

SC '24, Nov 17–22, 2024, Atlanta, USA

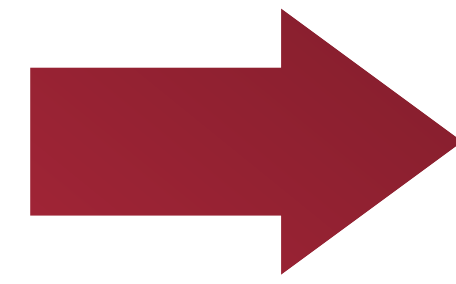


Stellaris: Staleness-Aware Distributed Reinforcement Learning with Serverless Computing

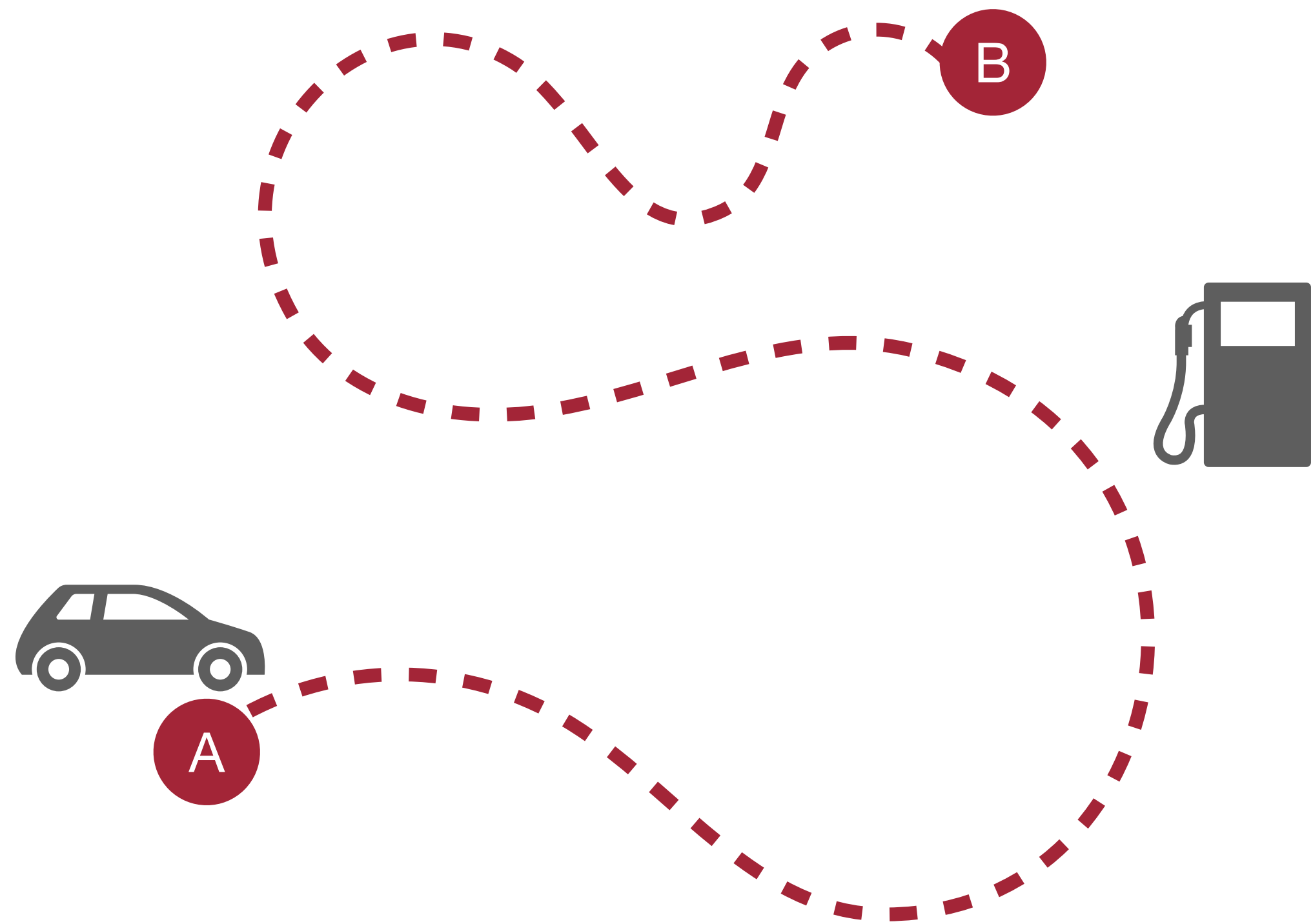
Hanfei Yu¹, Hao Wang¹, Devesh Tiwari², Jian Li³, Seung-Jong Park⁴

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

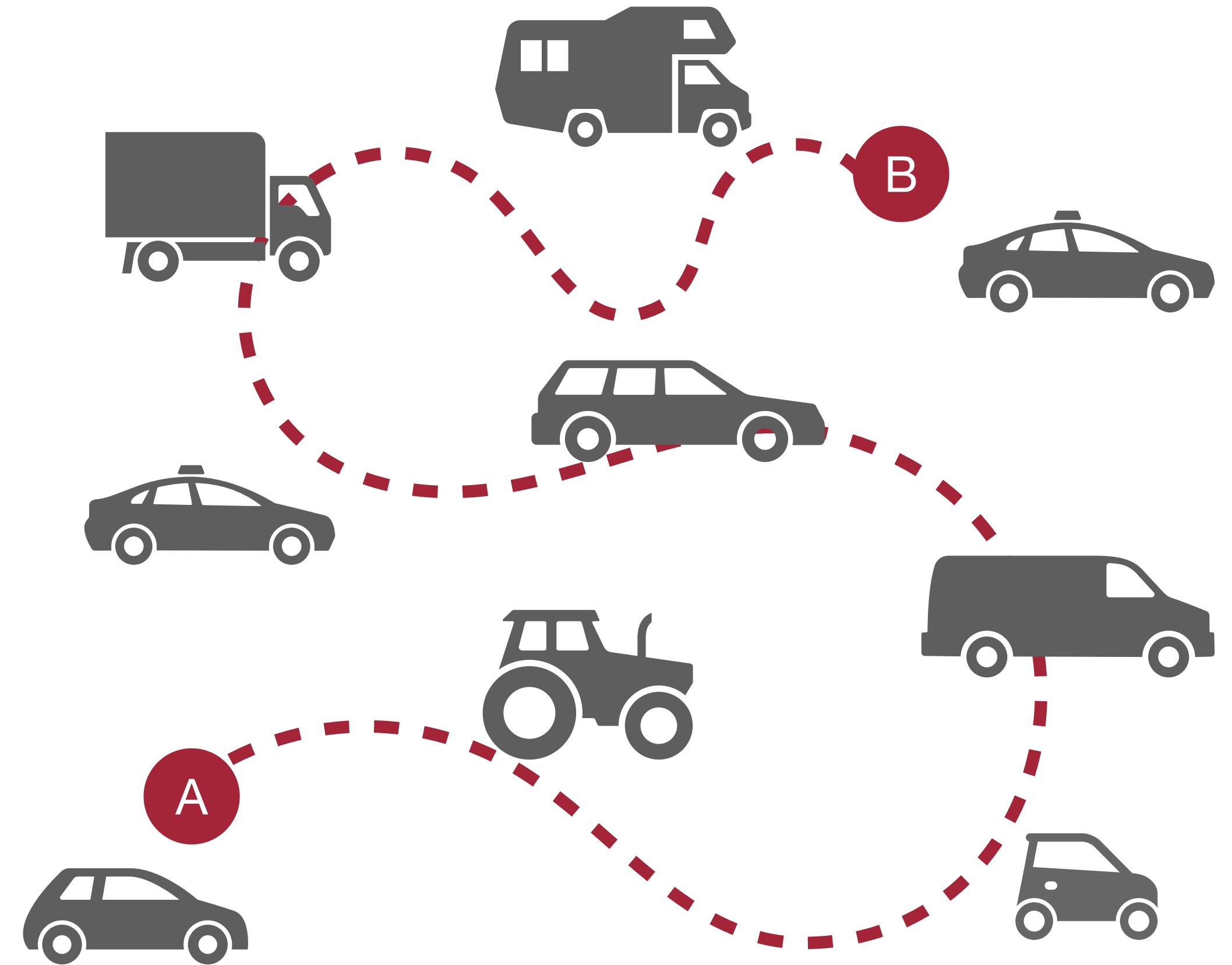
Cloud / HPC



Serverless

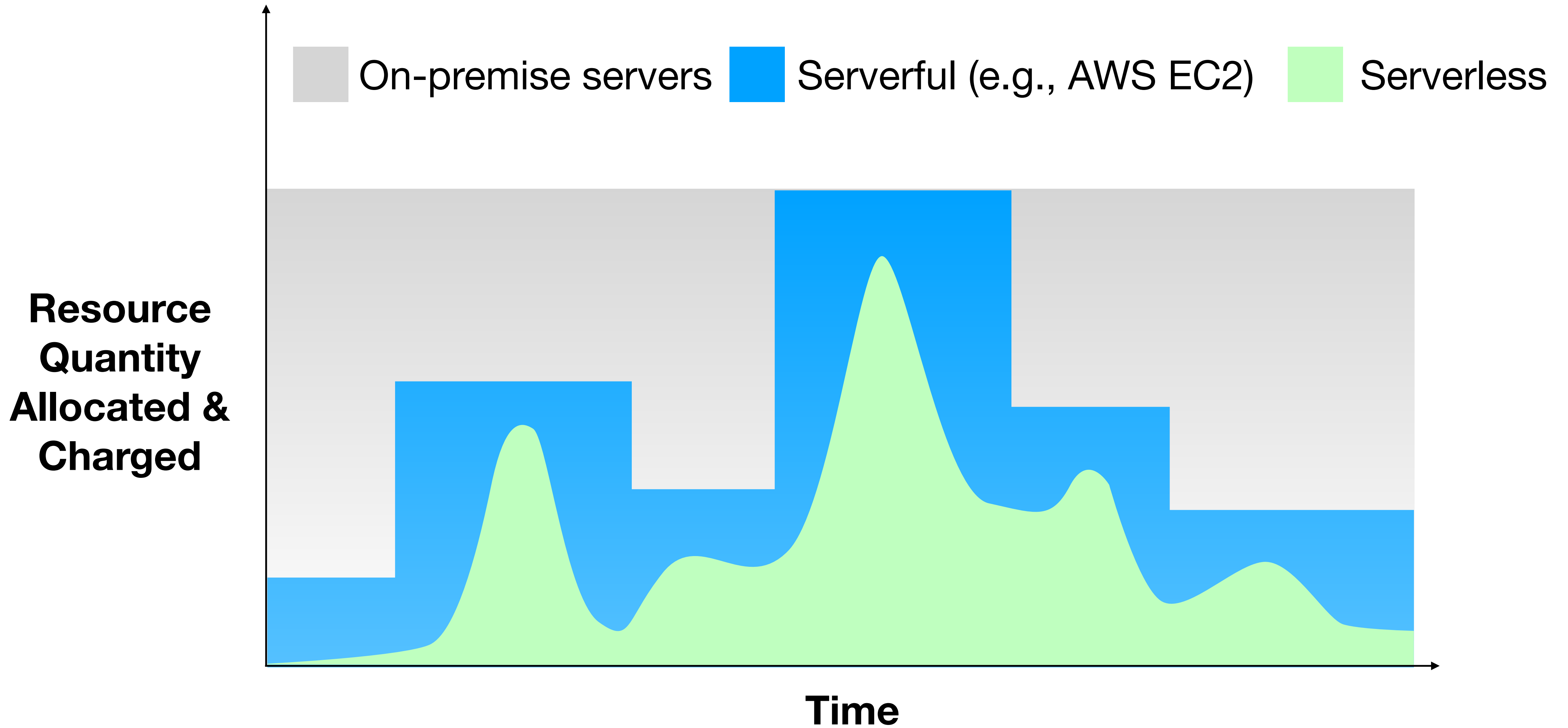


Car rental



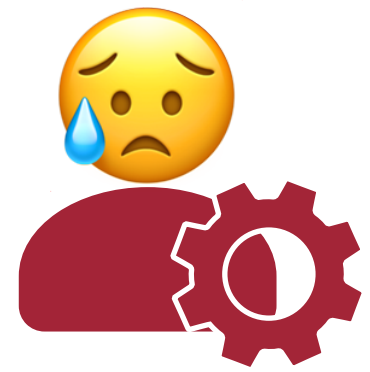
Uber/Lyft

From *A Serverless Vision for Cloud Computing*
by Prof. Ana Klimovic



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.

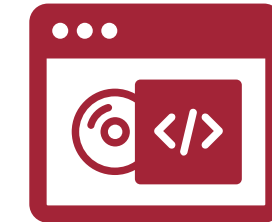
Cloud / HPC



Provider



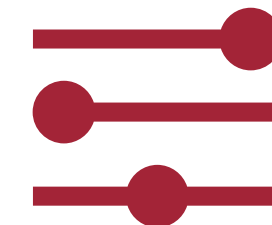
Server / VM hosting



Software Installation



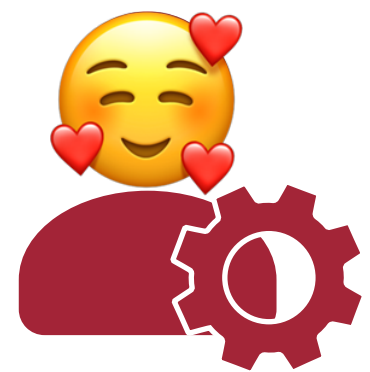
Workload management



Resource configuration



User



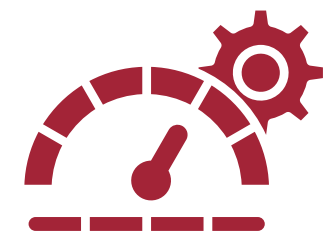
Provider



Software standardization



Workload scheduling



Resource optimization



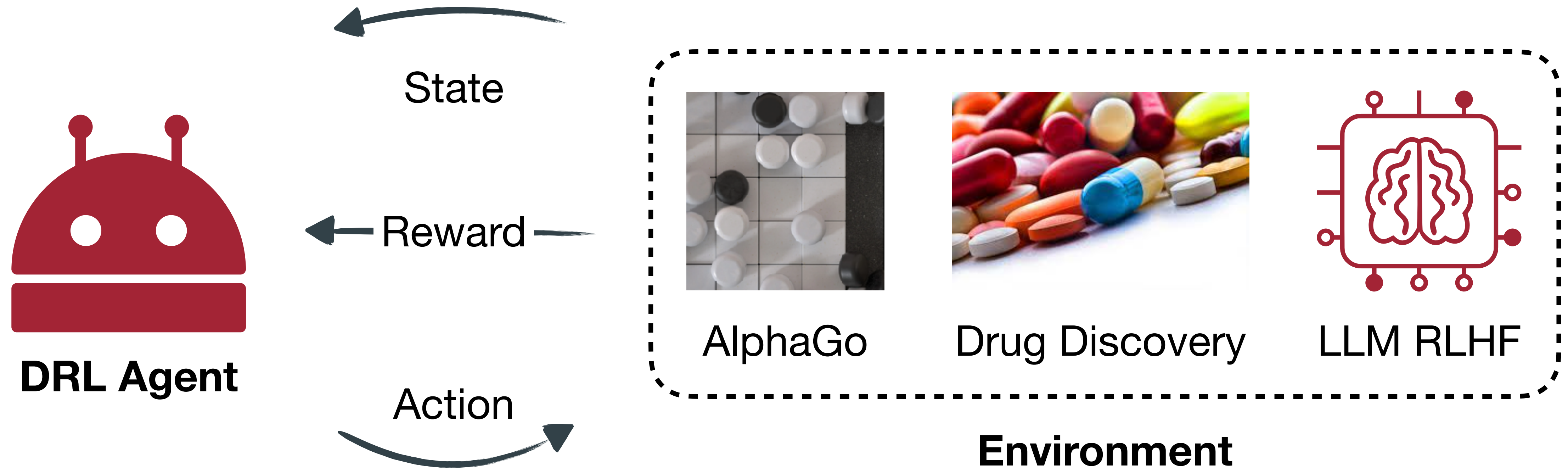
Function invocation



User

Serverless

Deep Reinforcement Learning (DRL)



DRL in HPC Community

RLScheduler: An Automated HPC Batch Job Scheduler Using Reinforcement Learning

Di Zhang¹, Dong Dai¹, Youbiao He², Forrest Sheng Bao², and Bing Xie³

¹Computer Science Department, University of North Carolina at Charlotte, {dzhang16, ddai}@uncc.edu

²Computer Science Department, Iowa State University, {yh54, fsb}@iastate.edu

³Oak Ridge Leadership Computing Facility, Oak Ridge National Laboratory. xieb@ornl.gov

Facilitating DRL

Reinforcement Learning Strategies for Compiler Optimization in High level Synthesis

Hafsah Shahzad*, Ahmed Sanallah**, Sanjay Arora**, Robert Munafo*, Xiteng Yao*, Ulrich Drepper**, and Martin Herbordt*

*CAAD Lab, ECE Department, Boston University

**Red Hat Inc.

Optimizing DRL

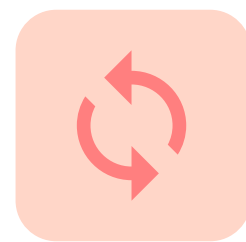
Optimizing and Extending the Functionality of EXARL for Scalable Reinforcement Learning

SAI CHENNA*, KATHERINE COSBURN*, UCHENNA EZEObI*, MAXIM MORARU*, HYUN LIM*, JULIEN LOISEAU*, JAMAL MOHD-YUSOF*, ROBERT PAVEL*, VINAY RAMAKRISHNAIAH*, ANDREW REISNER*, and KAREN TSAI*, Los Alamos National Laboratory, USA

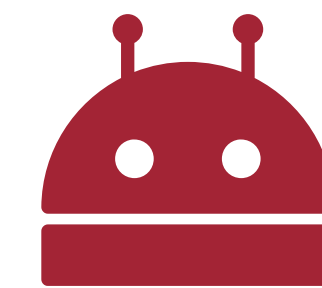
Scaling DRL Training: Actor-Learner



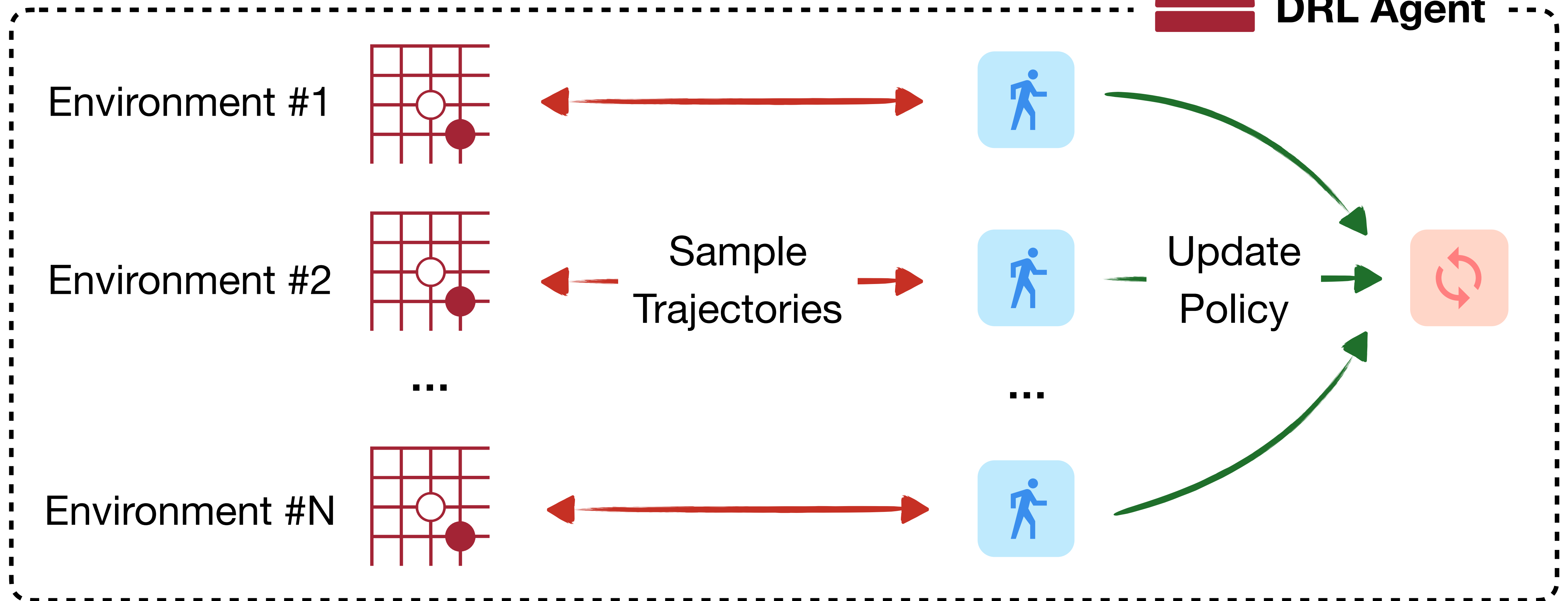
Actor



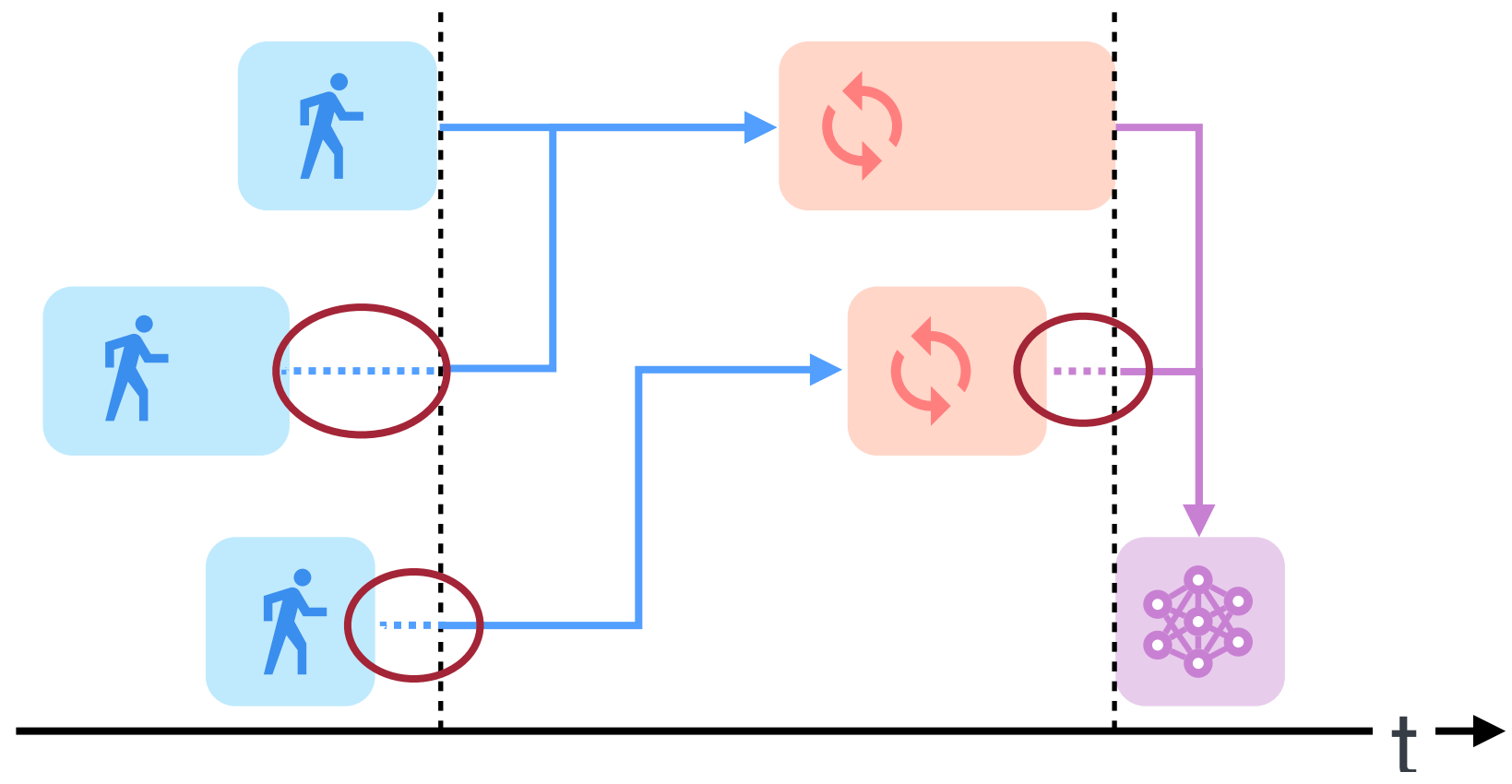
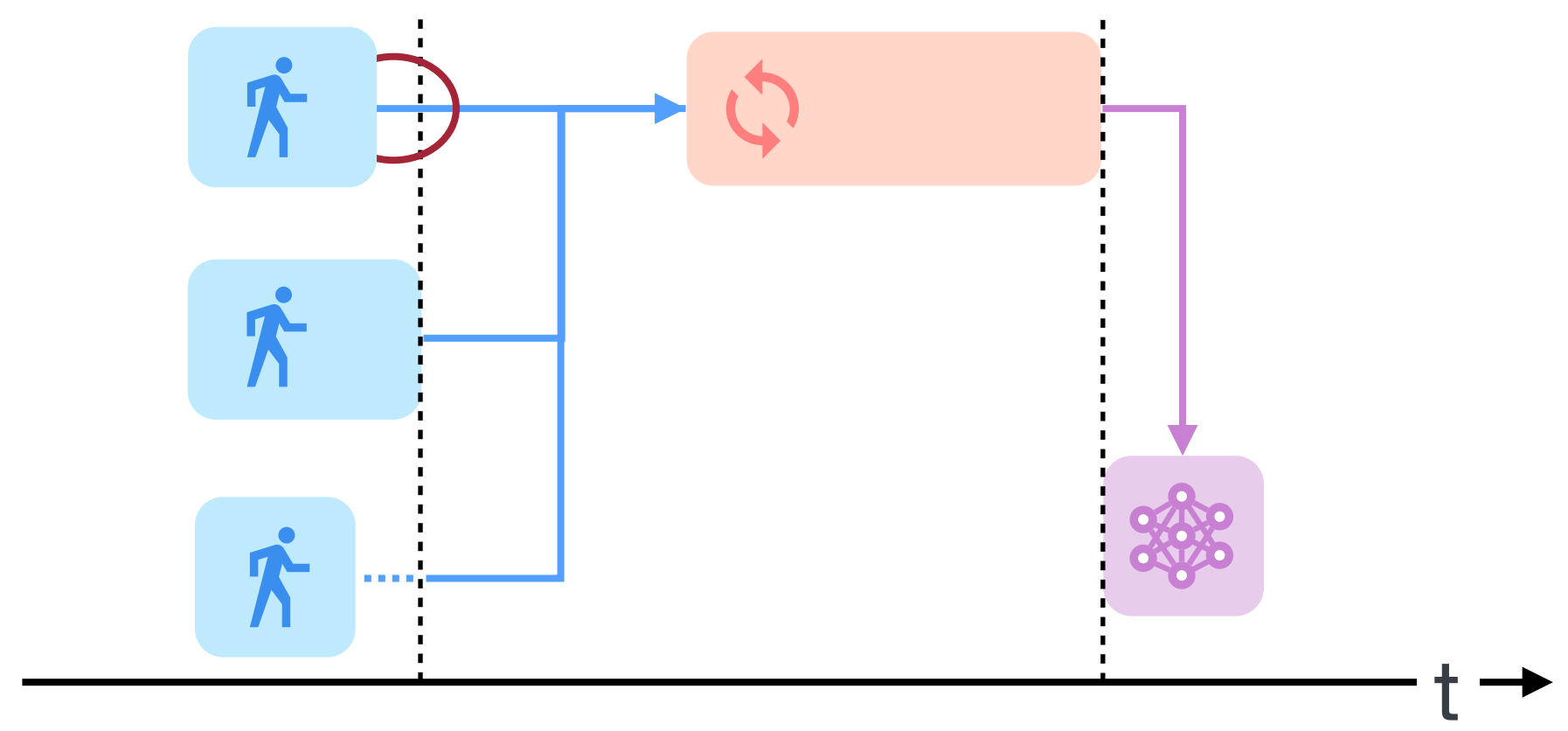
Learner



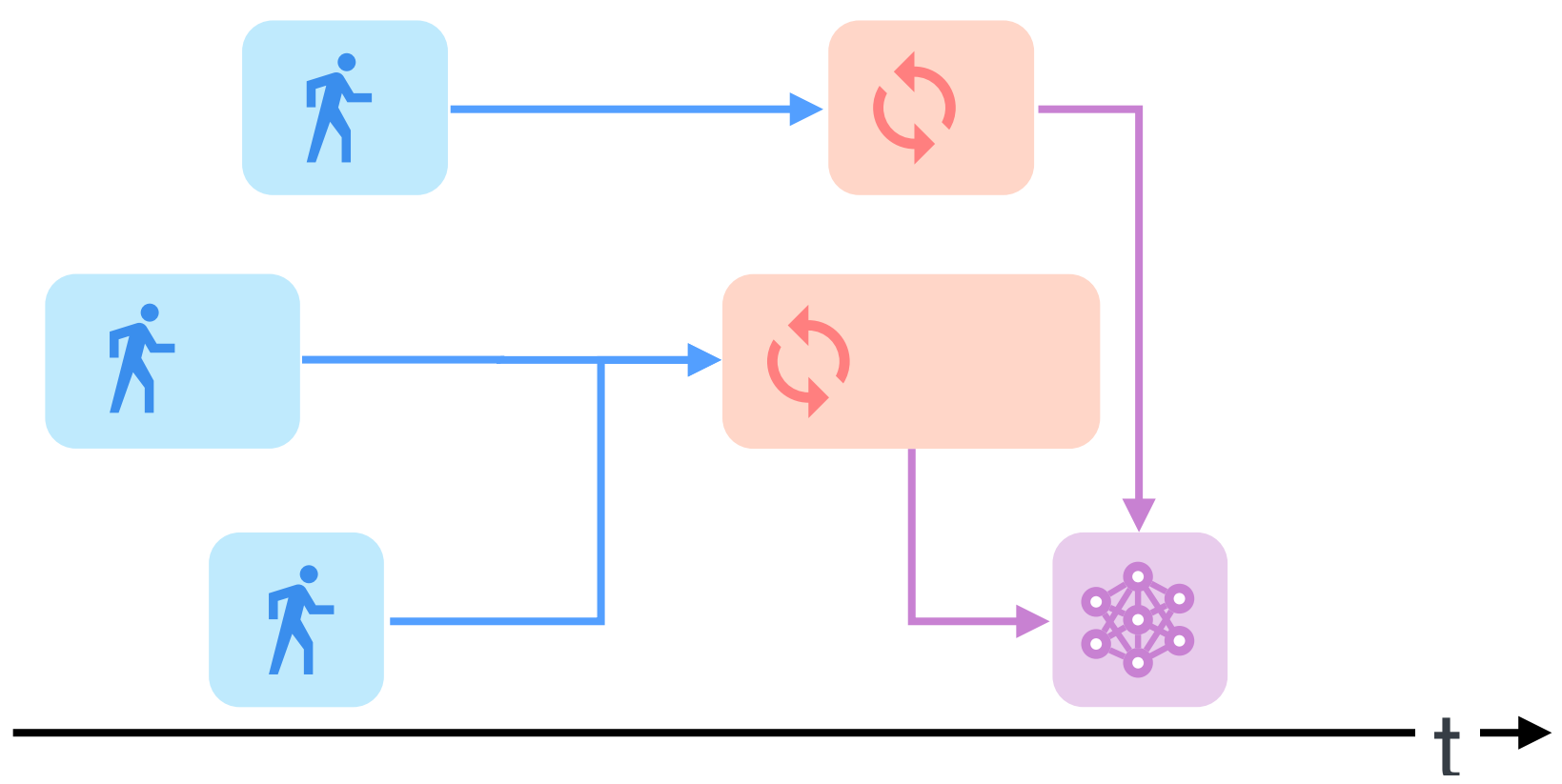
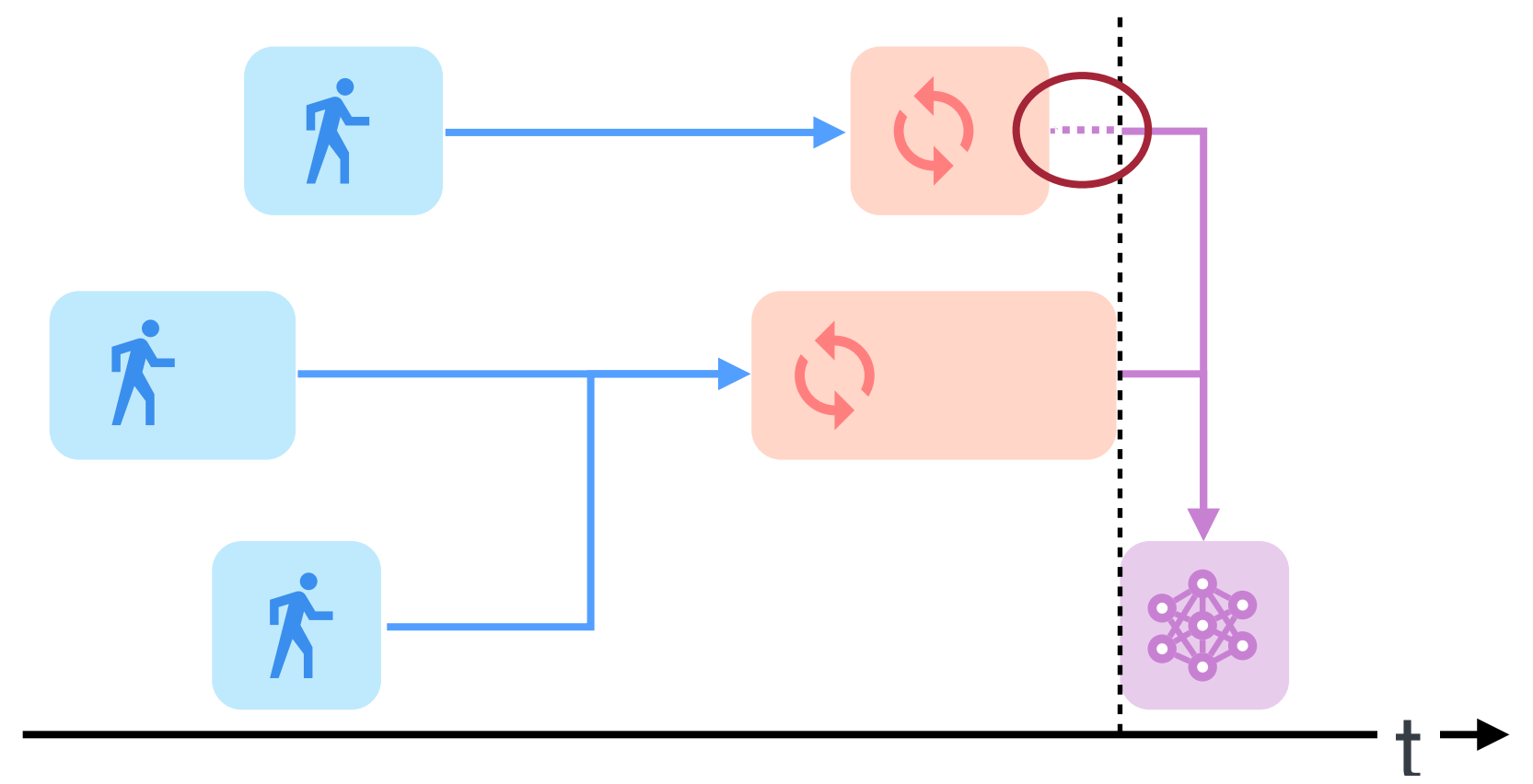
DRL Agent



Barrier  Actor  Learner  Policy update






Synchronous Actors + Centralized Sync. Learning | **Synchronous Actors + Synchronous Learners**



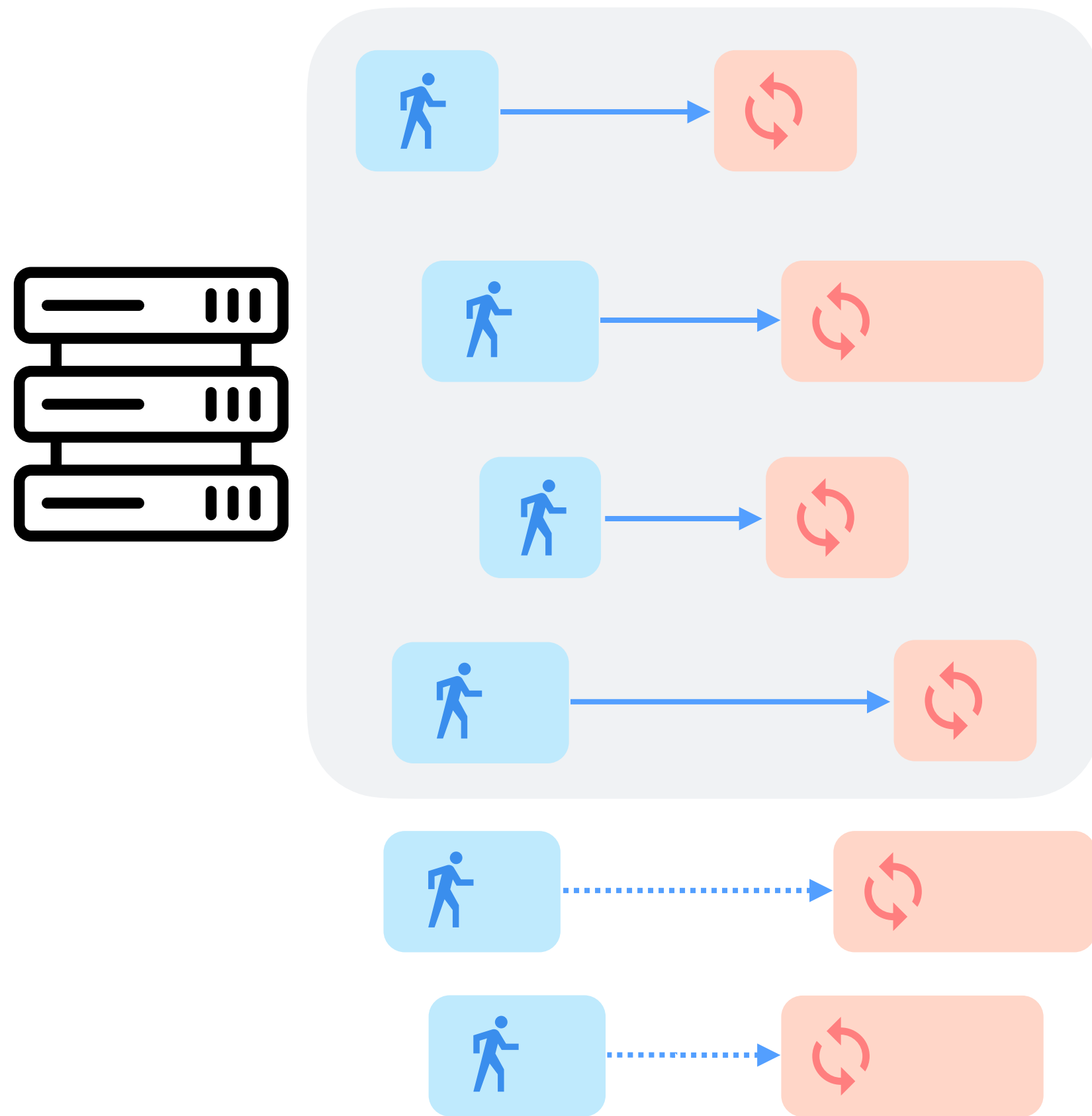
Asynchronous Actors + Synchronous Learning | **Synchronous Actors + Asynchronous Learning**



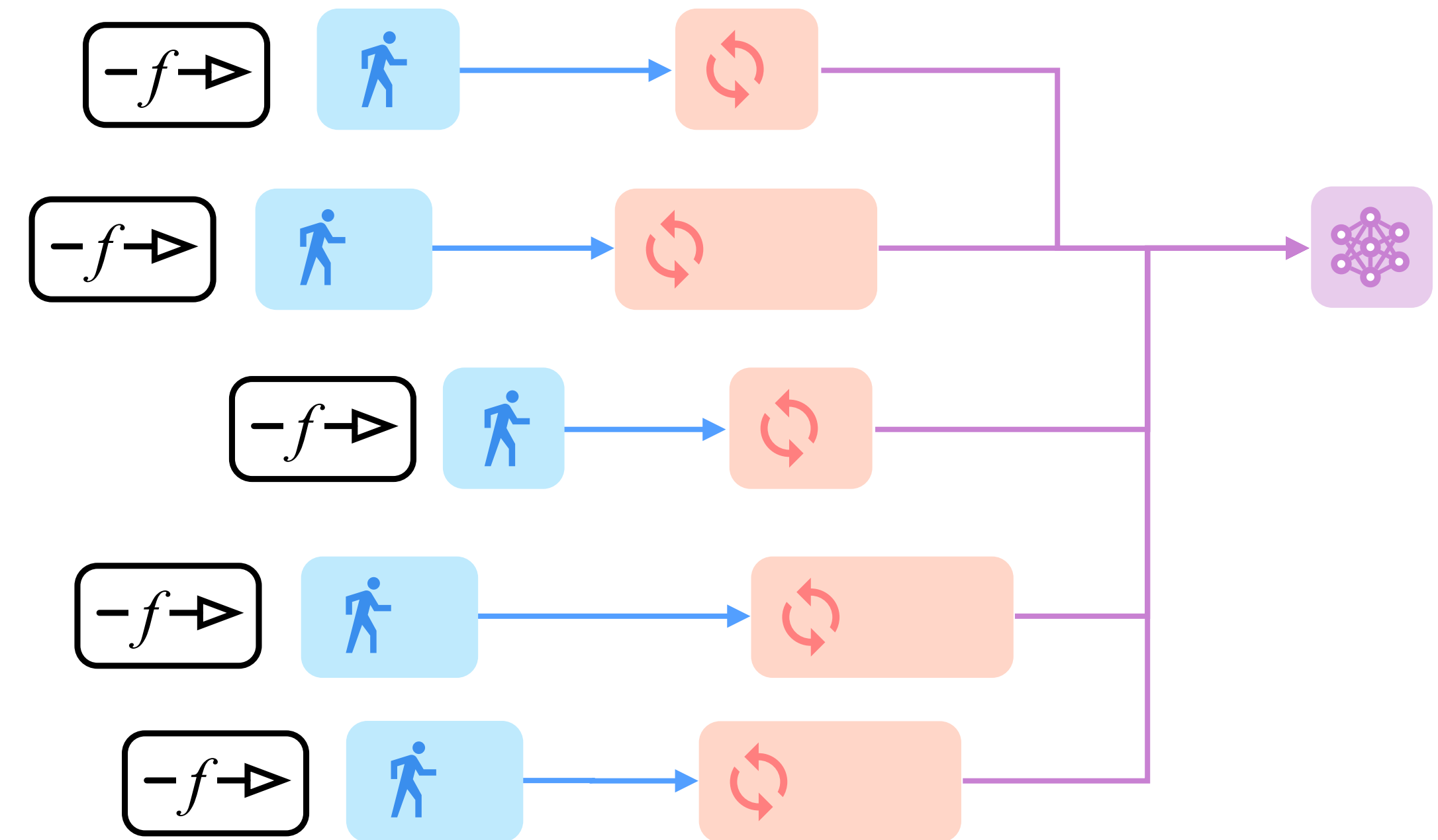
How Does Serverless Benefit DRL?

-  Actor
-  Learner
-  Policy update

 VMs or servers reaching capacity



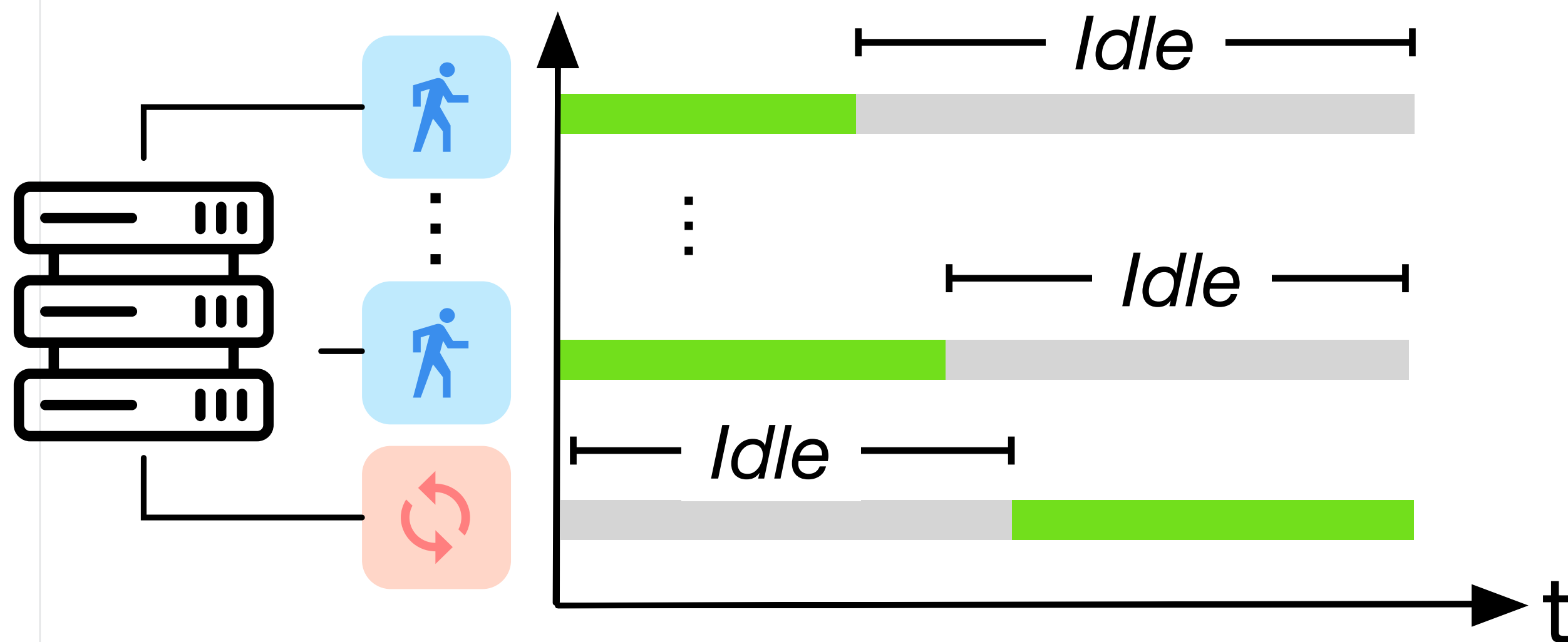
Serverful DRL



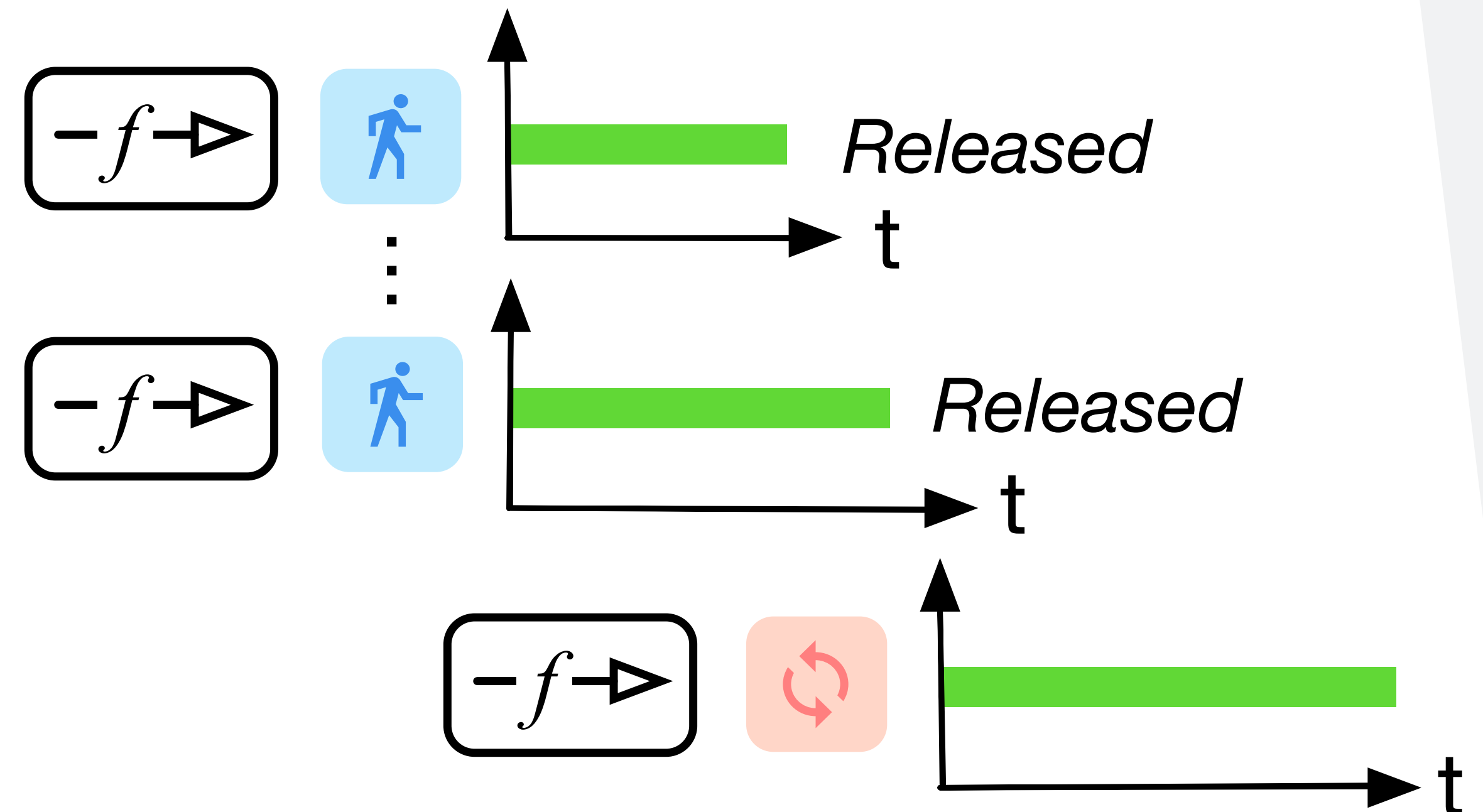
 Serverless functions instantly **scale up**

Serverless DRL

Serverful vs. Serverless DRL Training



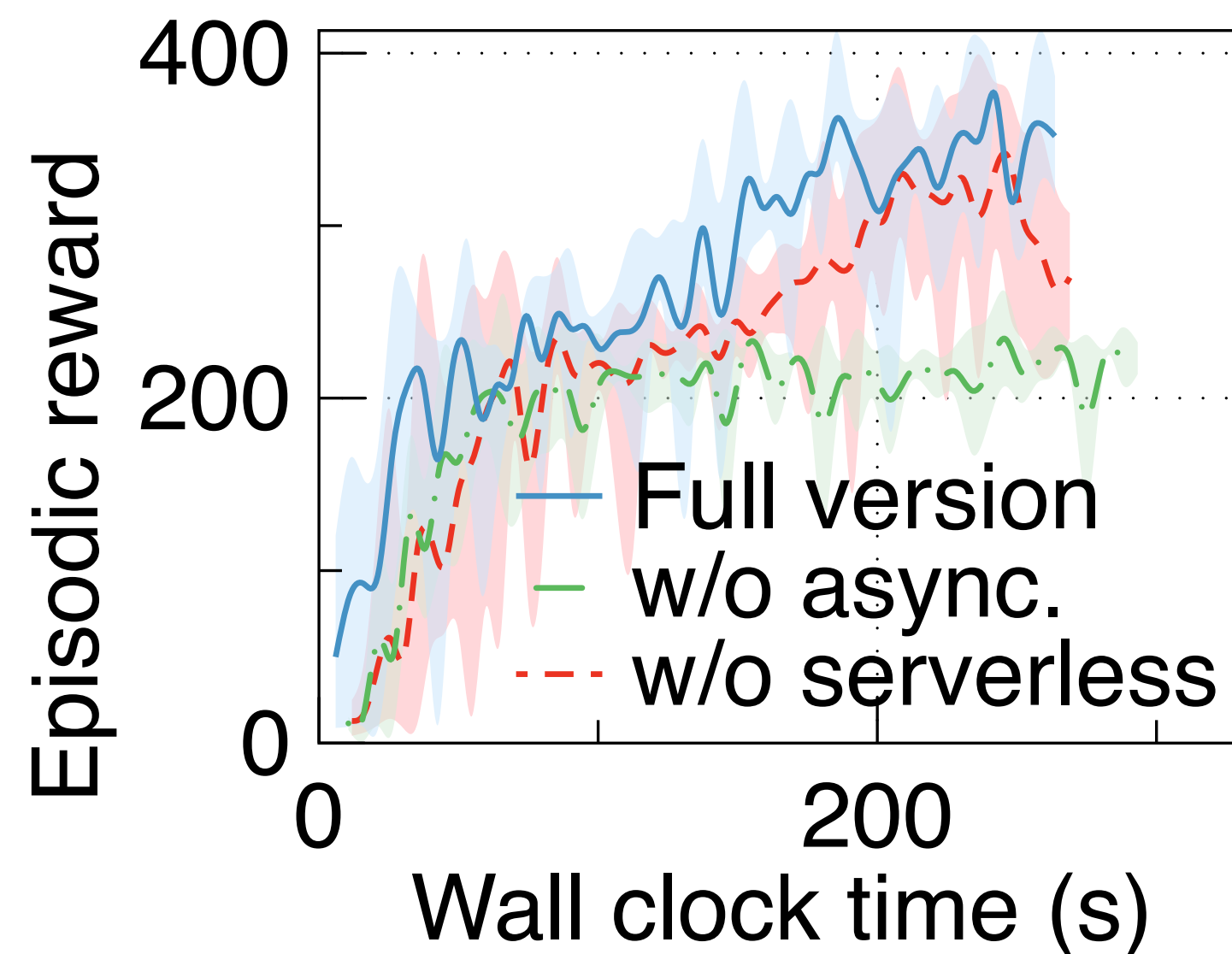
Serverful Training



Serverless Training

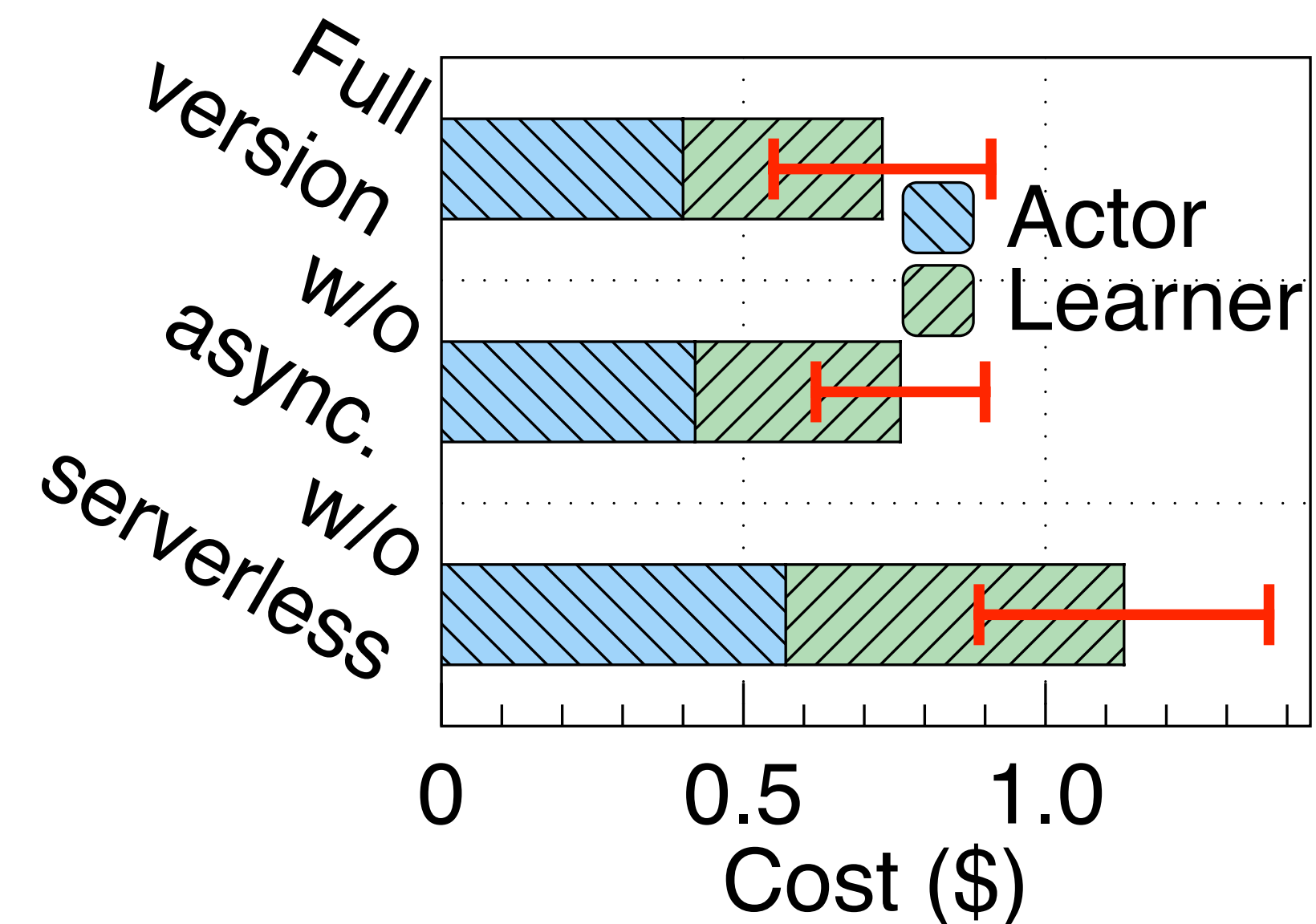
Asynchronous and Serverless DRL Training

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 2024)	✓	✓	✓	✓

Challenges

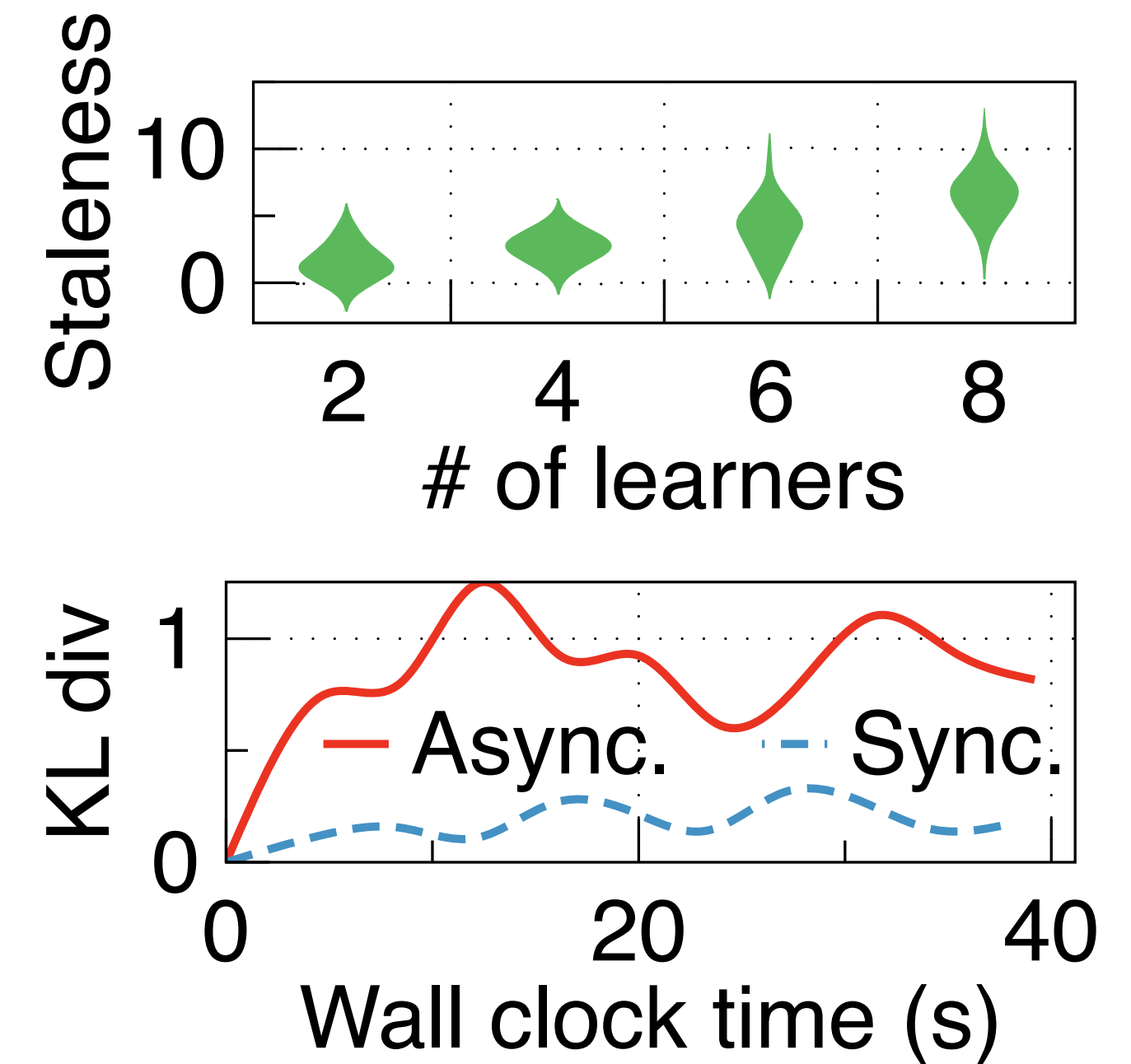
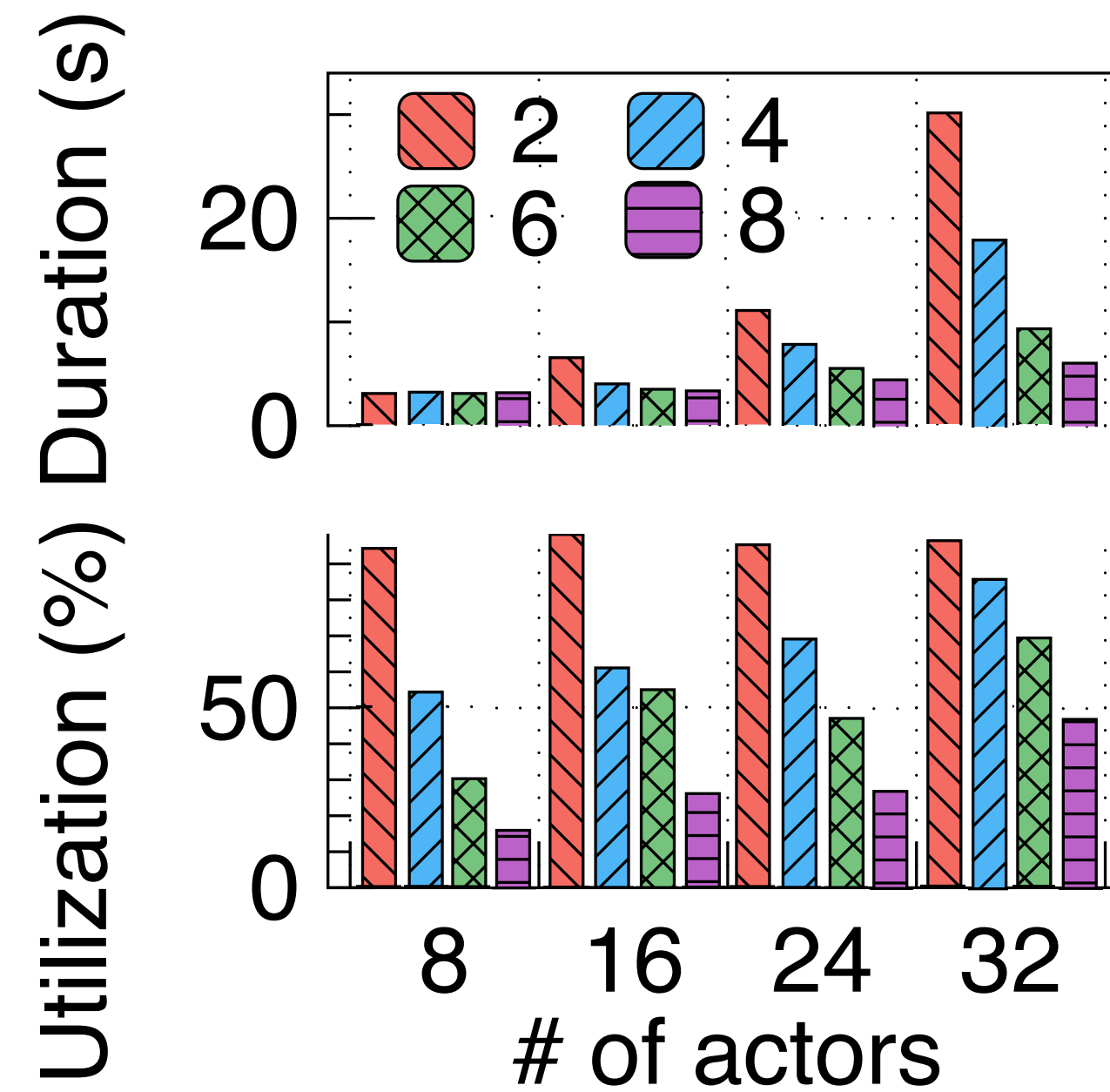
Dynamic learner orchestration and usage



Dynamic staleness



Unstable policy updates



Design Goals

Dynamic learner
orchestration



On-Demand Serverless Learners

Dynamic staleness



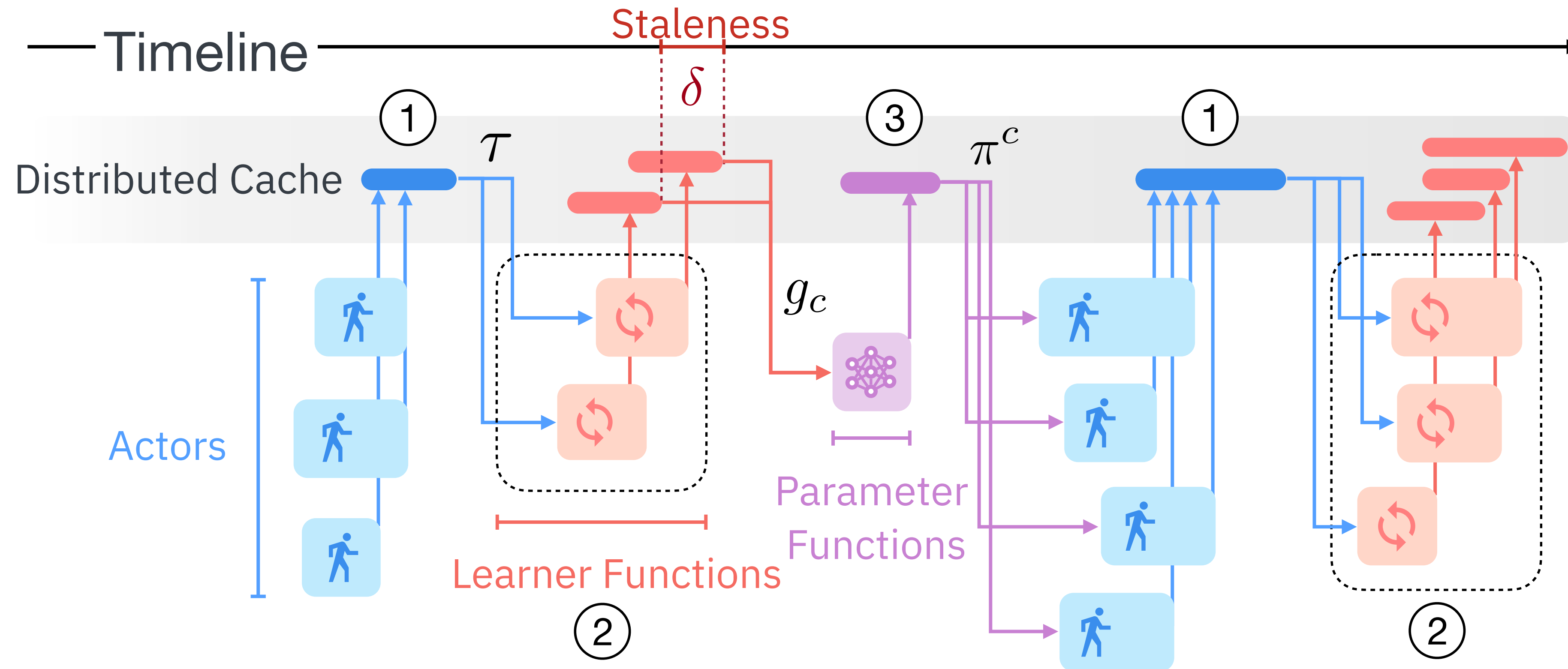
Staleness-Aware Gradient
Aggregation

Unstable policy updates



Global Importance Sampling
Truncation

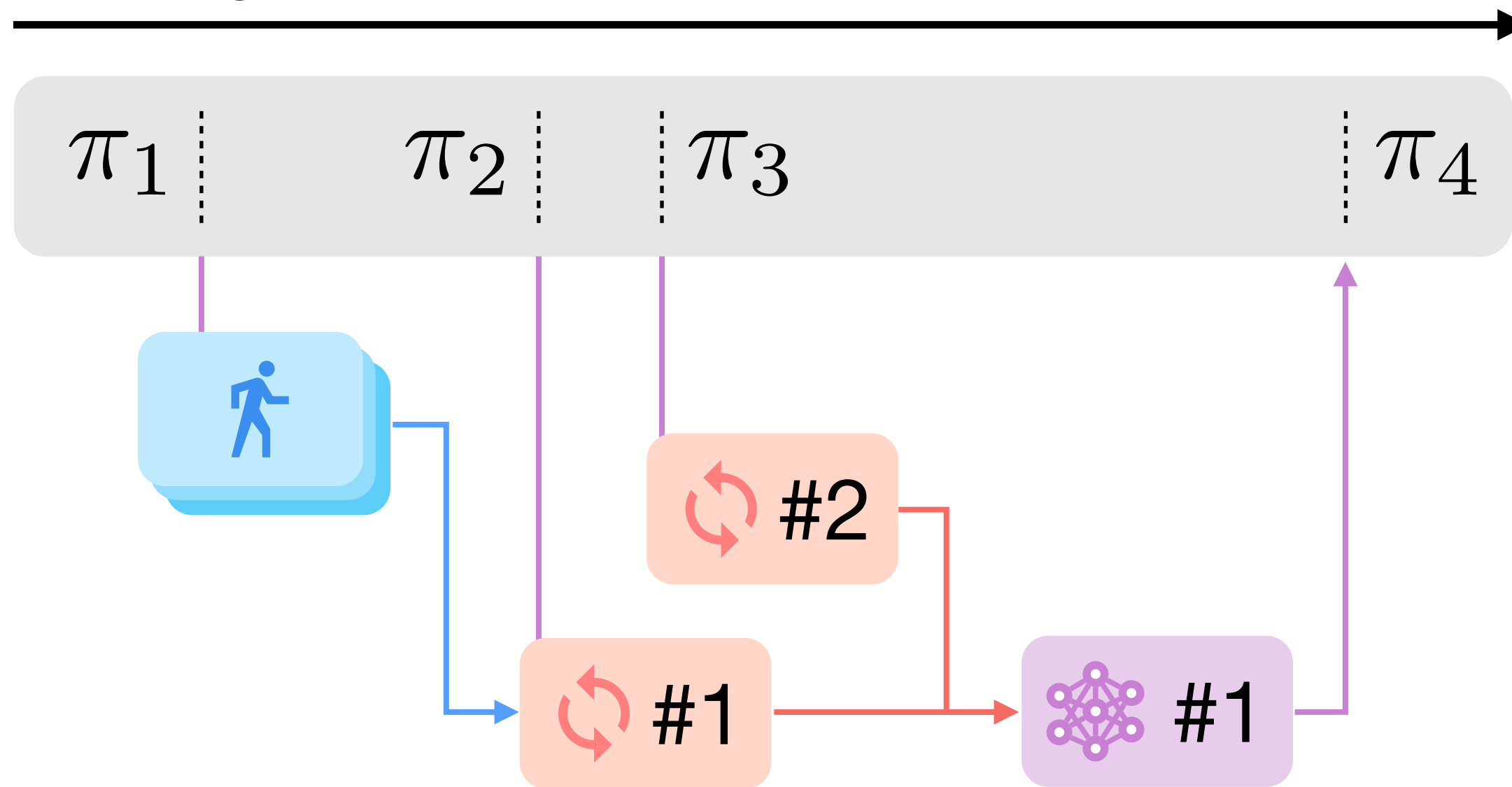
Stellaris







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

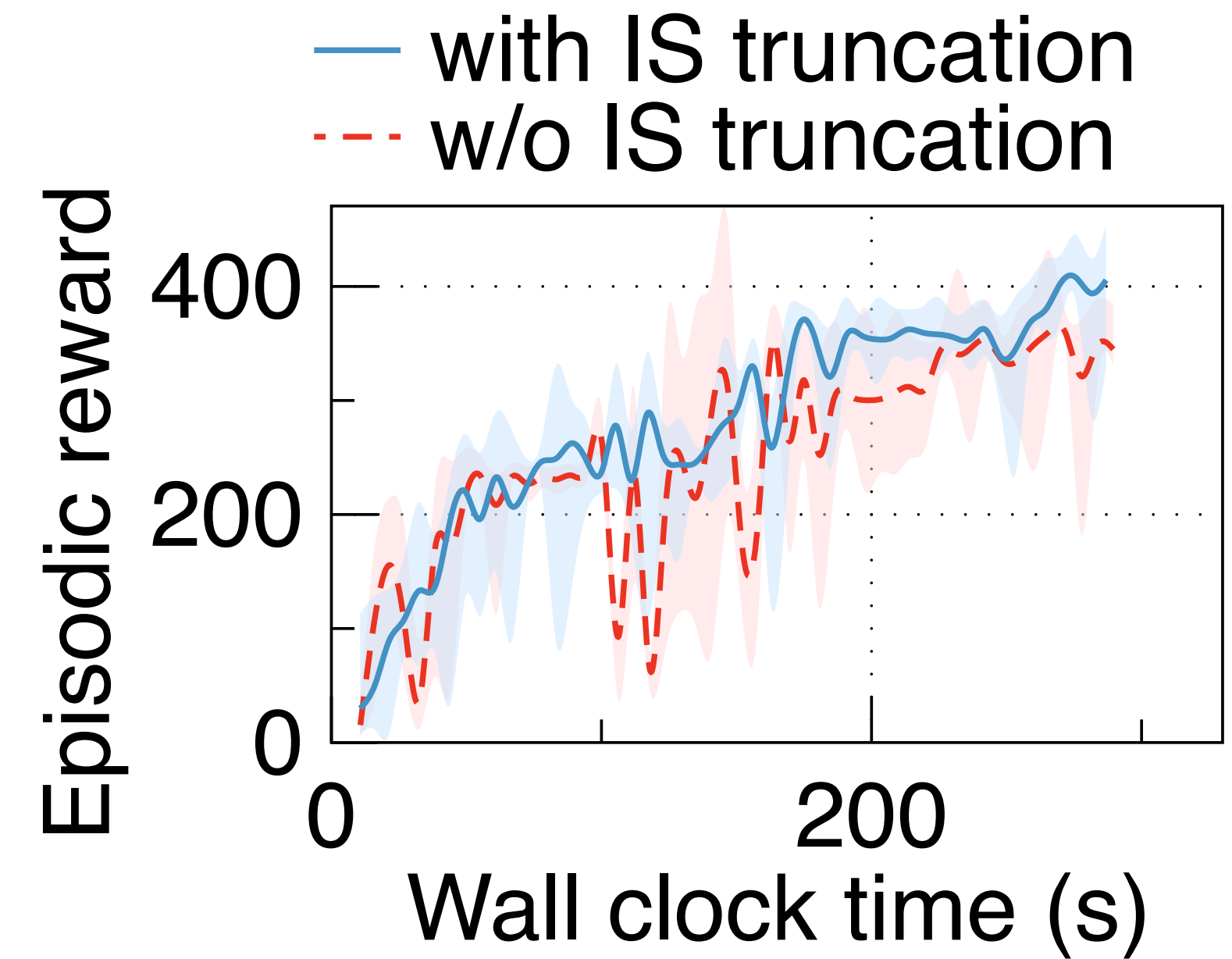
Global Importance Sampling Truncation

Policy Version Timeline



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

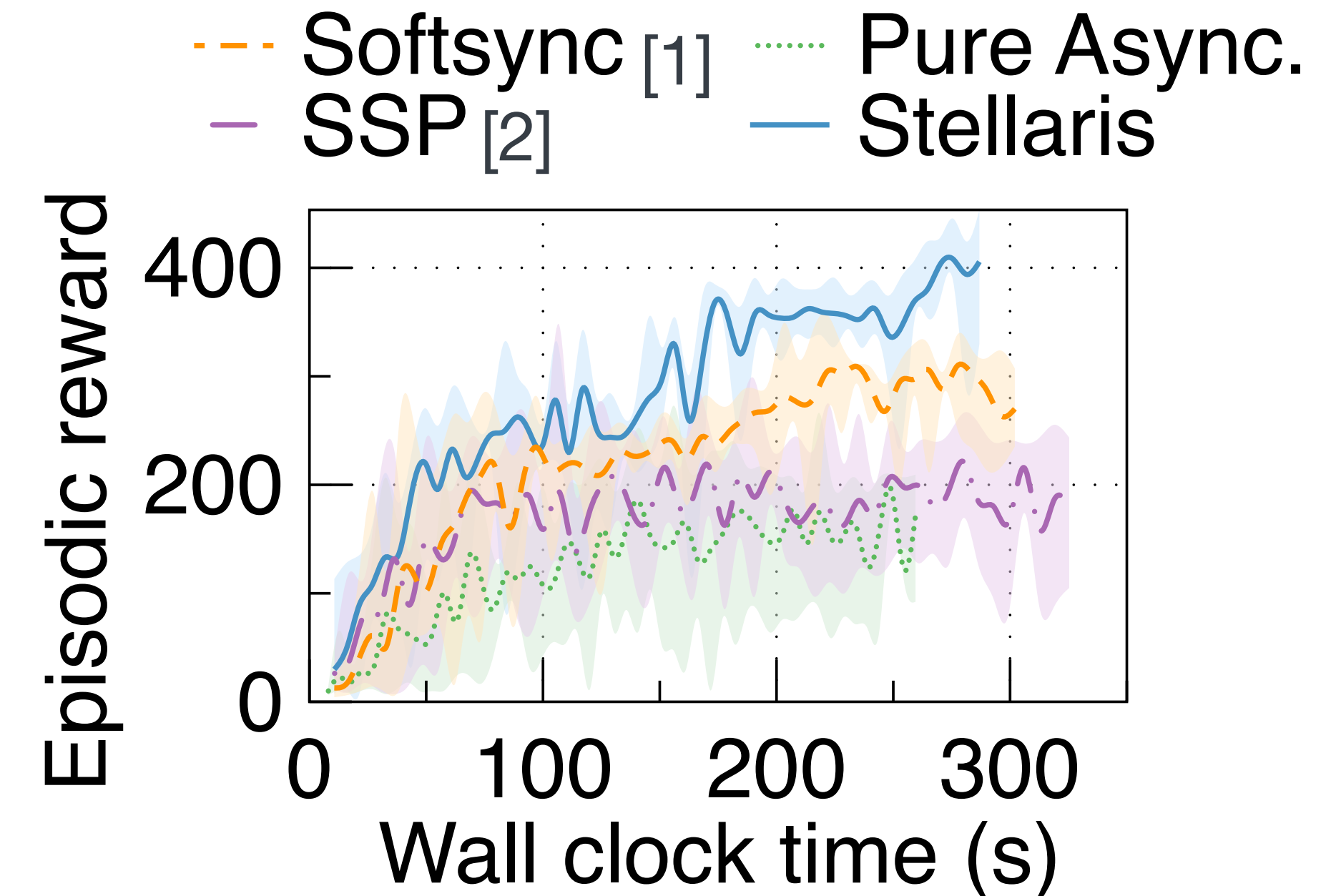
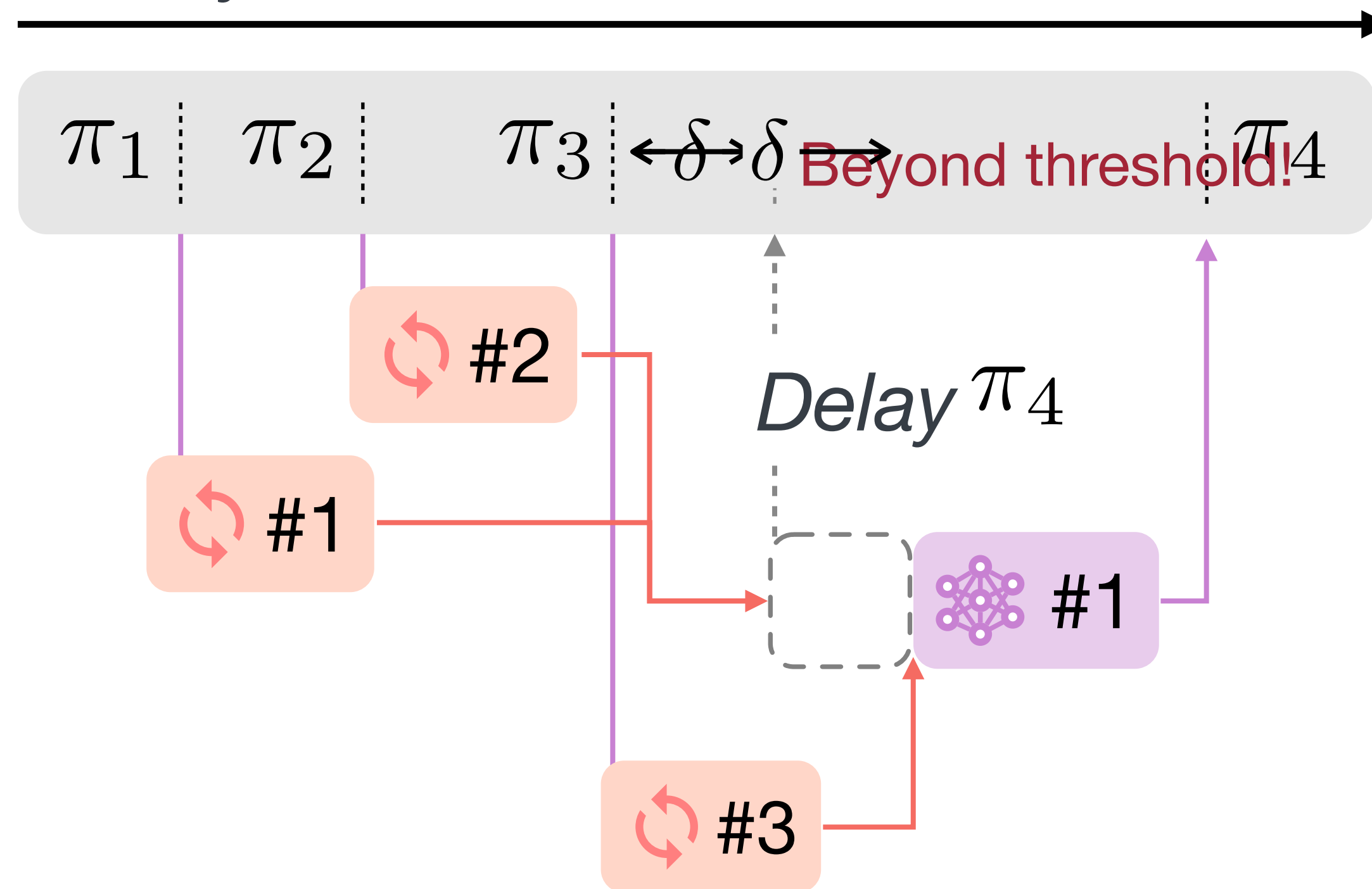
 Actor
  Learner Function
  Parameter Function
  Delay



(b) IS truncation

Staleness-Aware Gradient Aggregation

Policy Version Timeline



(a) Gradient aggregation

🚶 Actor
 ↻ Learner Function
 🧠 Parameter Function
 Delay

[1] Zhang, W. et al. "Staleness-Aware Async-SGD for Distributed Deep Learning". IJCAI. 2016.

[2] Ho, Qirong, et al. "More Effective Distributed ML via a Stale Synchronous Parallel Parameter Server." NeurIPS. 2013.

Theoretical Guarantees

Importance Sampling Truncation

Lower bound on monotonic
reward improvement [1]

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

Gradient Aggregation

Near-linear convergence rate [2]

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

[1] Joshua, A. et al. "Constrained Policy Optimization." *ICML*. 2017.

[2] Zhang, W. et al. "Staleness-Aware Async-SGD for Distributed Deep Learning". *IJCAI*. 2016.

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

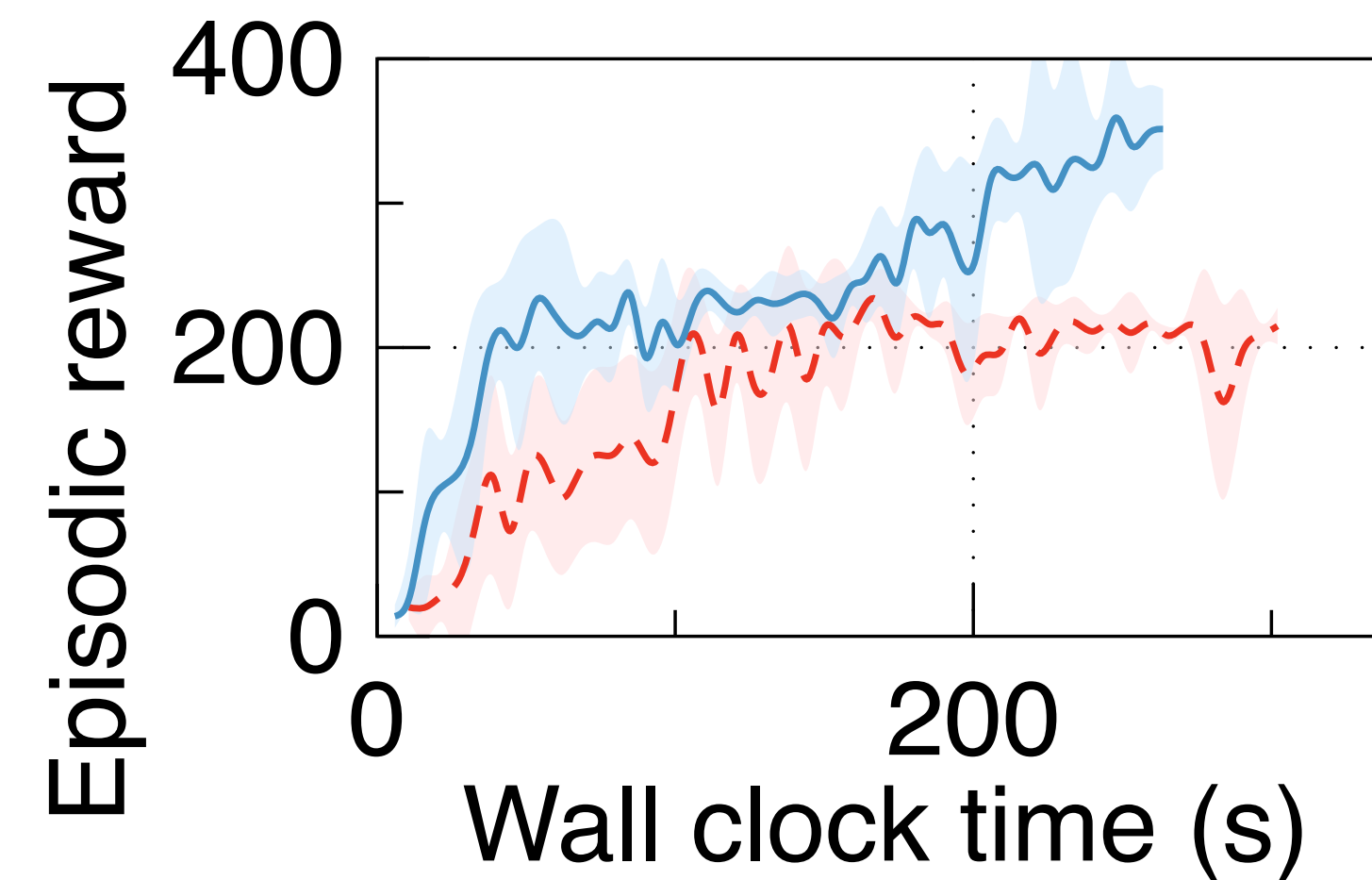
Training Performance

Faster Training

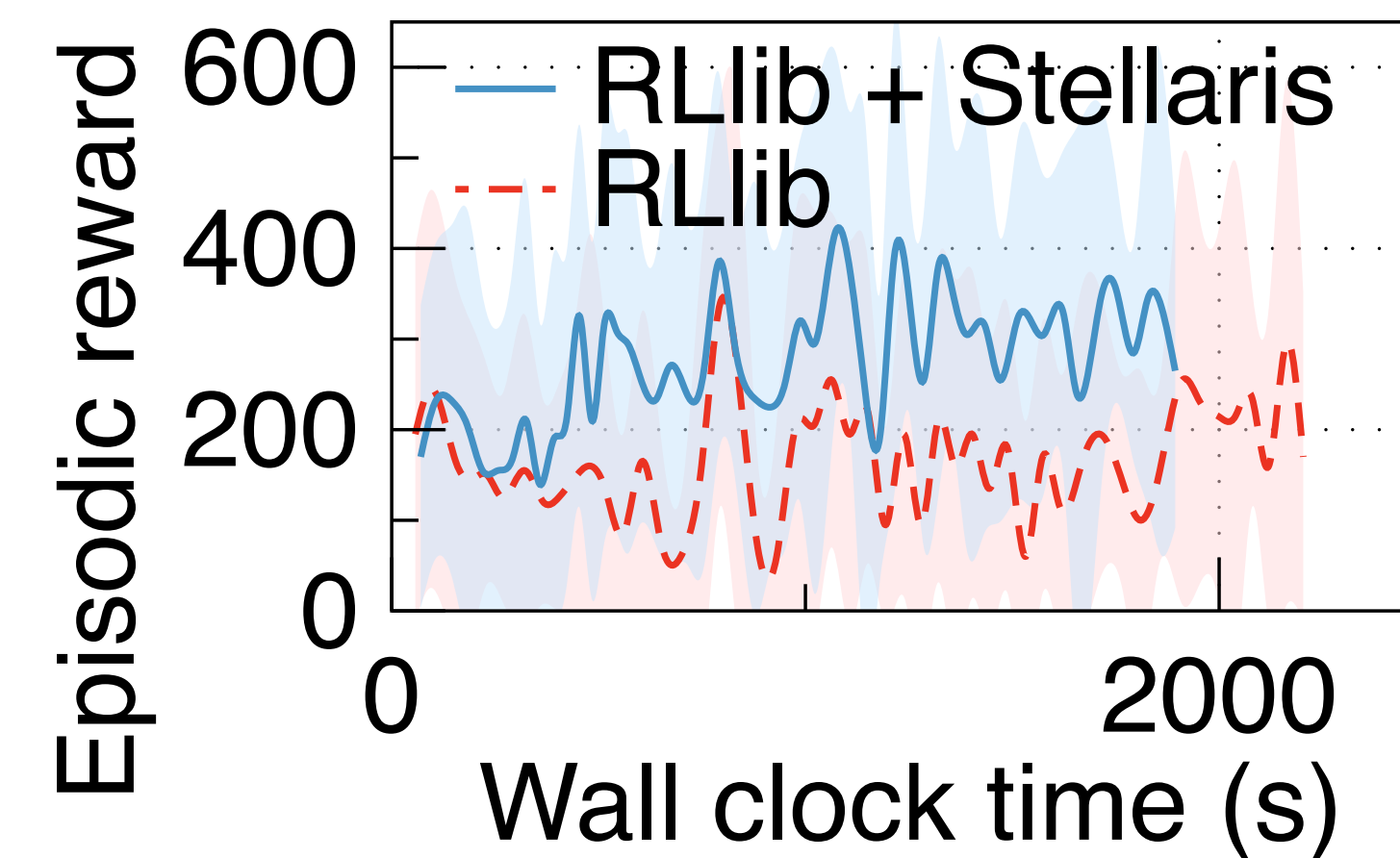
Higher Rewards

2.2x

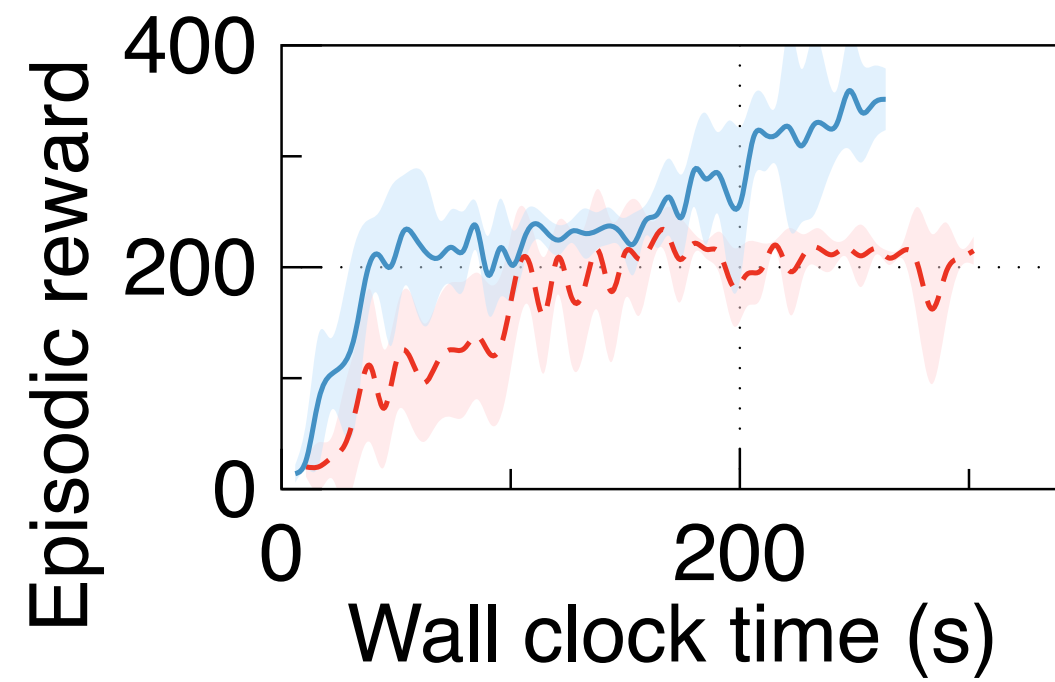
Training performance improvement



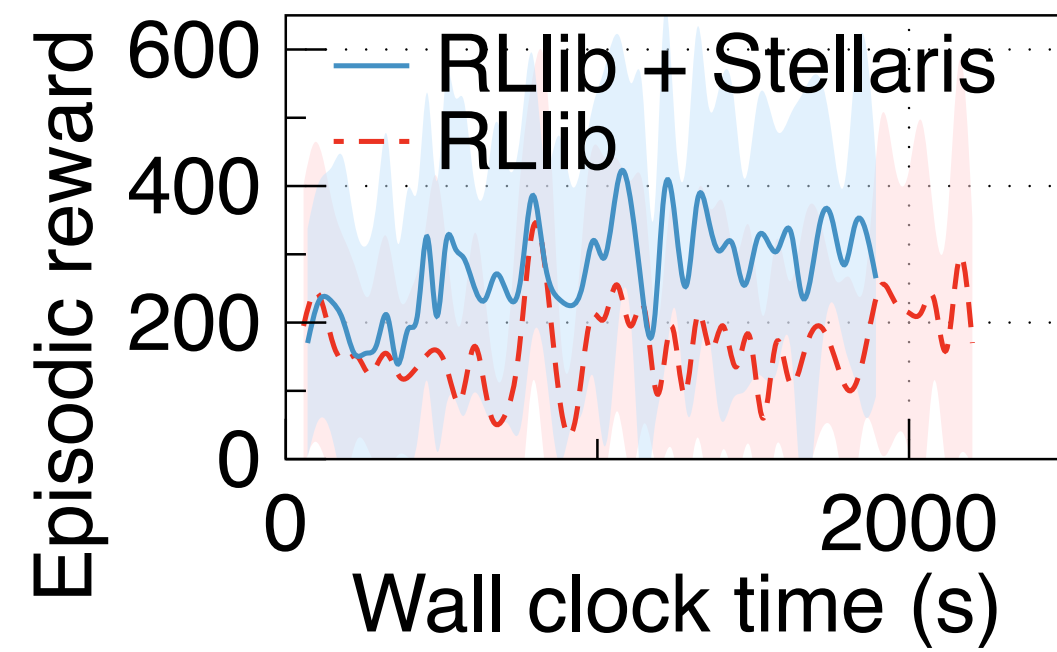
(a) Hopper



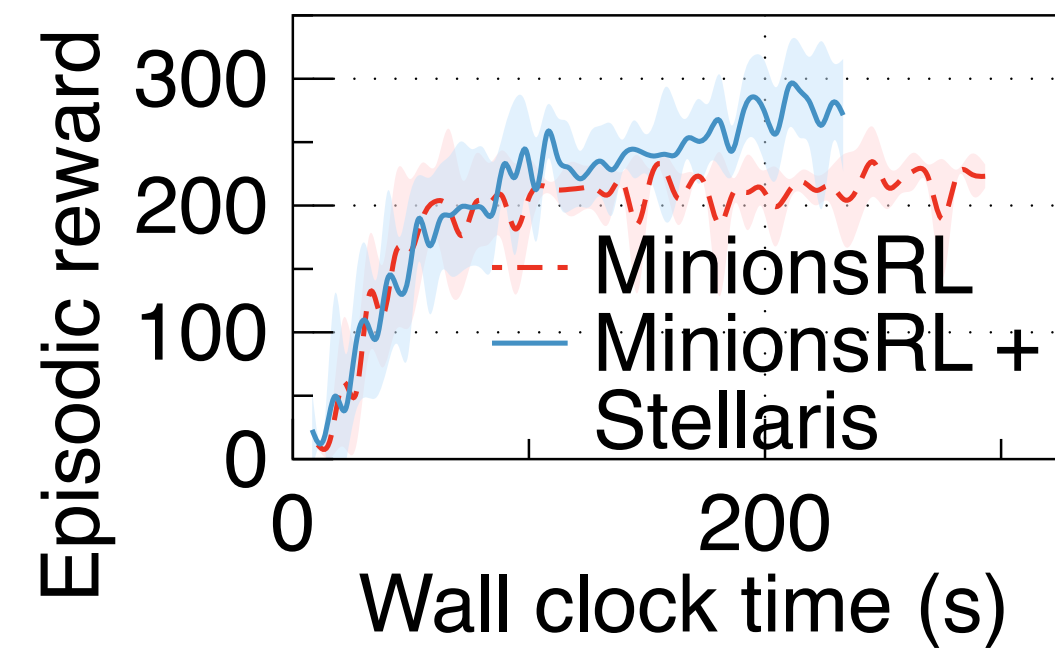
(b) Gravitar



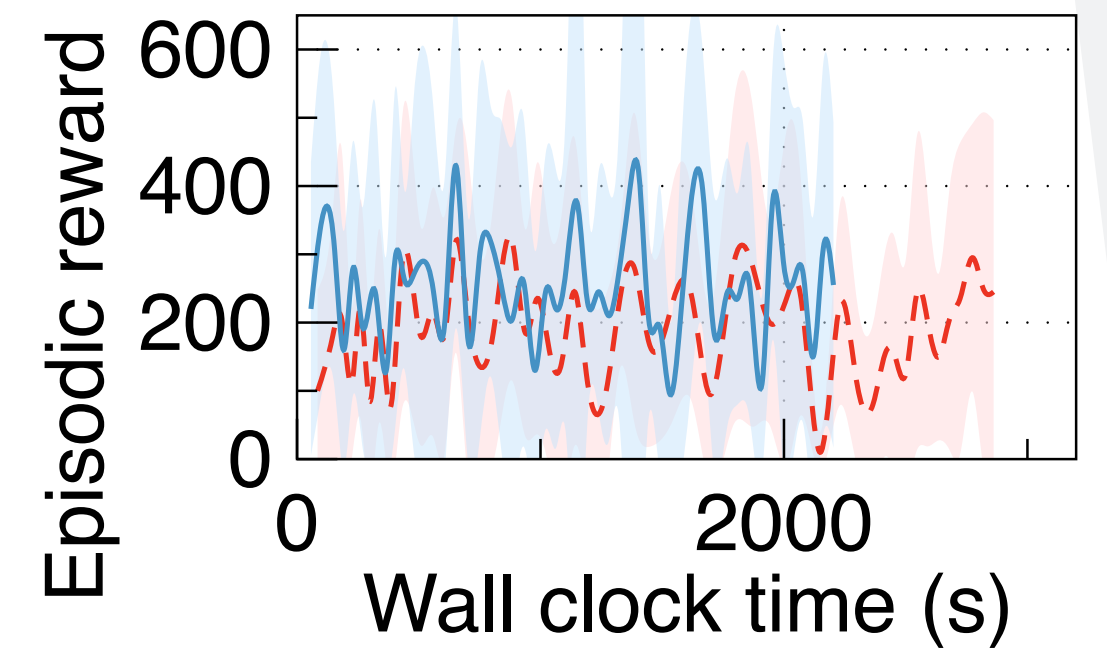
(a) Hopper



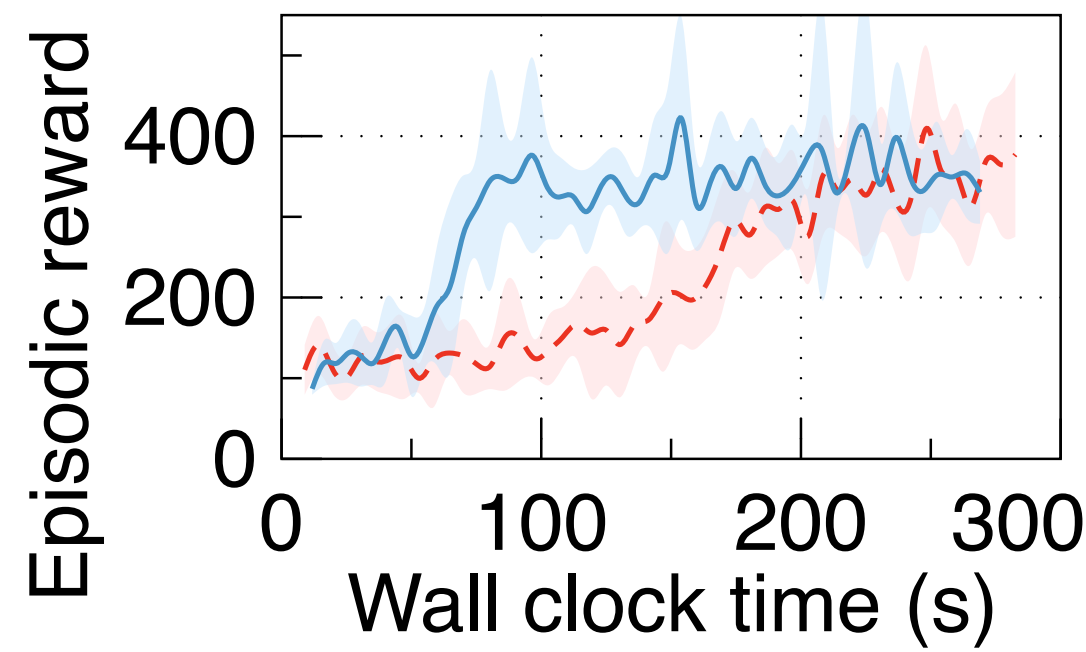
(d) Gravitar



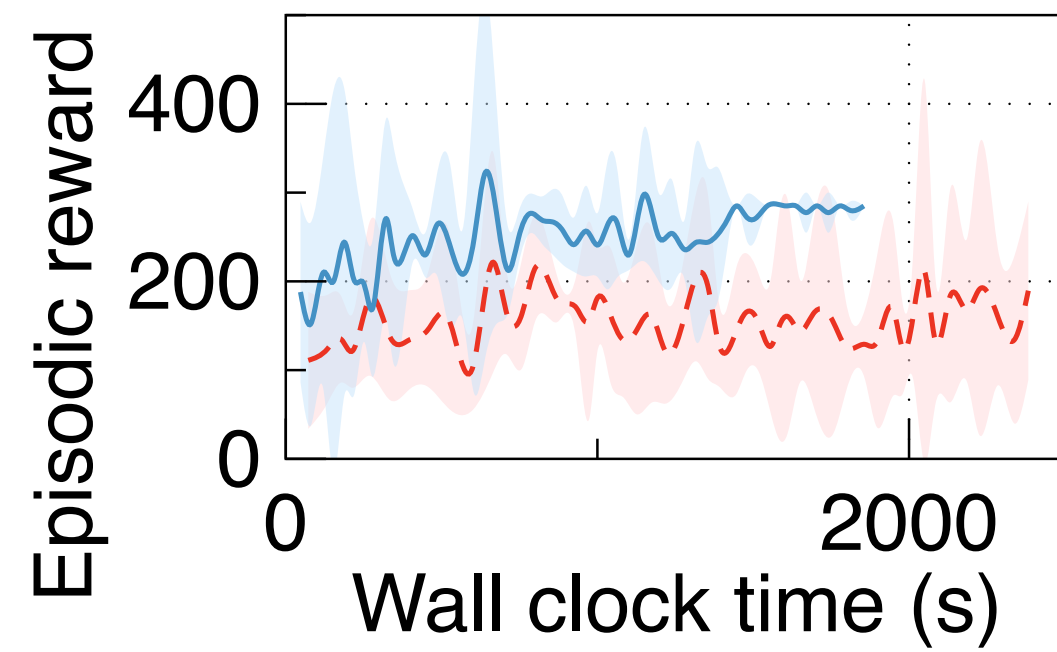
(a) Hopper



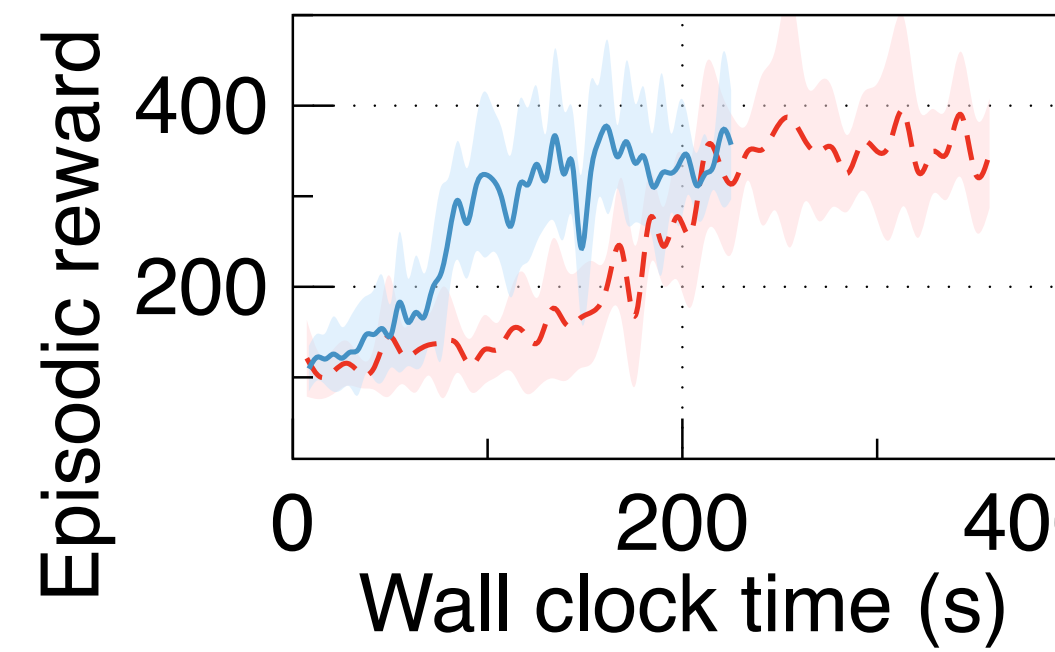
(d) Gravitar



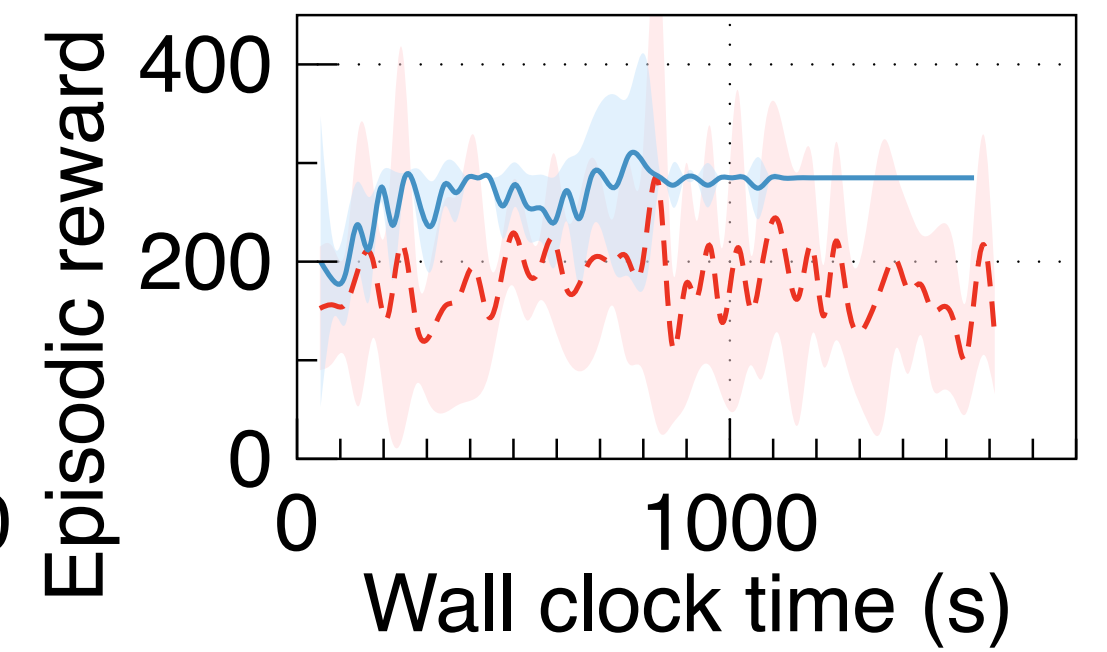
(b) Humanoid



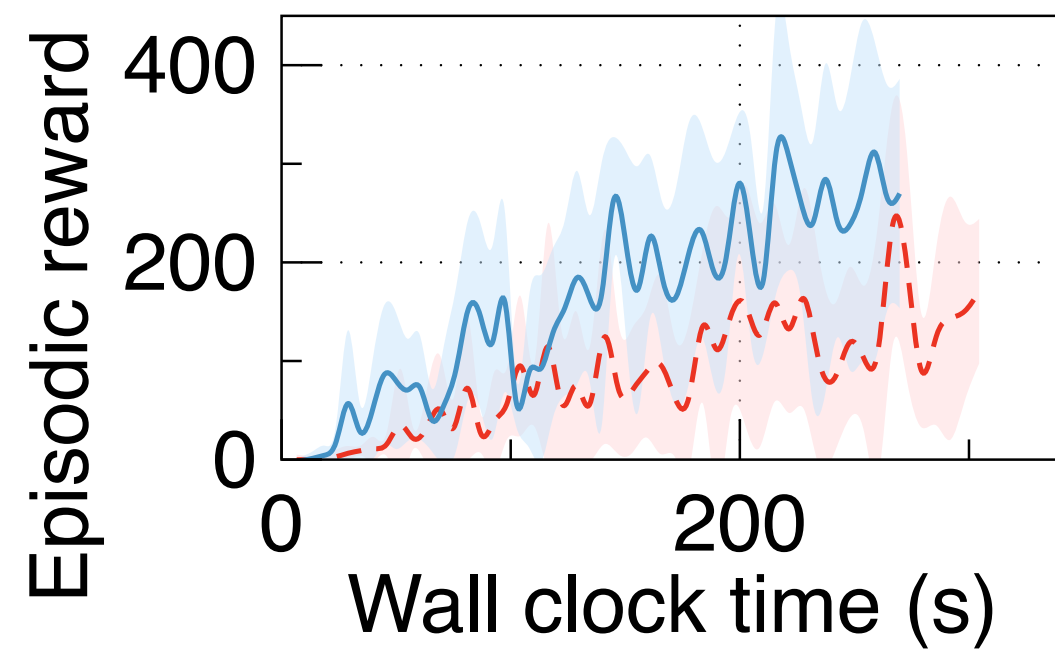
(e) SpaceInvaders



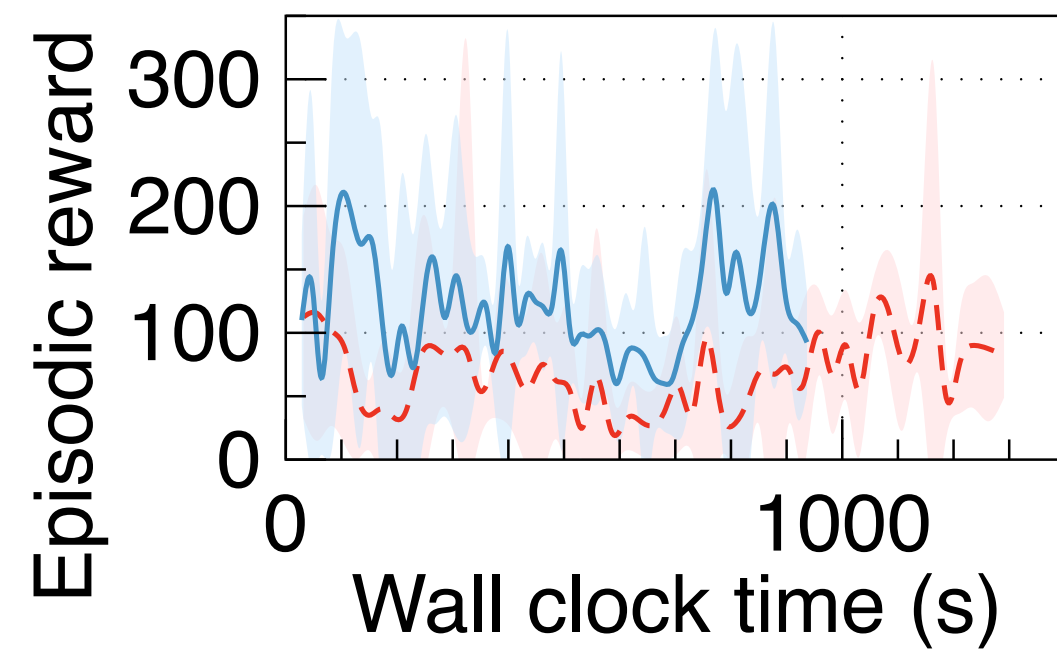
(b) Humanoid



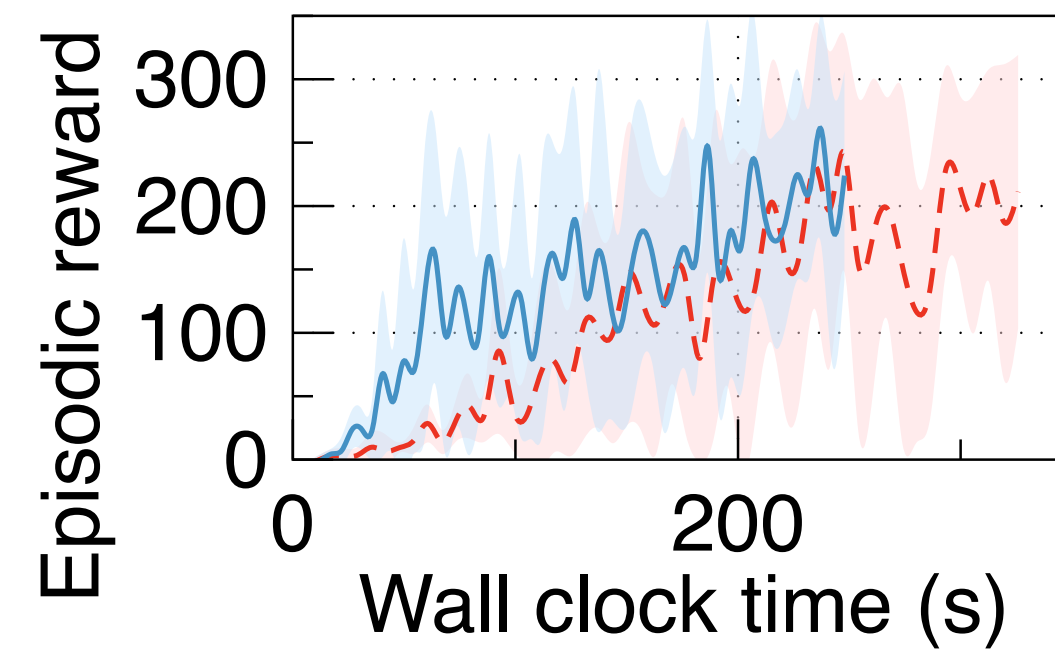
(e) SpaceInvaders



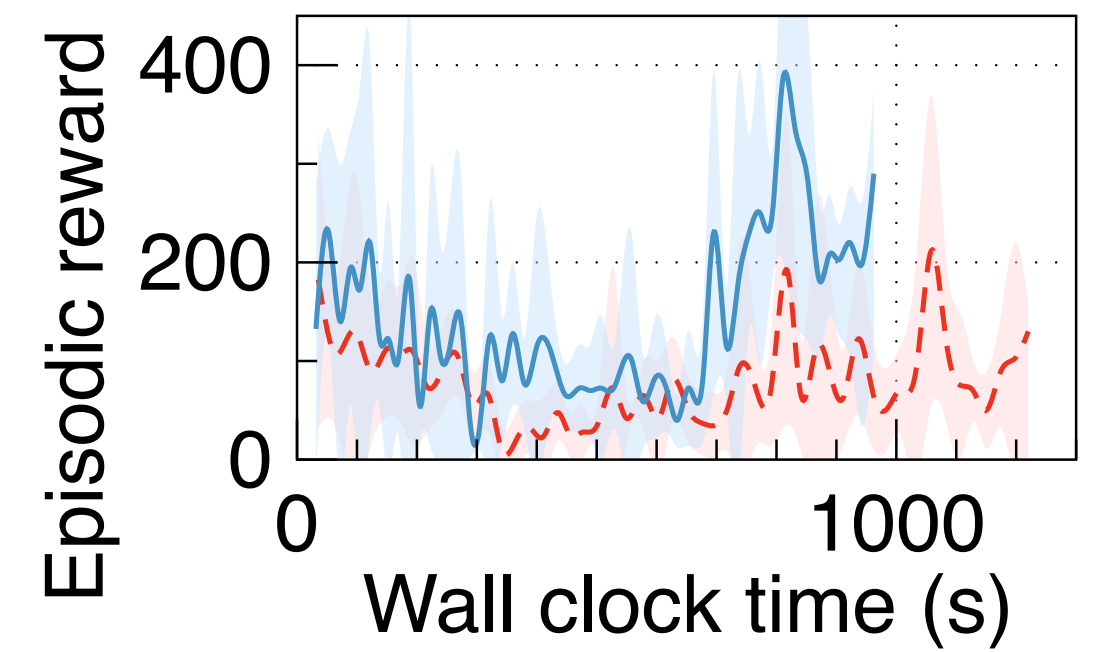
(c) Walker2d



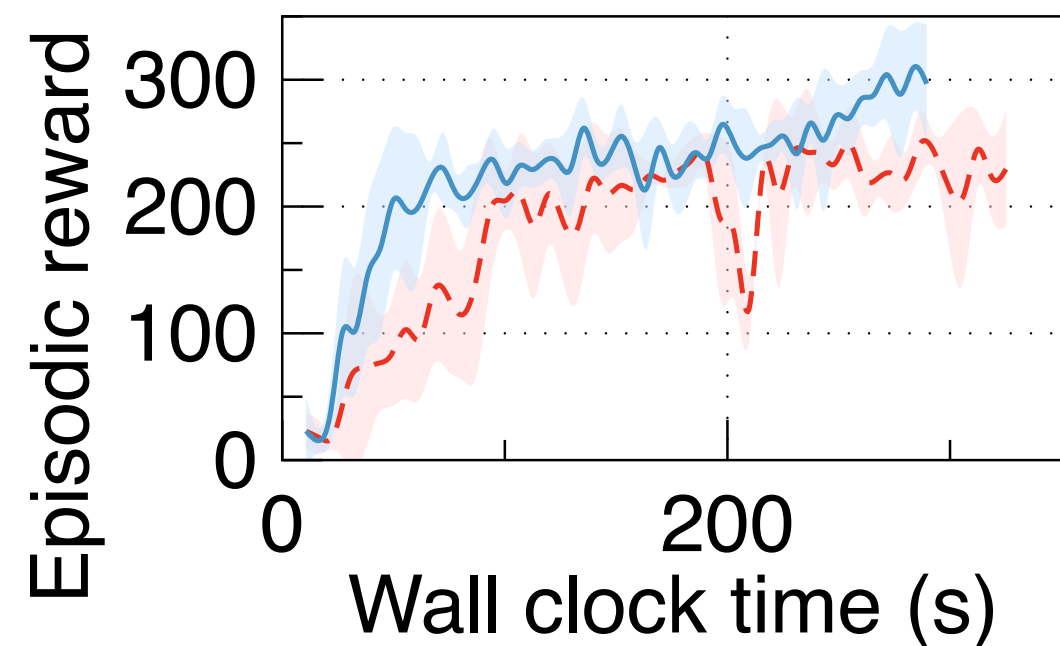
(f) Qbert



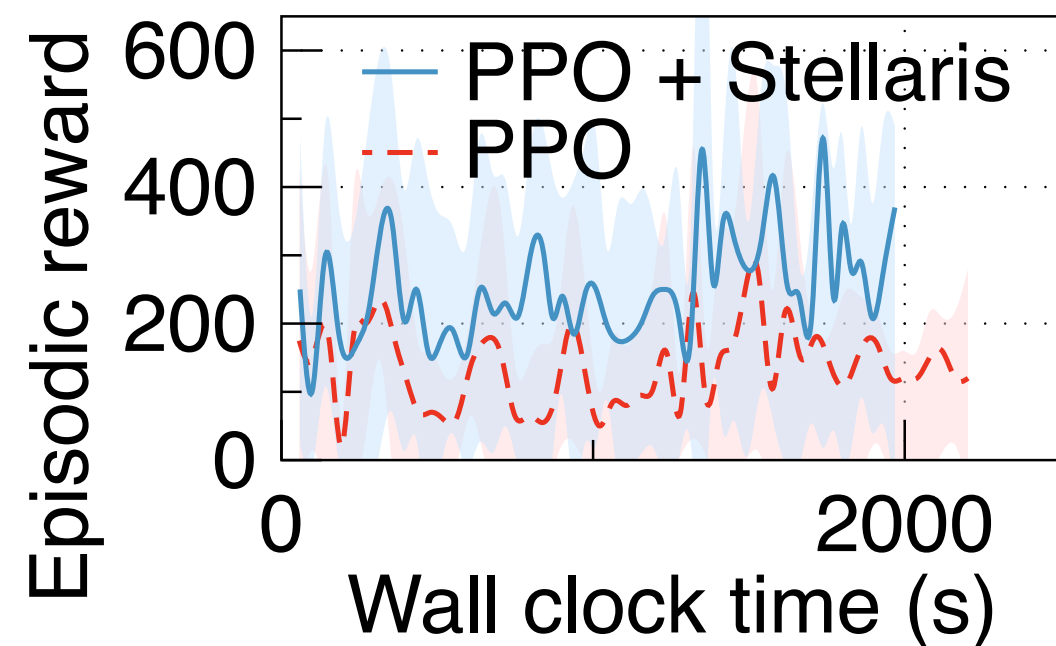
(c) Walker2d



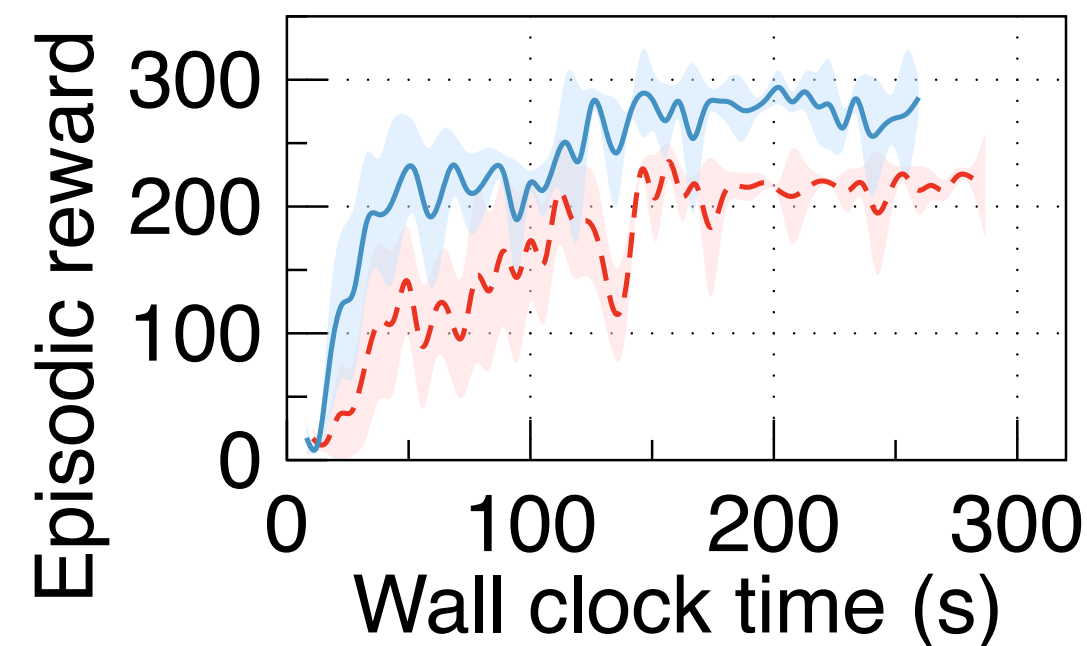
(f) Qbert



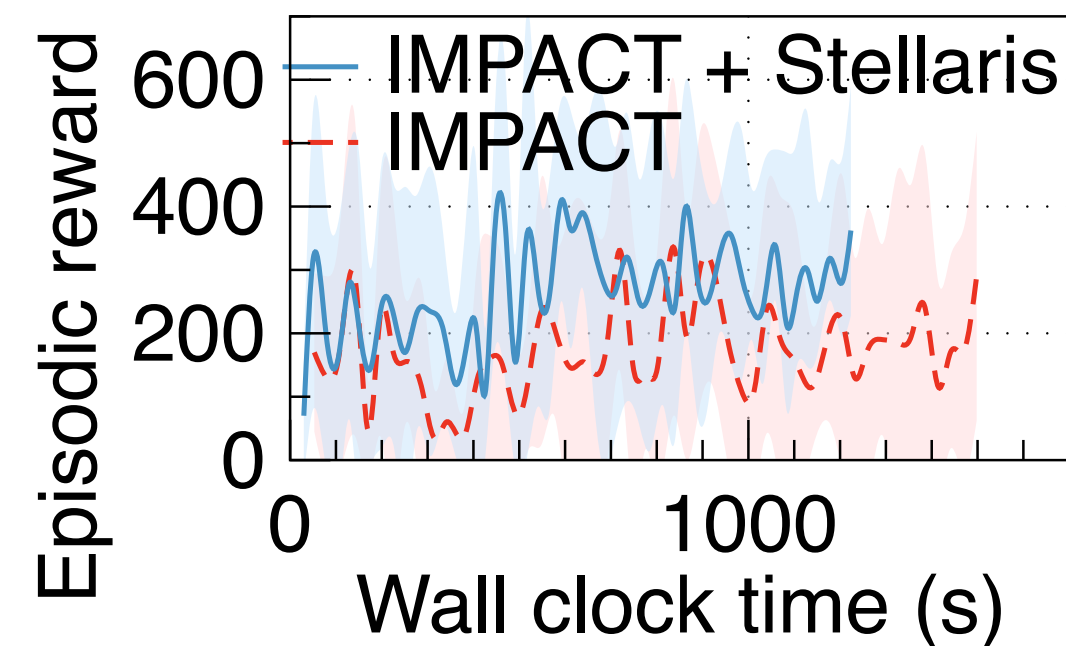
(a) Hopper



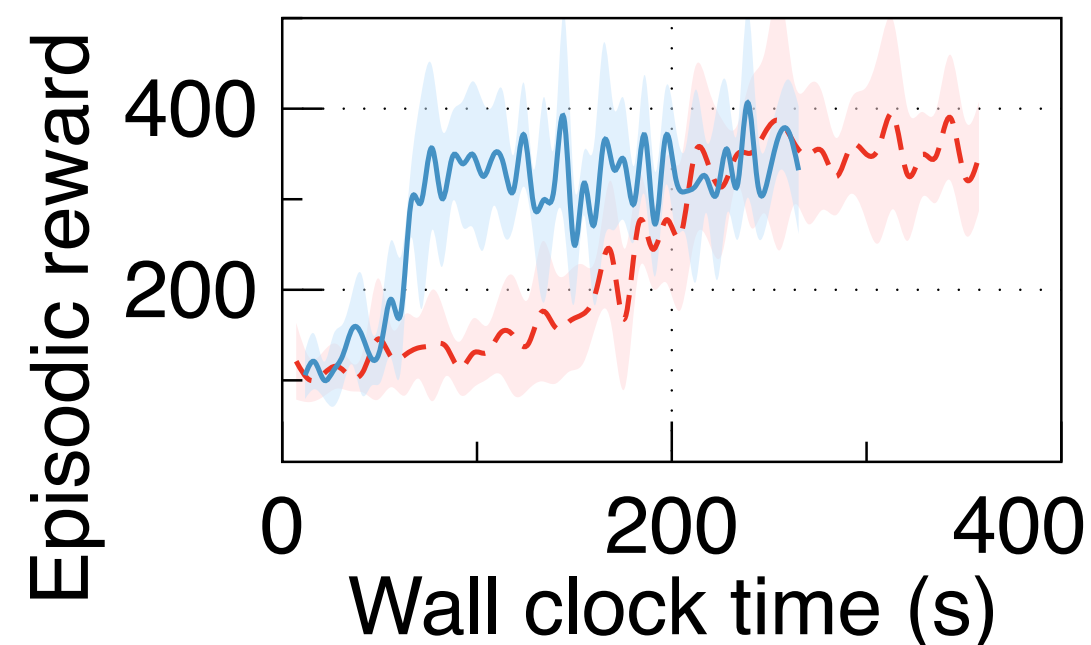
(d) Gravitar



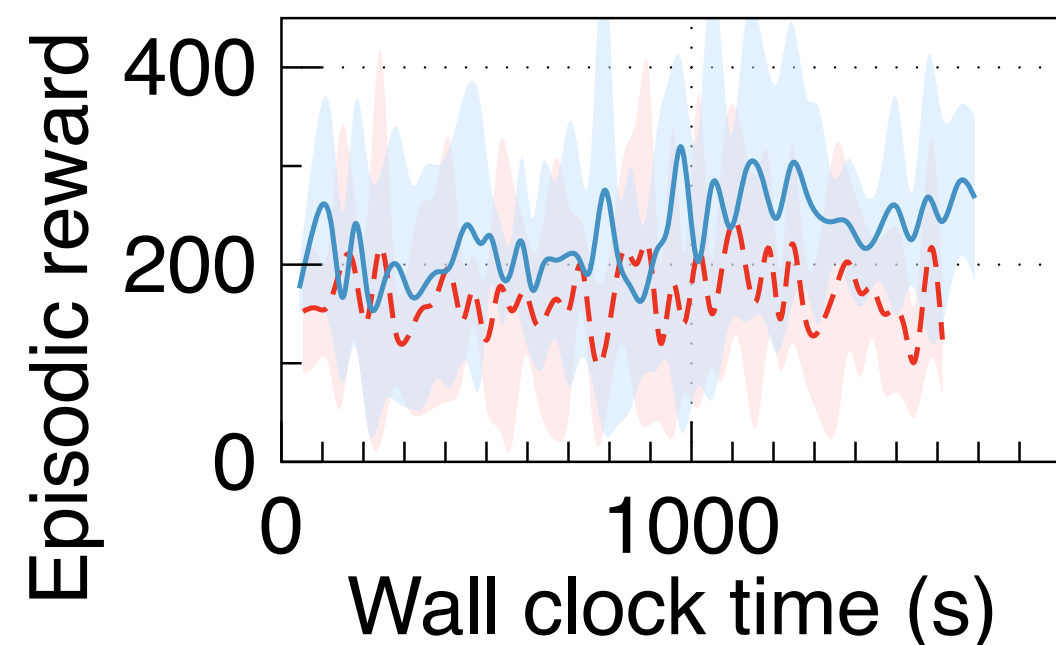
(a) Hopper



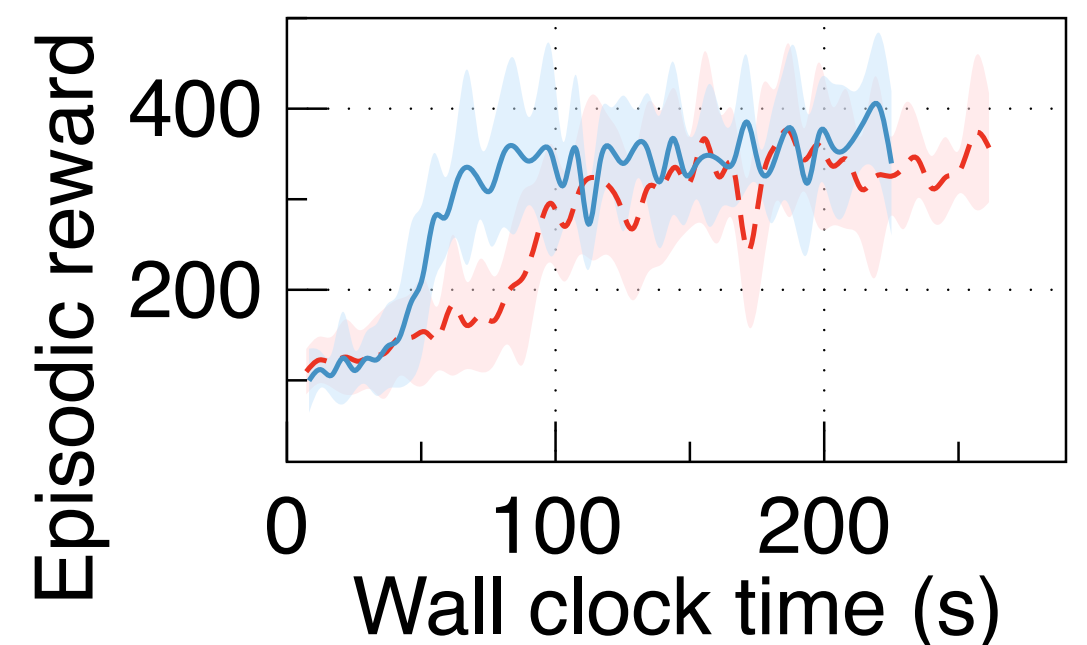
(d) Gravitar



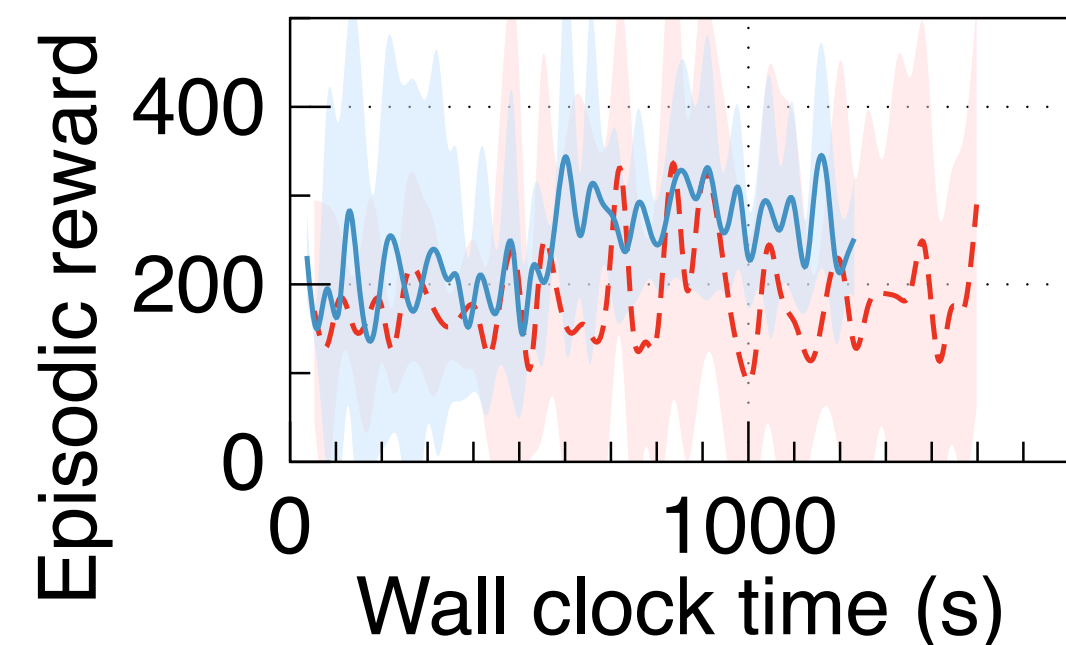
(b) Humanoid



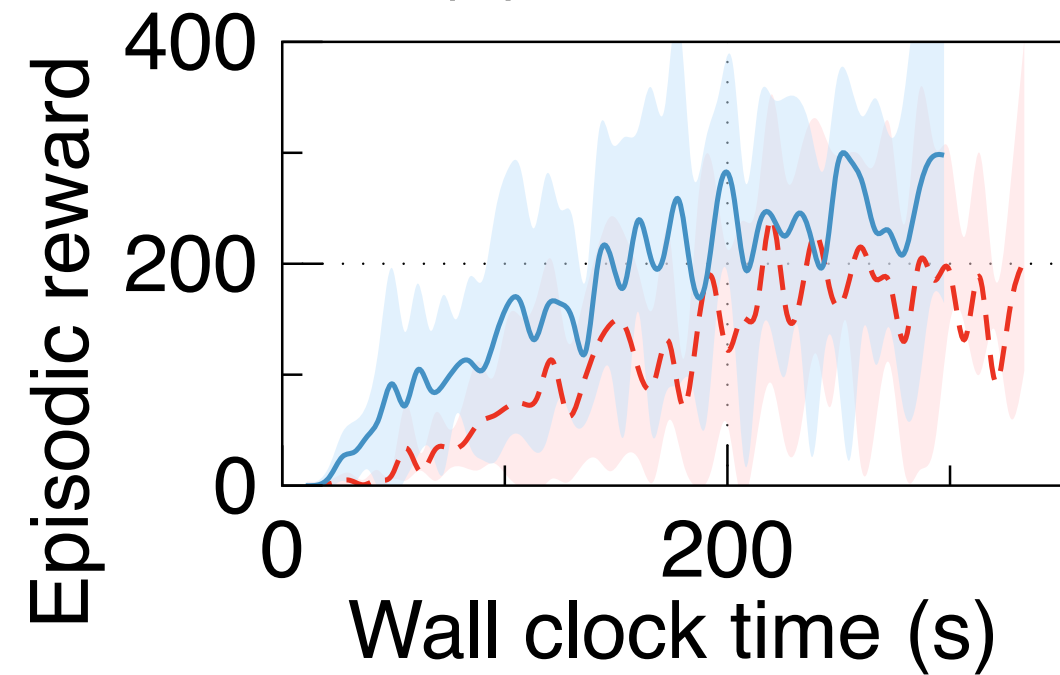
(e) SpaceInvaders



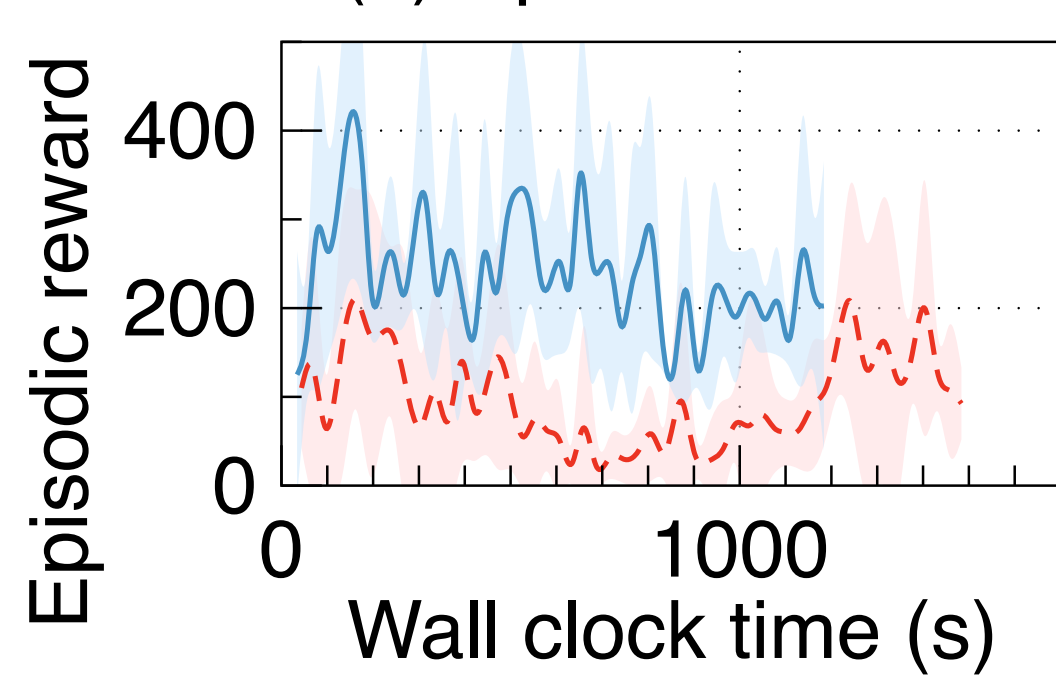
(b) Humanoid



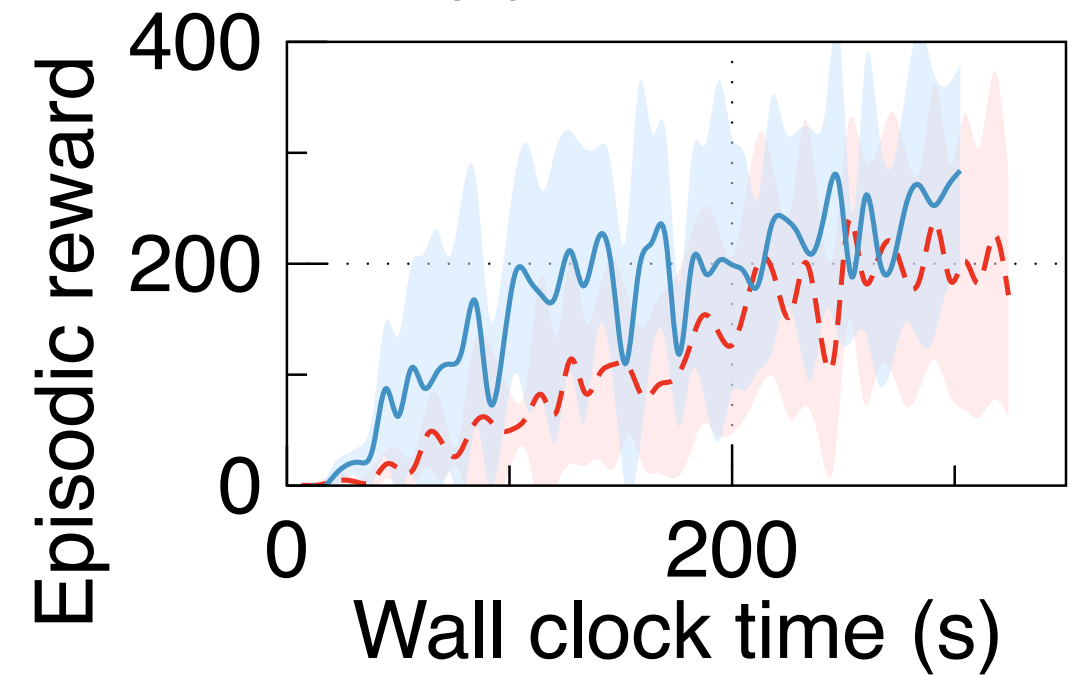
(e) SpaceInvaders



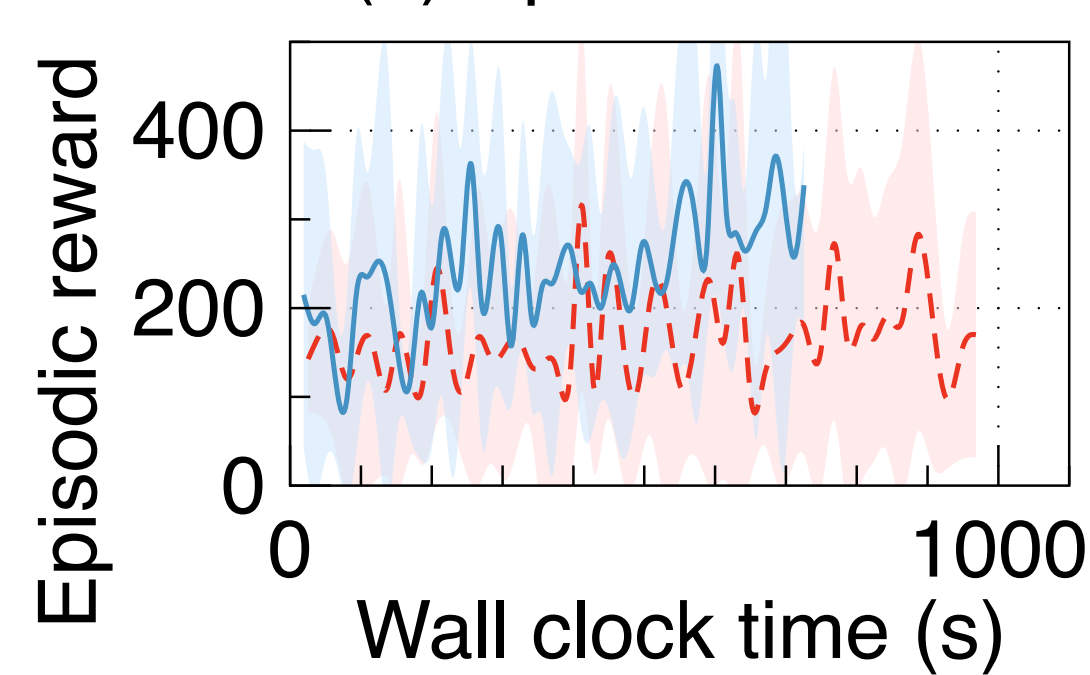
(c) Walker2d



(f) Qbert

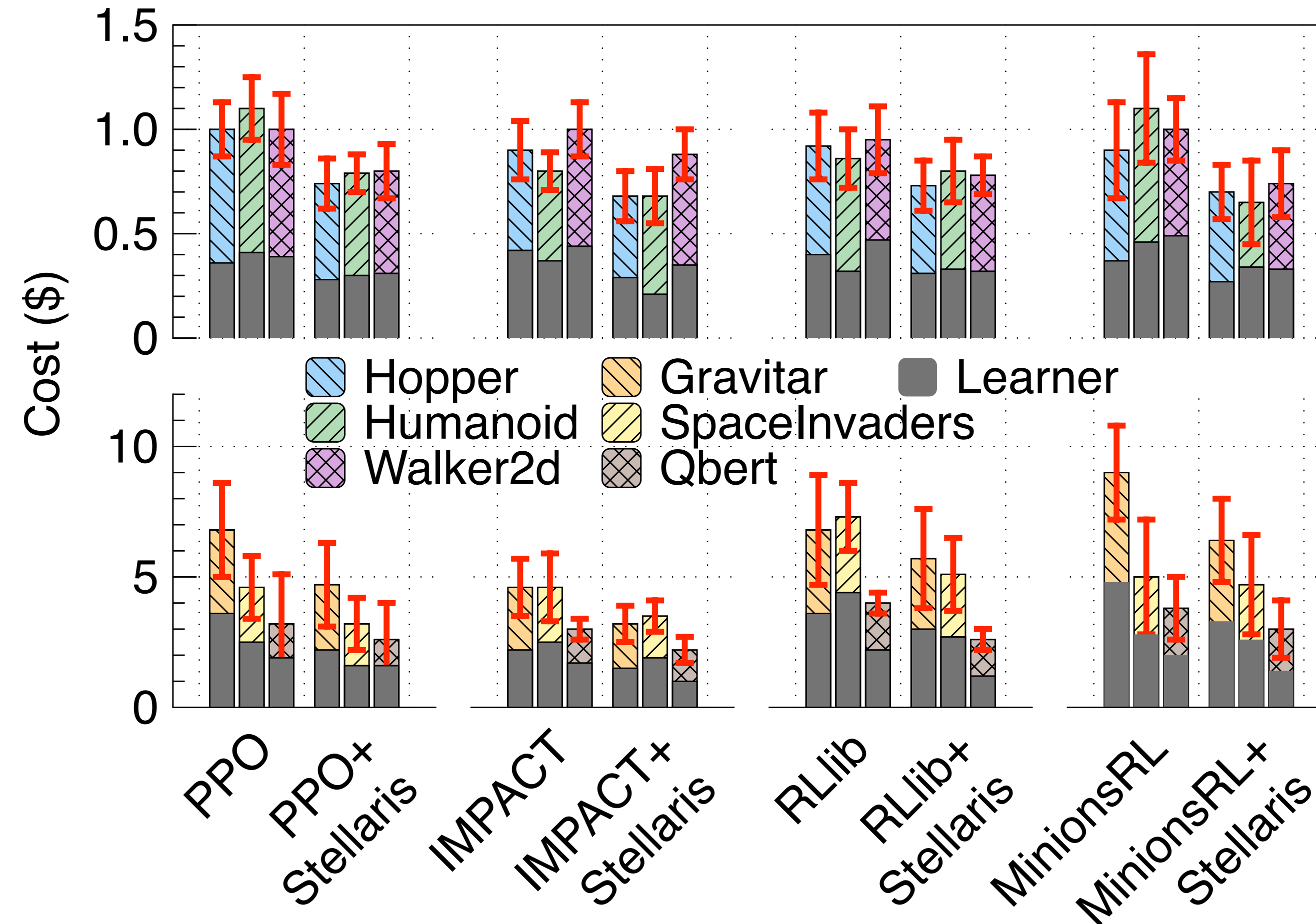


(c) Walker2d



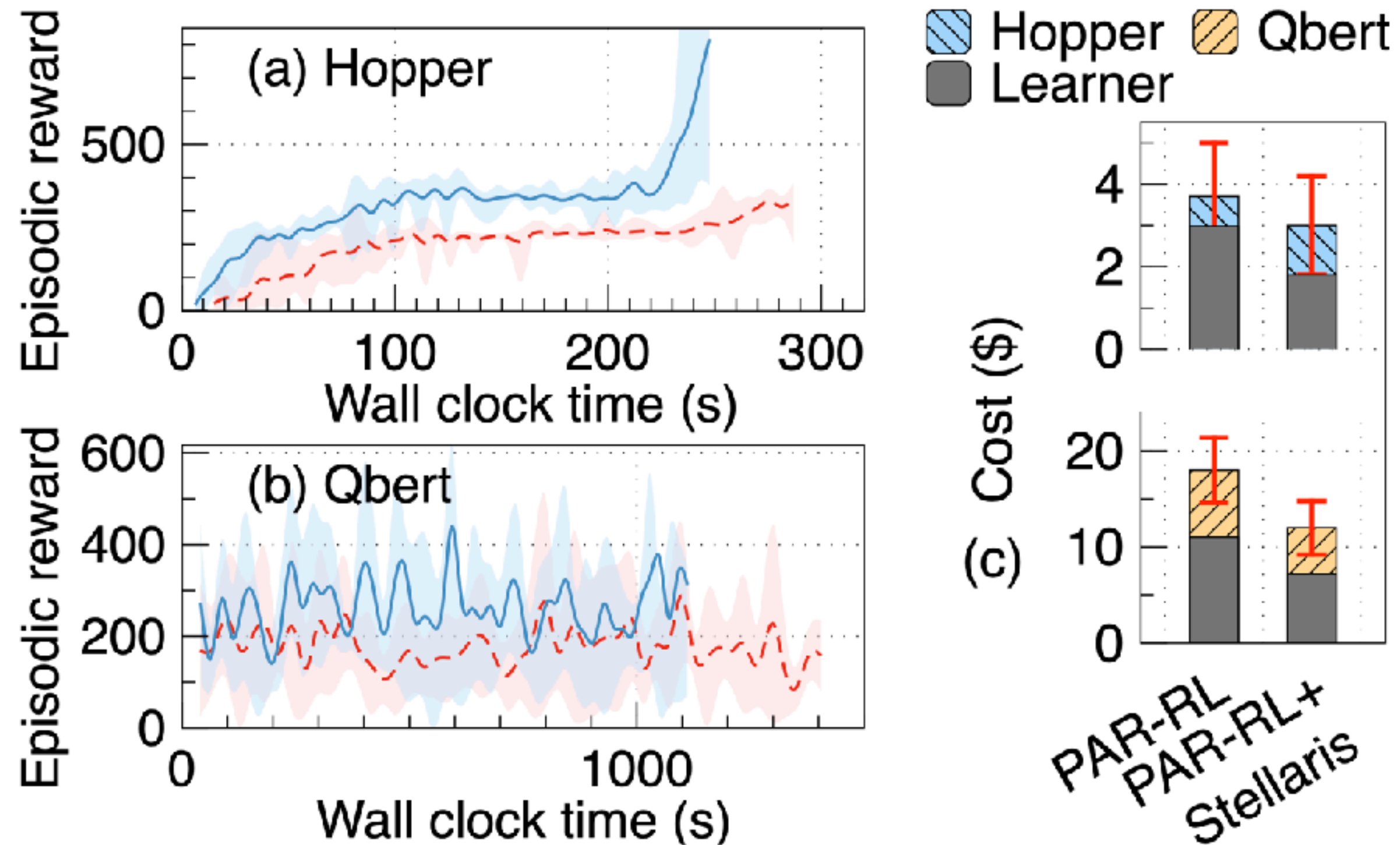
(f) Qbert

Training Cost



41%
Training cost reduction

Improving PAR-RL [1] on HPC Testbeds



2.4x

Training performance improvement

34%

Training cost reduction

Asynchronous
Serverless Learners

Global Importance
Sampling Truncation

Staleness-Aware
Gradient Aggregation

Stellaris

2.2x

Training performance improvement

41%

Training cost reduction



Stellaris Code Repository:

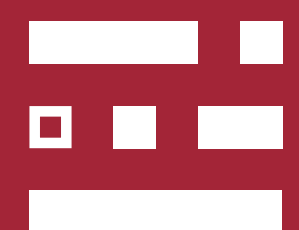
<https://github.com/IntelliSys-Lab/Stellaris-SC24>

Corresponding Author:

Hanfei Yu <hyu42@stevens.edu>

Hao Wang <hwang9@stevens.edu>

I'm looking for a research internship in summer 2025!



IntelliSys Lab



GitHub Code



About Me



THANK YOU

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