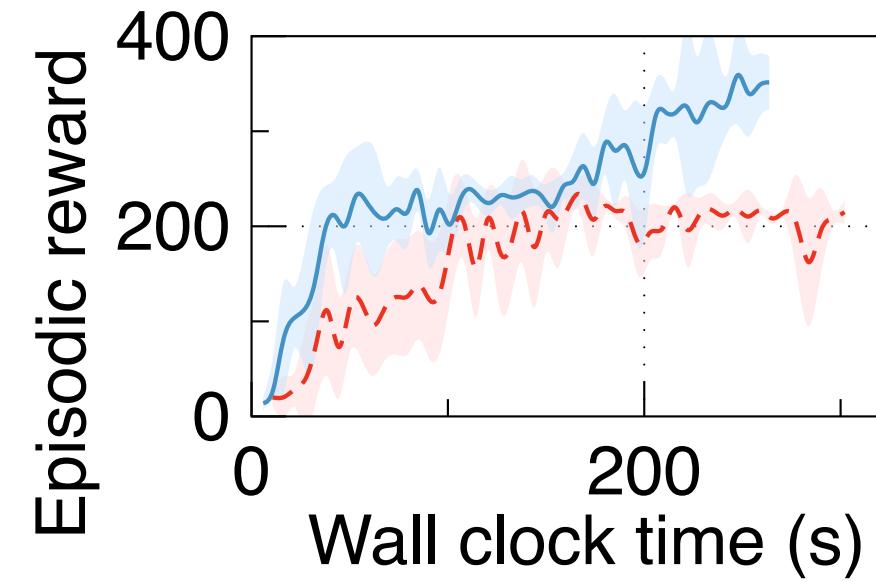
GPU Testbed

3 nodes
128 AMD EPYC 7R13 CPU cores
2 V100 GPUs

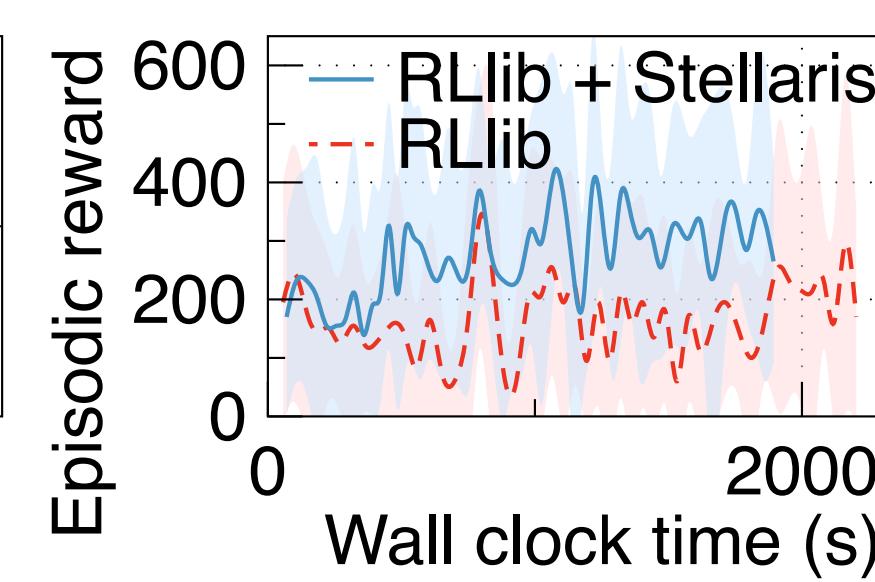
HPC Testbed

7 nodes
960 AMD EPYC 9R14 CPU cores
16 V100 GPUs

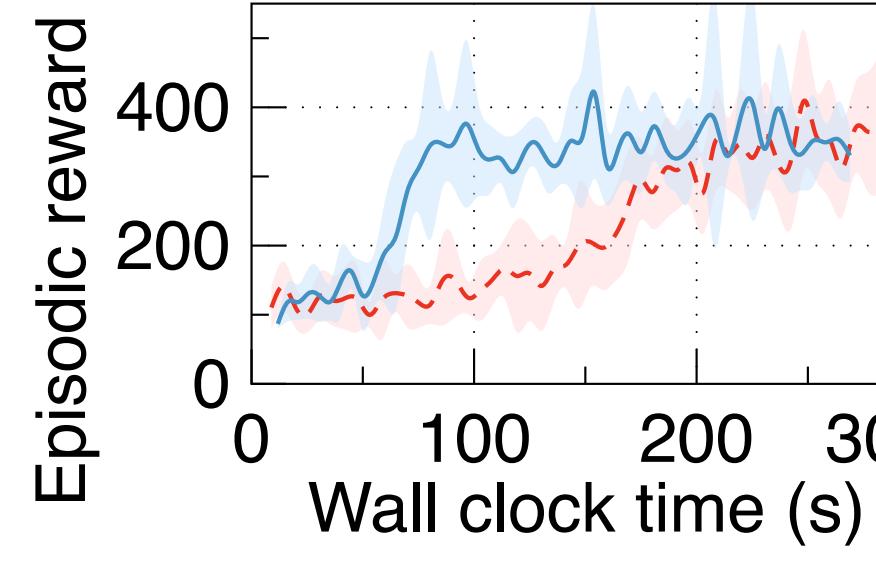
Training Performance



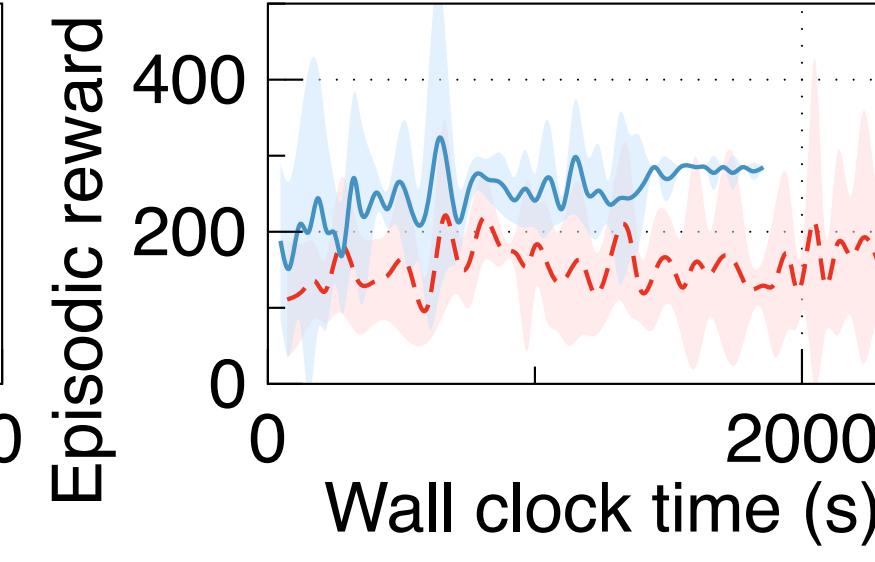
(a) Hopper



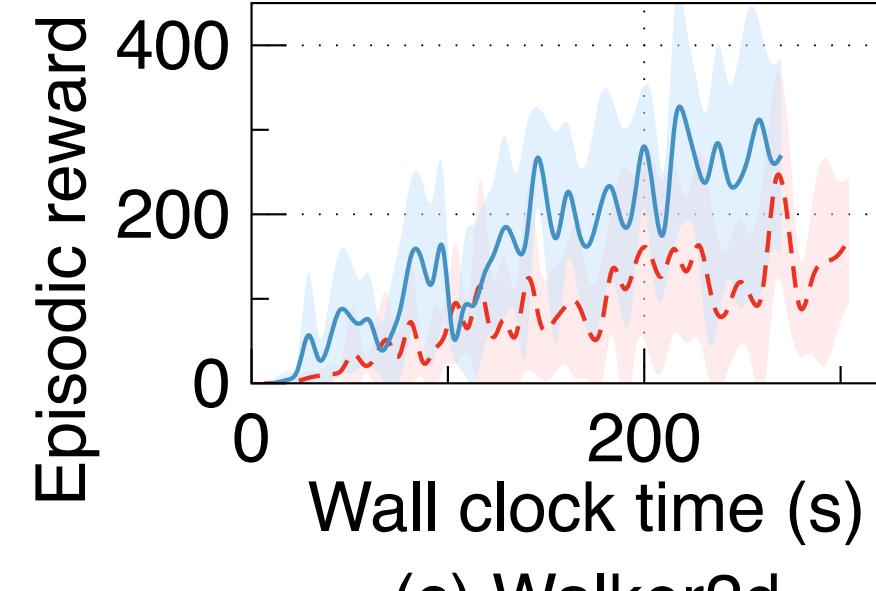
(d) Gravitar



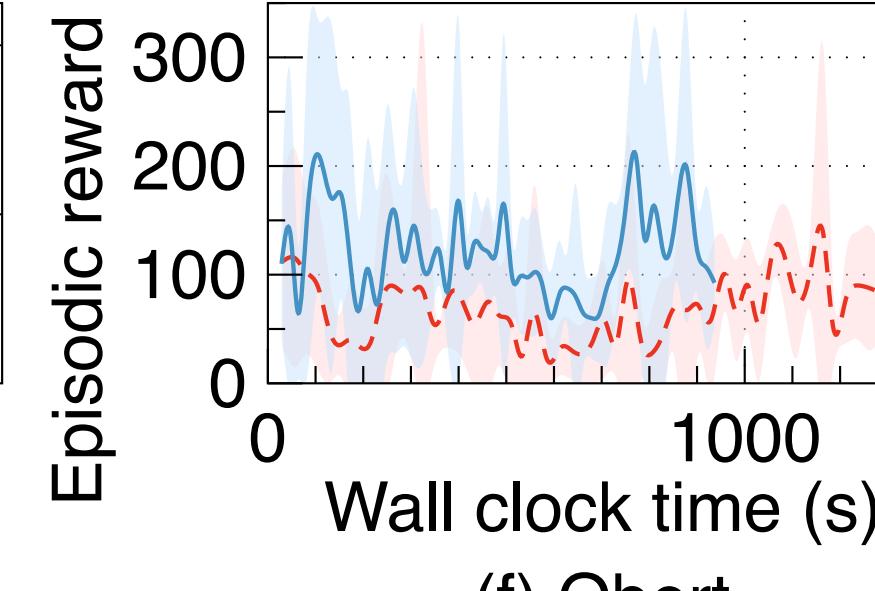
(b) Humanoid



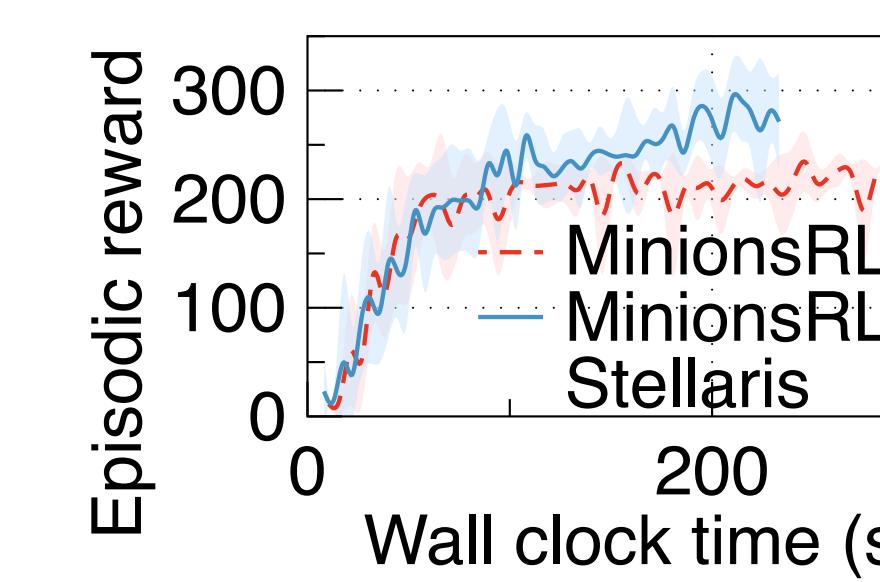
(e) SpacelInvaders



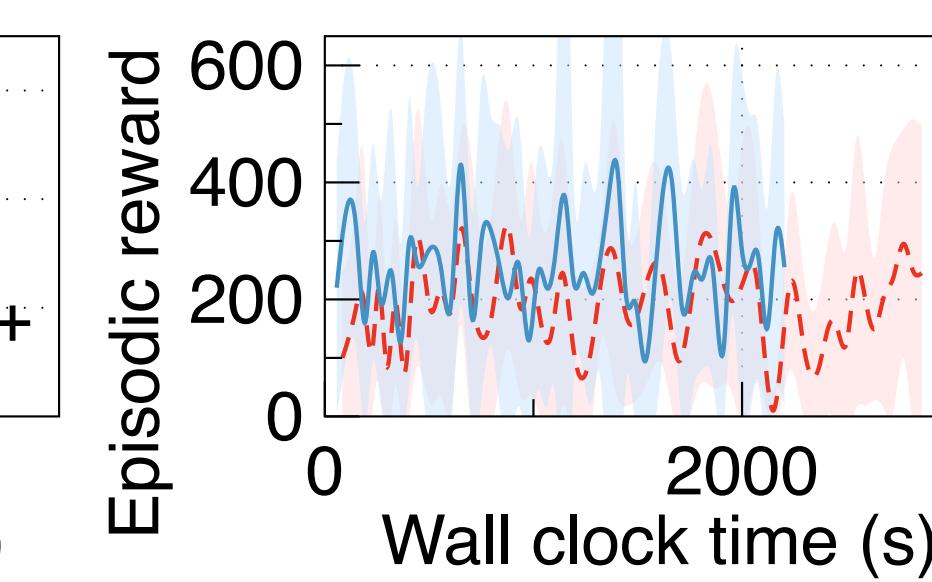
(c) Walker2d



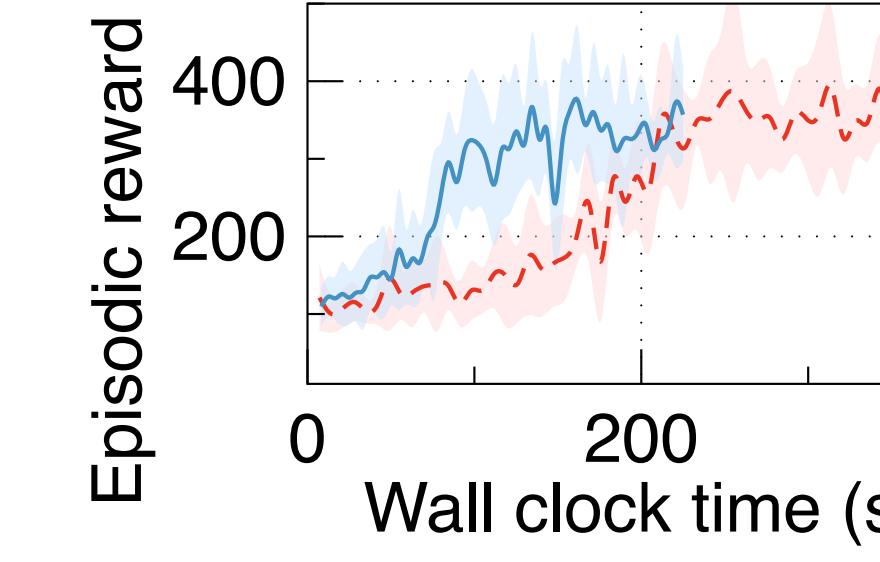
(f) Qbert



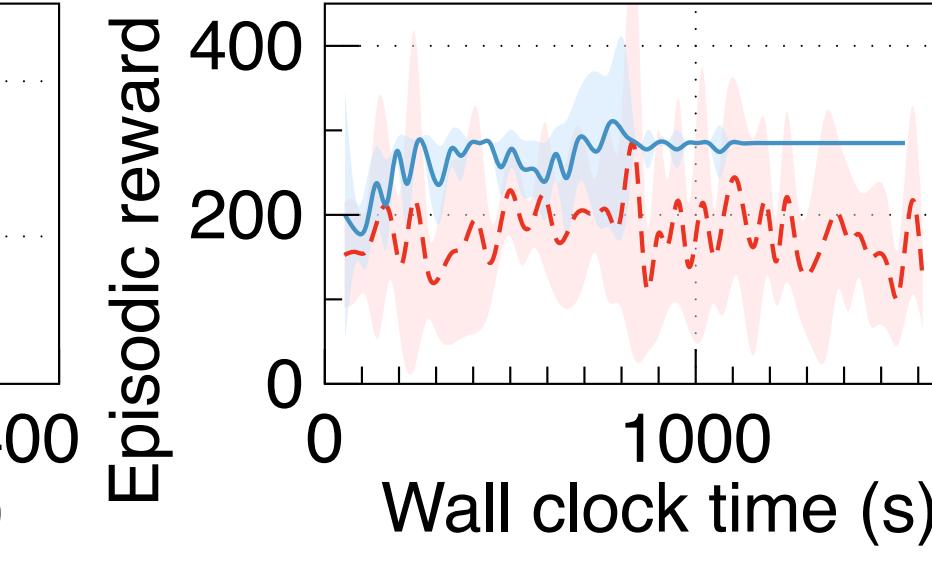
(a) Hopper



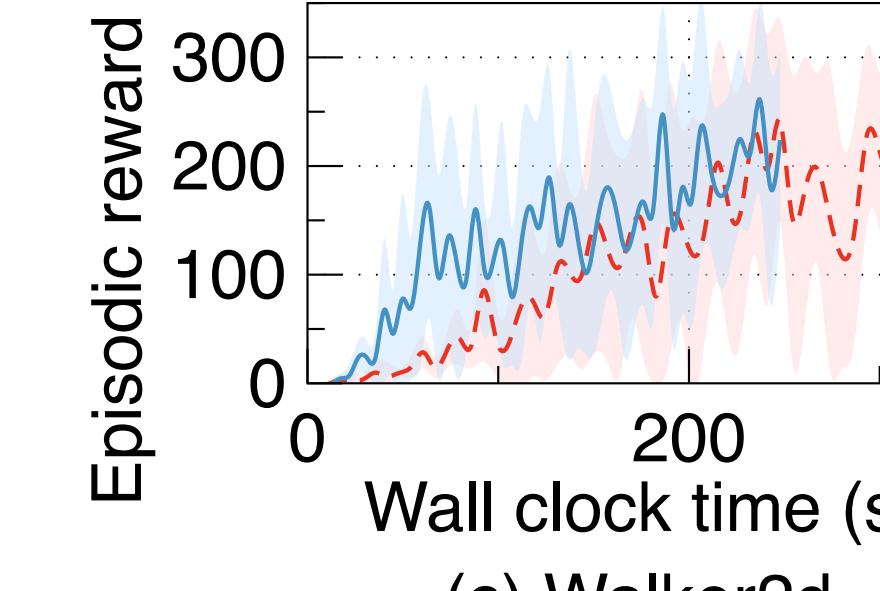
(d) Gravitar



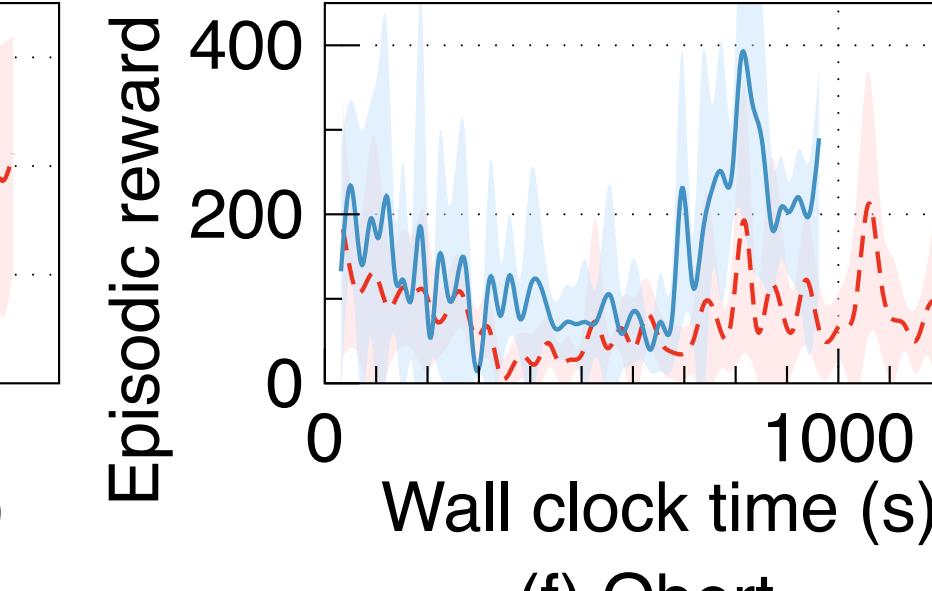
(b) Humanoid



(e) SpacelInvaders

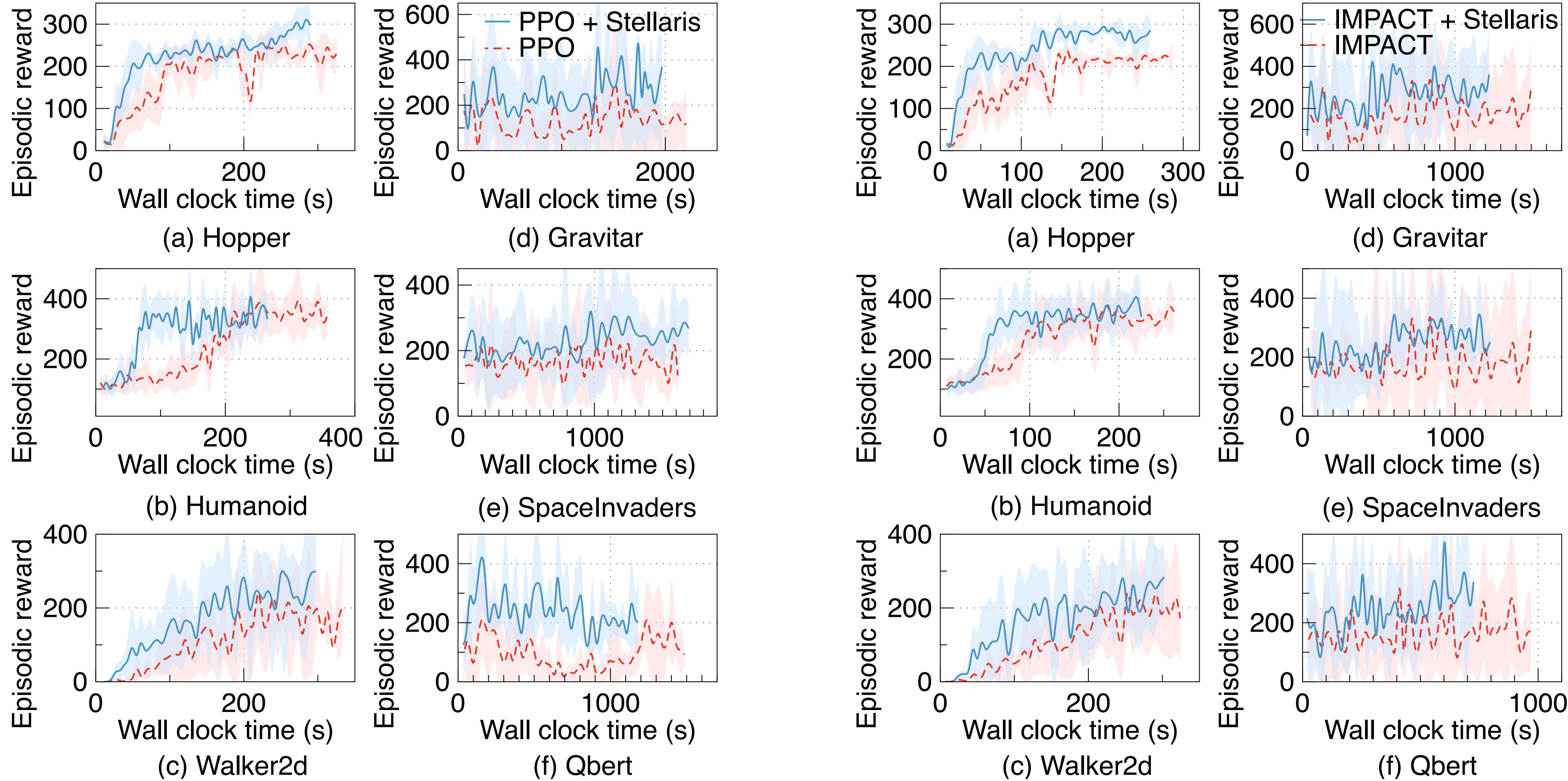


(c) Walker2d

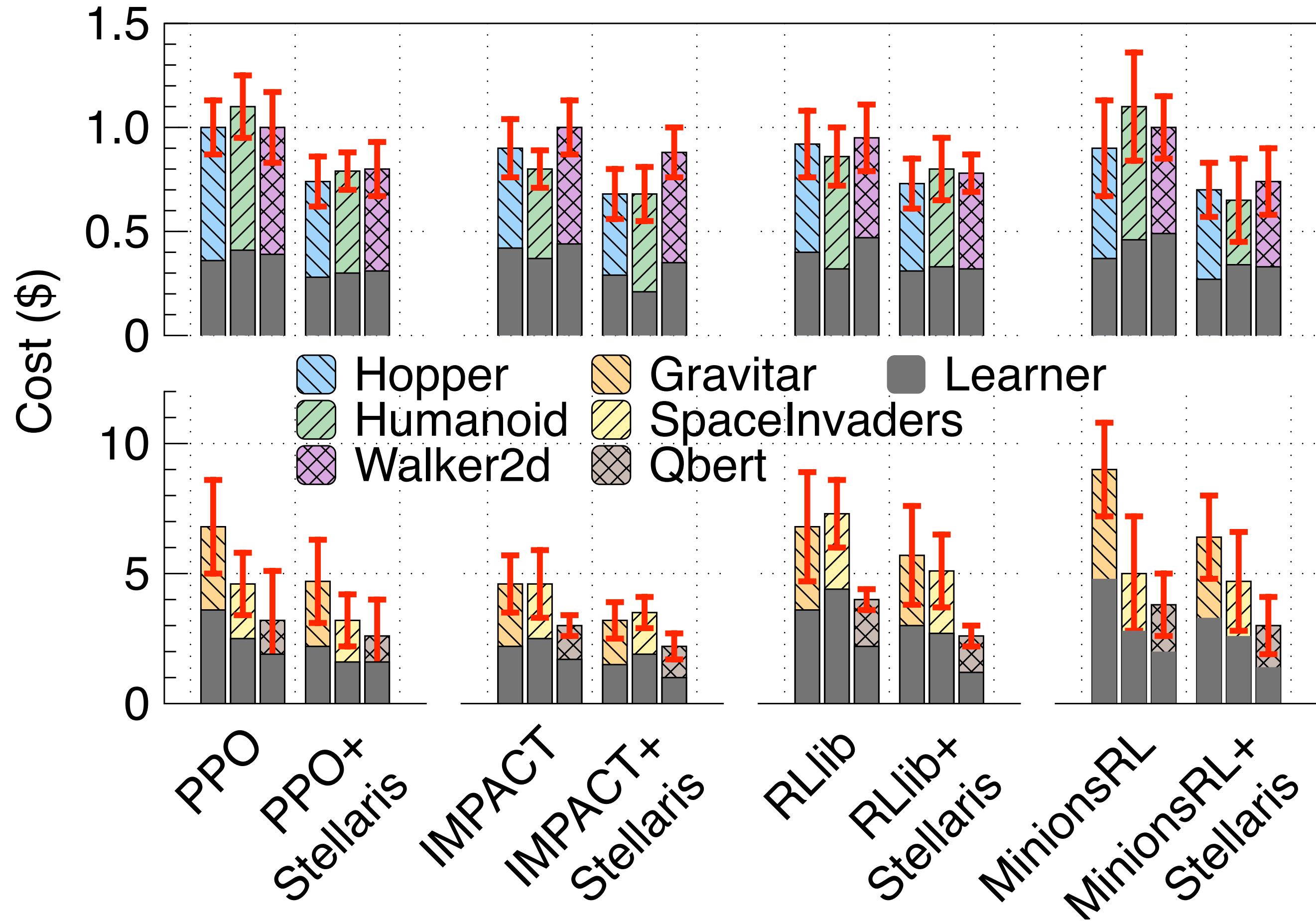


(f) Qbert

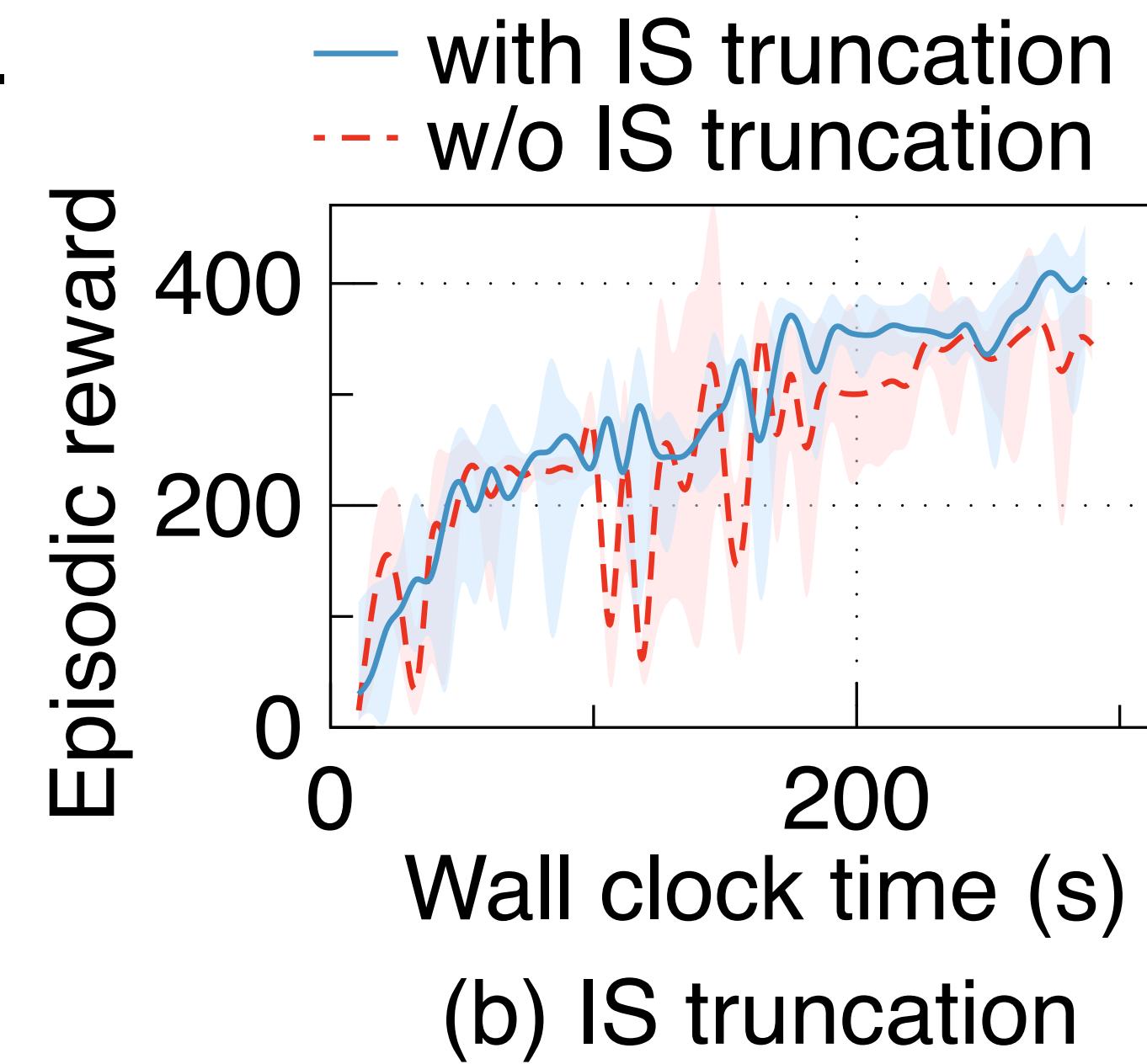
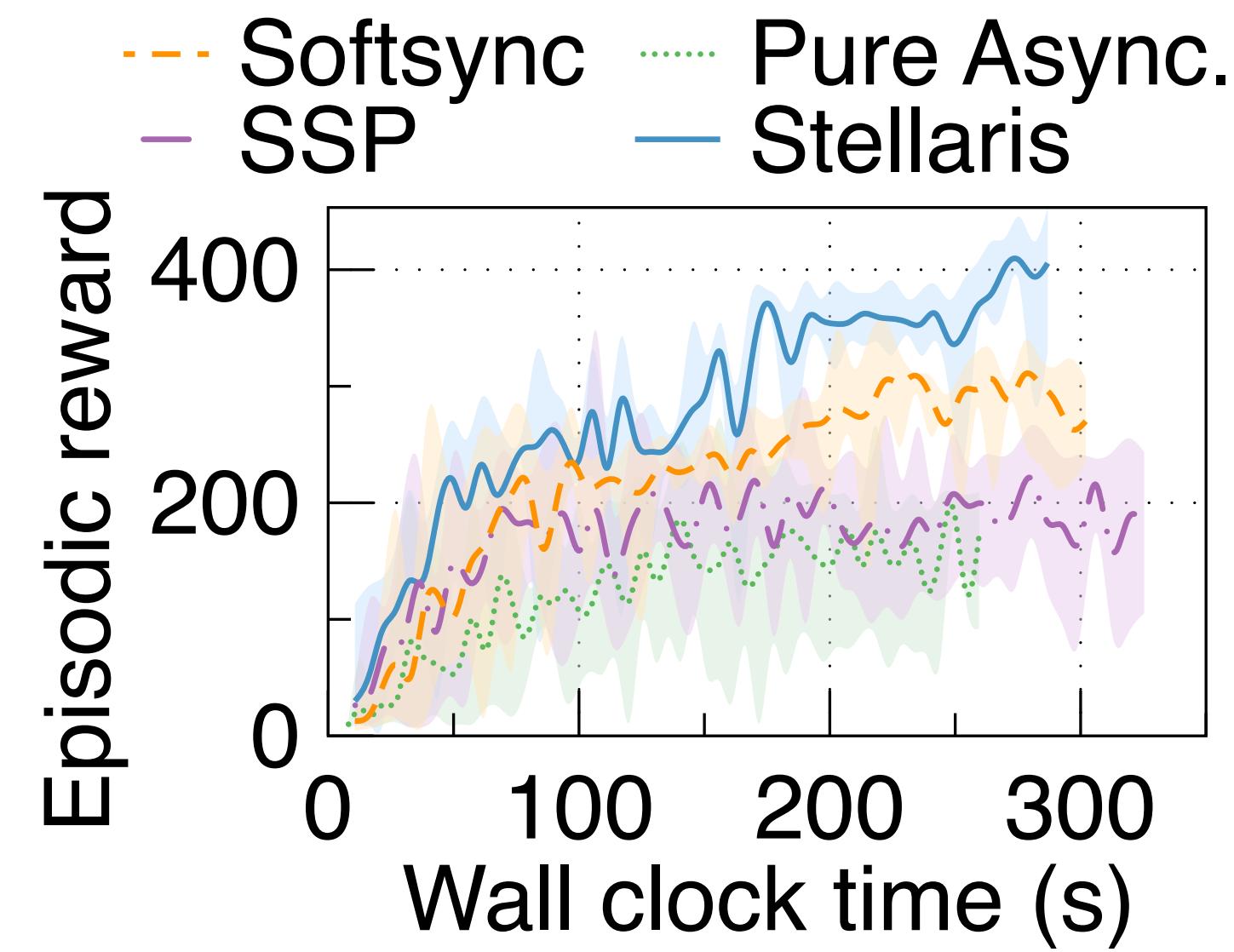
Training Performance



Training Cost



Ablation Study



Asynchronous
Serverless Learners

Global Importance
Sampling Truncation

Staleness-Aware
Gradient Aggregation

Stellaris

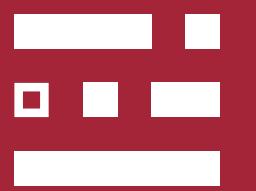
2.2X

Training performance improvement

41%

Training cost reduction



 IntelliSys Lab

THANK YOU

Stevens Institute of Technology
1 Castle Point Terrace, Hoboken, NJ 07030