

Learn To Race

**Distributed RL and Optimization Thrust
- Final Presentation**

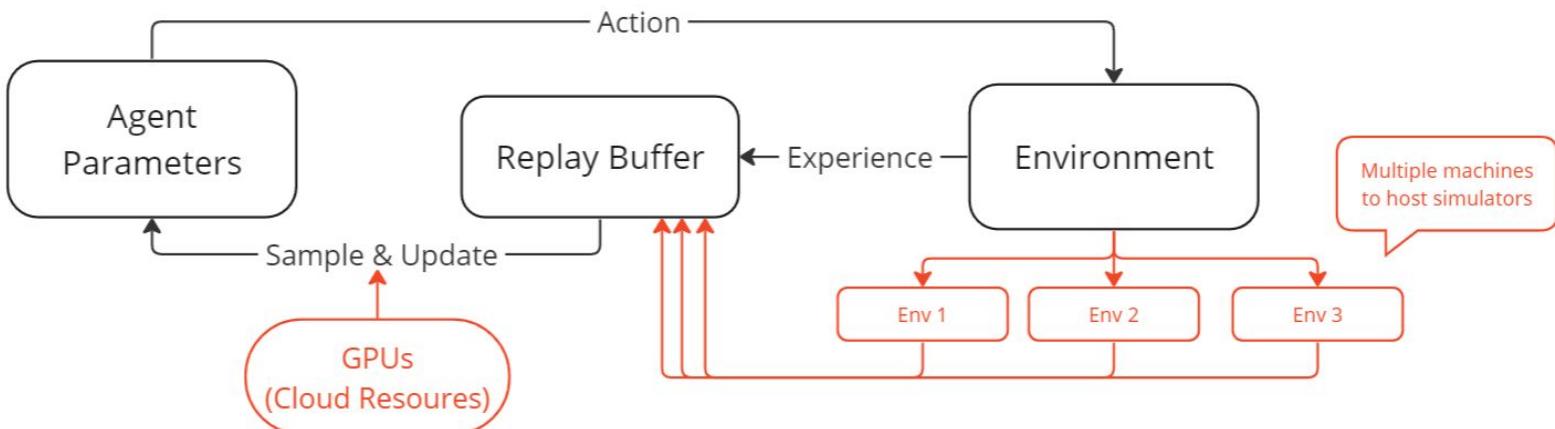
Shijie (Brandon) Bian

Advisors: Jon Francis, Eric Nyberg, Arav Agarwal

Introduction

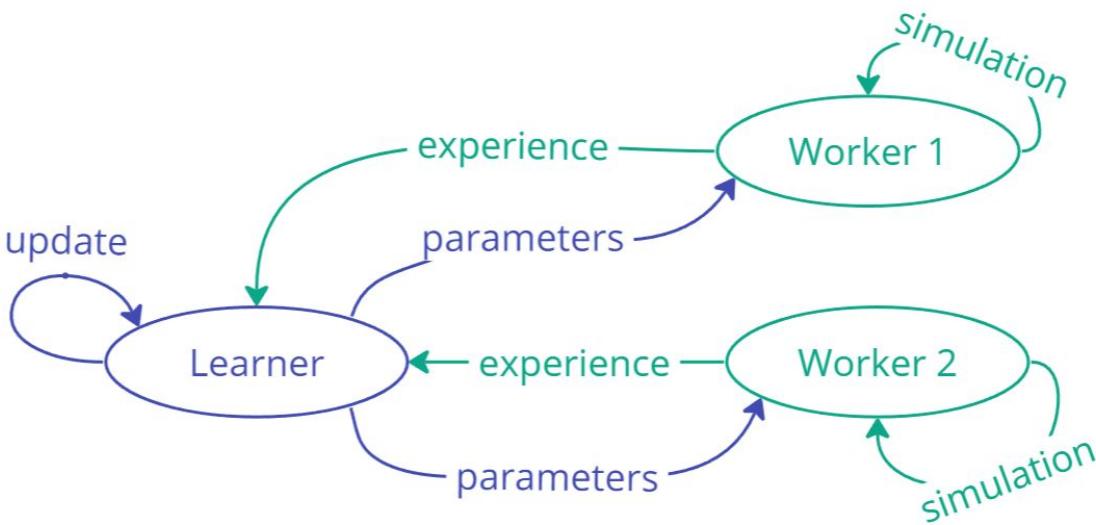
Introduction

- **Reinforcement Learning (RL)** may take days or weeks to train.
- Can we use **multiple GPUs** to accelerate training?
- Can we collect experience from **diverse environments** at the same time?



Relation to Prior Work

- Previous iteration: IMPALA architecture [1]
- Learner: performs parameter update
- Worker (i.e., Actor): performs simulation and experience collection



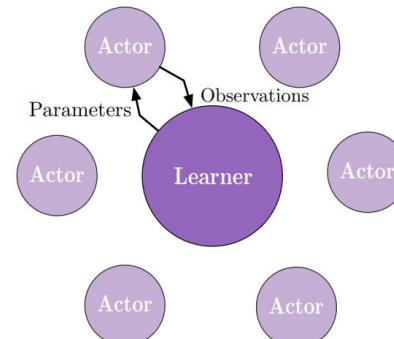
Problem Description

Problem Description

1. Performance does not scale

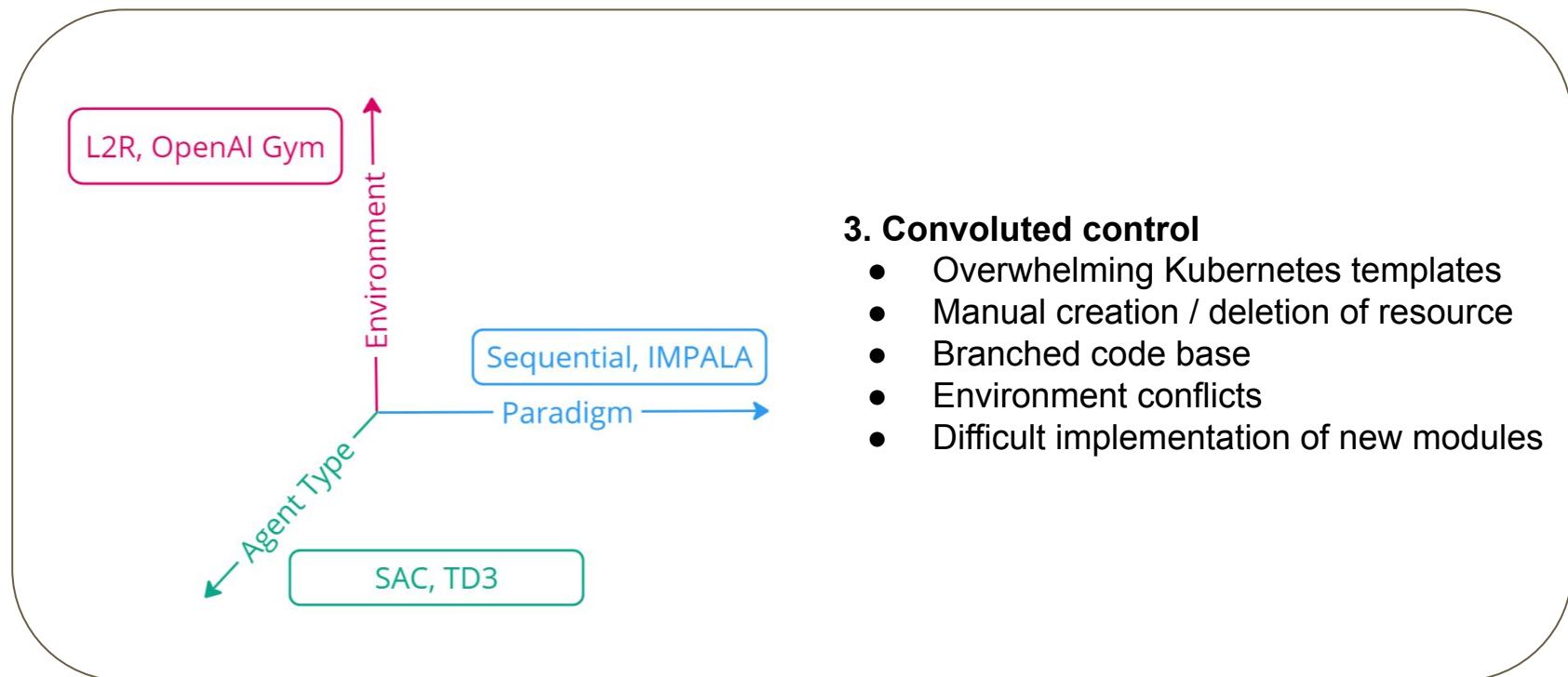


2. Learner overwhelmed



[1]

Problem Description



Proposed Solution

<u>Solution</u> Distributed Training Processes & Resource Control Optimizations	
Challenge	Solution
<ul style="list-style-type: none">• Performance not scaling	<ul style="list-style-type: none">• Utilizing OpenAI Gym environments and SAC implementations as baselines
<ul style="list-style-type: none">• Learner node being overwhelmed	<ul style="list-style-type: none">• Distribute training tasks to worker nodes
<ul style="list-style-type: none">• Convoluted control	<ul style="list-style-type: none">• Automated and unified resource control and deletion

Hypothesis

1. Training efficiency not scaling

- Learner being a bottleneck
- Agent and parameter configurations are not suitable for the task

2. OpenAI's environment and agent configuration

- Serves as easier baselines to test the distributed system architecture

3. Training tasks can be allocated to workers

- Alleviate workload on learner

4. Modularization of the paradigms

- Helpful for trying out combinations of paradigms and RL environment

Features and Functionality

1. New distributed training paradigm

- Training tasks allocated to workers
- Learner responsible for controlling worker and distributing parameters

2. New baselines

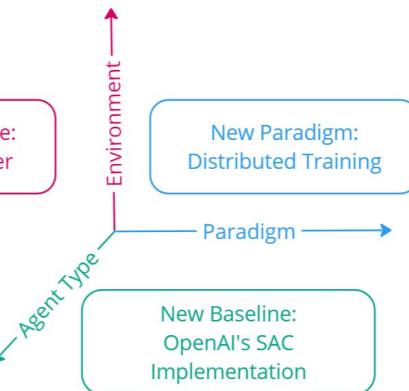
- OpenAI Spinning-Up SAC agent and replay buffers
- OpenAI Lunar Lander environment
- Full compatibility with the distributed system architecture

3. Automated resource control

- Unified modules and Kubernetes templates
- Resource launching and shutdown based on user input
- Additional logging for both learner and worker across paradigms

4. Modularization

- Easier plug-in and implementation of new Paradigms / Agents / Environments



Resource Requirement

Data

- Arrival simulator: visual data input for encoding environment state
- Training and evaluation metrics: such as reward of agent, and training step time (simulation time)



ARRIVAL

Hardware

- Phoebe clusters: Nvidia GPU (> 10 GB memory) and CPU (> 8 cores) support, allowing multi-pod allocation
- Cloud resource: AWS EC2 (g4dn.2xlarge) is used for system framework debugging purposes

Software

- Pytorch: for model training with GPU
- Kubernetes: for resource creation and control
- WandB: for metrics logging and visualization
- Mamba/Conda: for automated environment creation

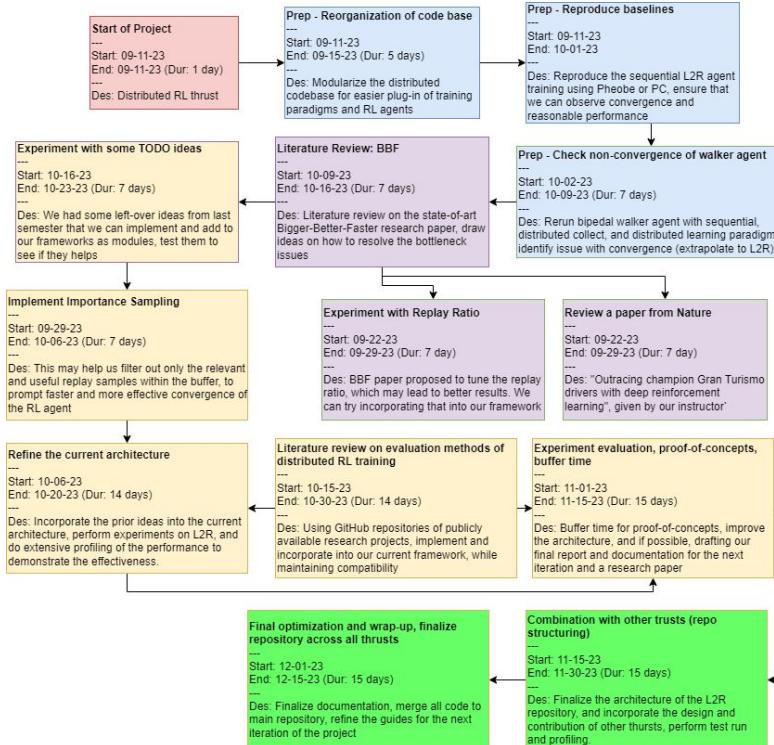


kubernetes



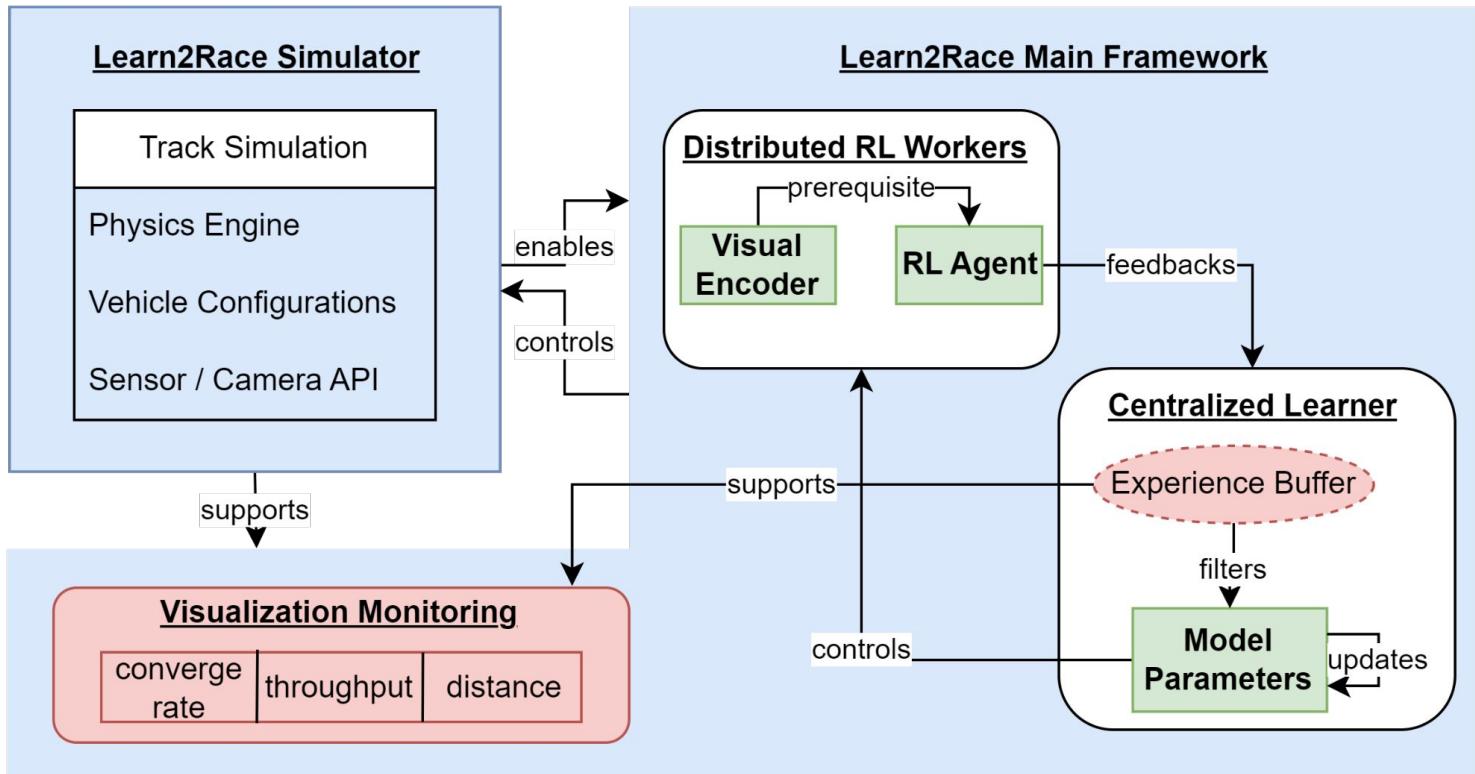
Weights & Biases

Plan

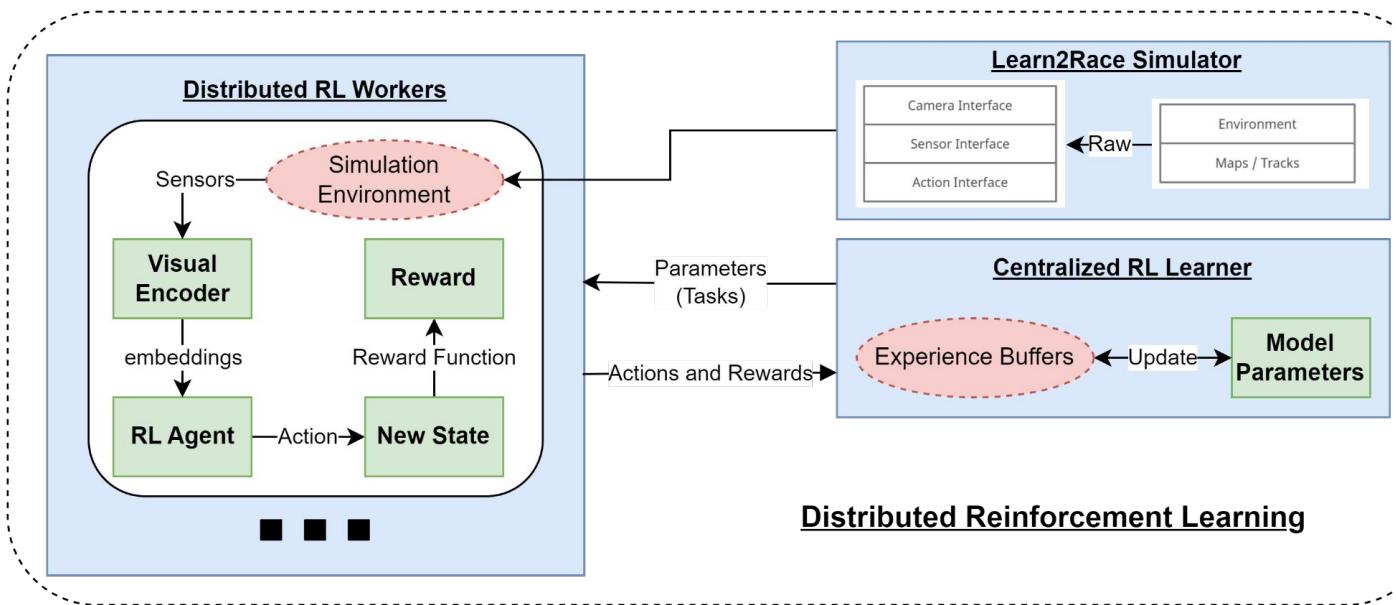


Design - Overview

Learn-to-Race Context



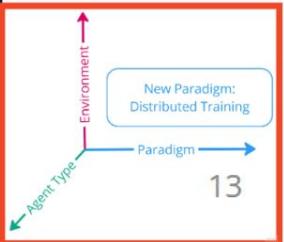
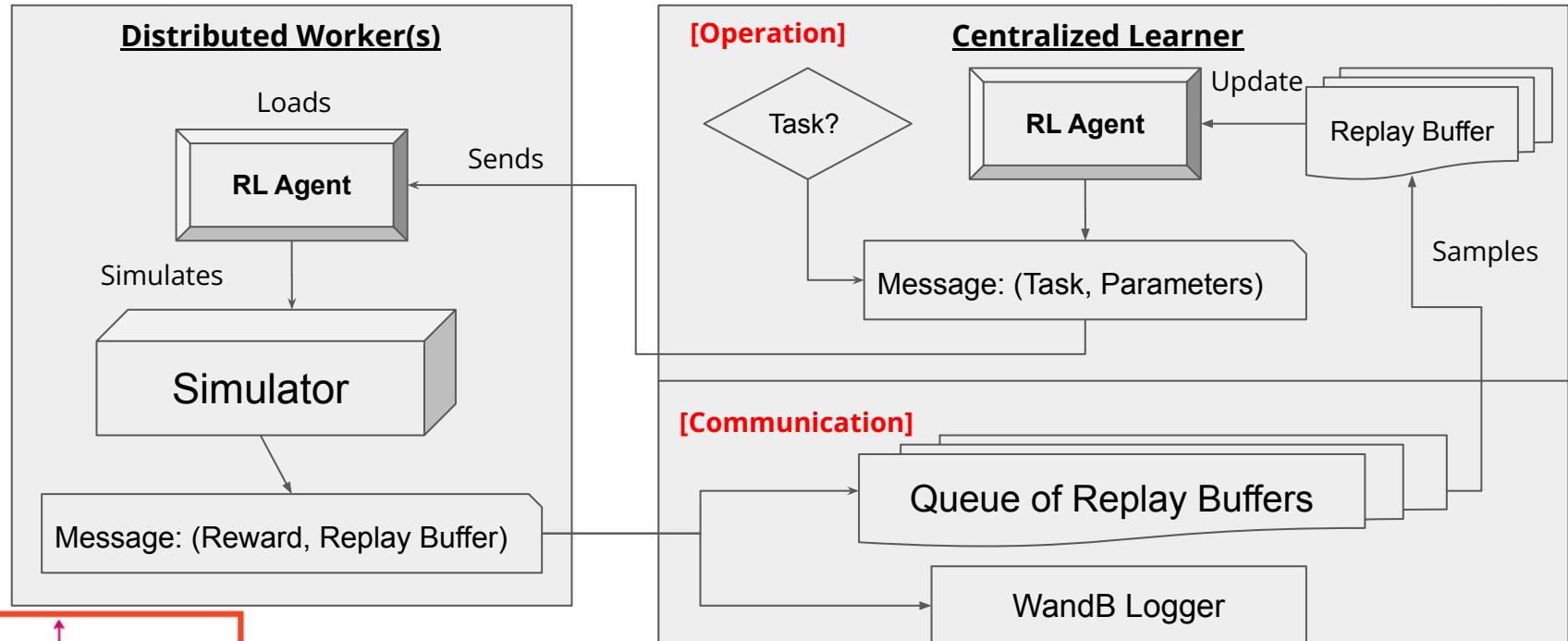
Paradigm - High-level Design



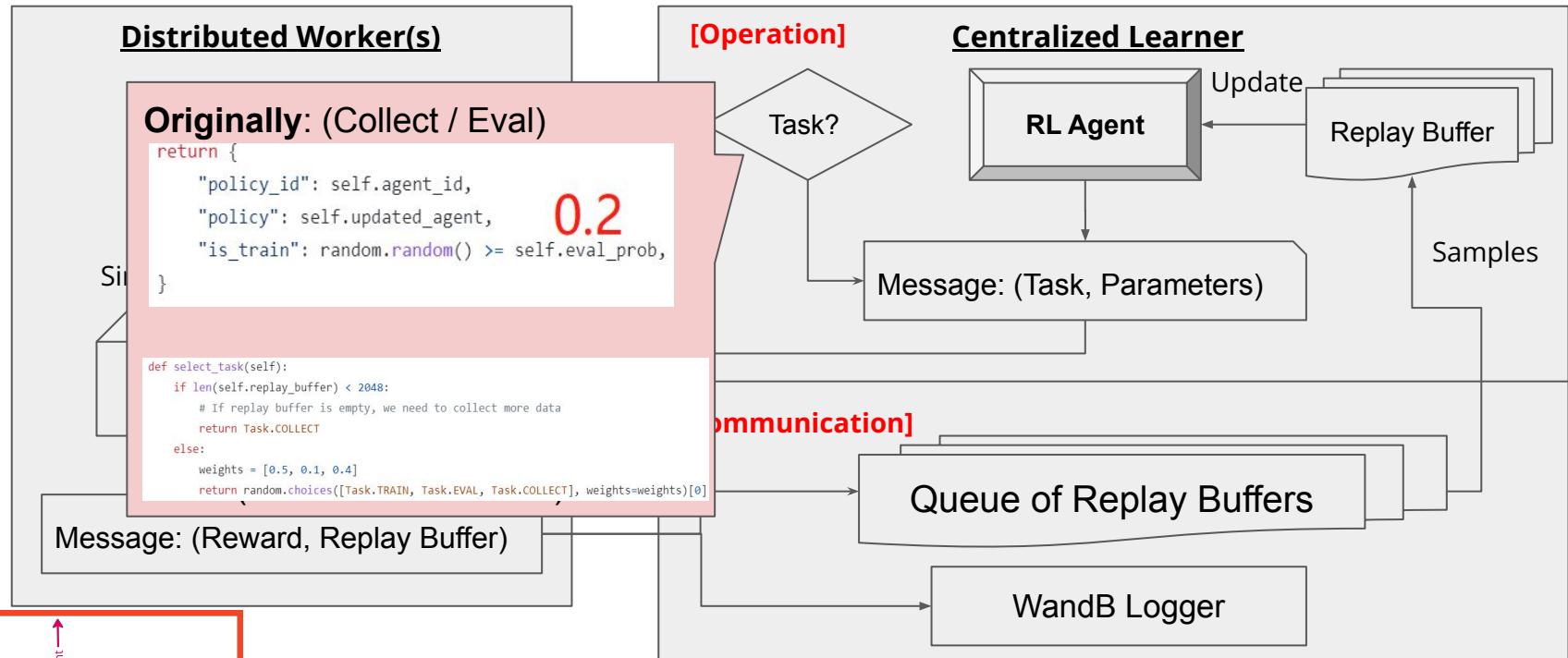


Design - Training Paradigms

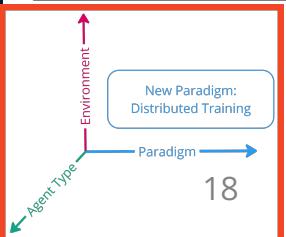
Paradigm - Distributed Training / Update



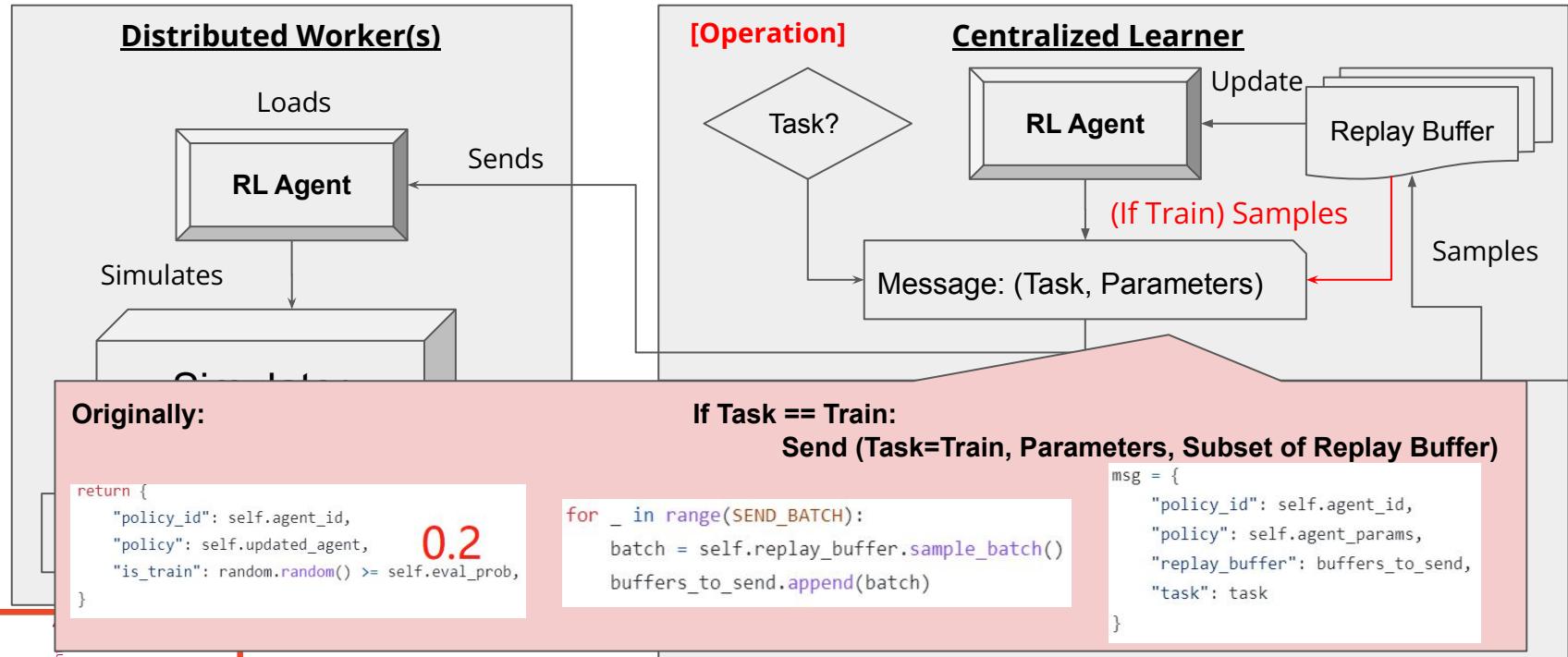
Paradigm - Distributed Training / Update



(1) Allocating training task to workers

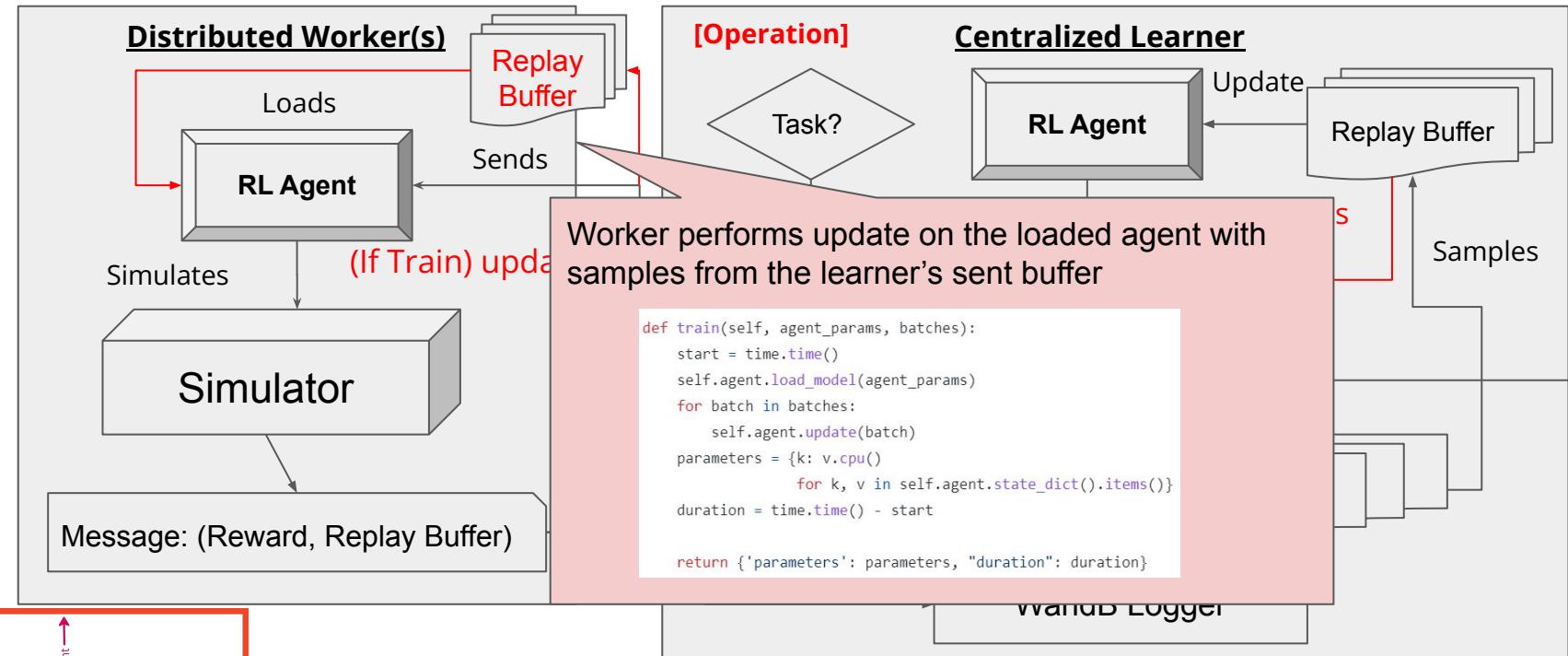


Paradigm - Distributed Training / Update



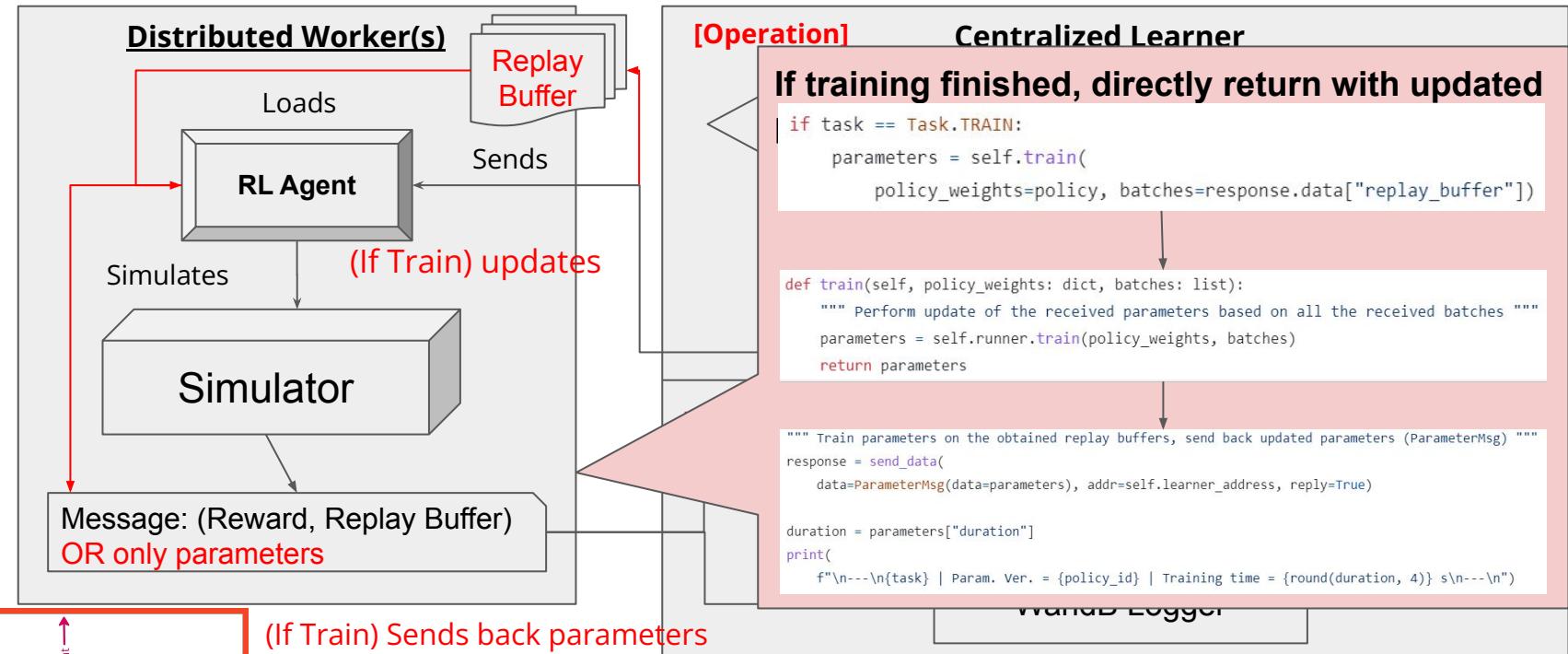
(2) For training tasks, send samples from learner's buffer

Paradigm - Distributed Training / Update



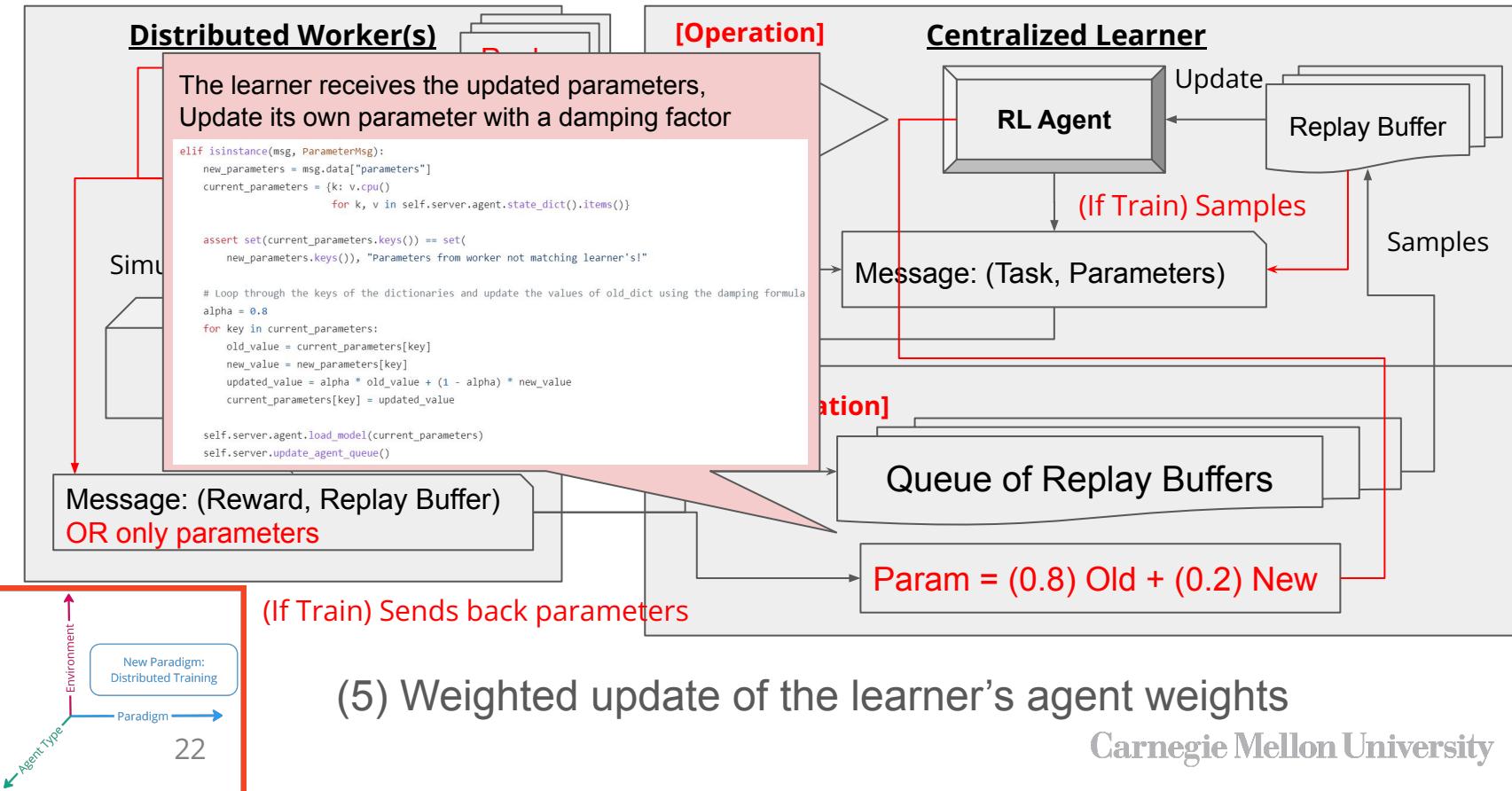
(3) Worker keeps local replay buffer for training local agent

Paradigm - Distributed Training / Update

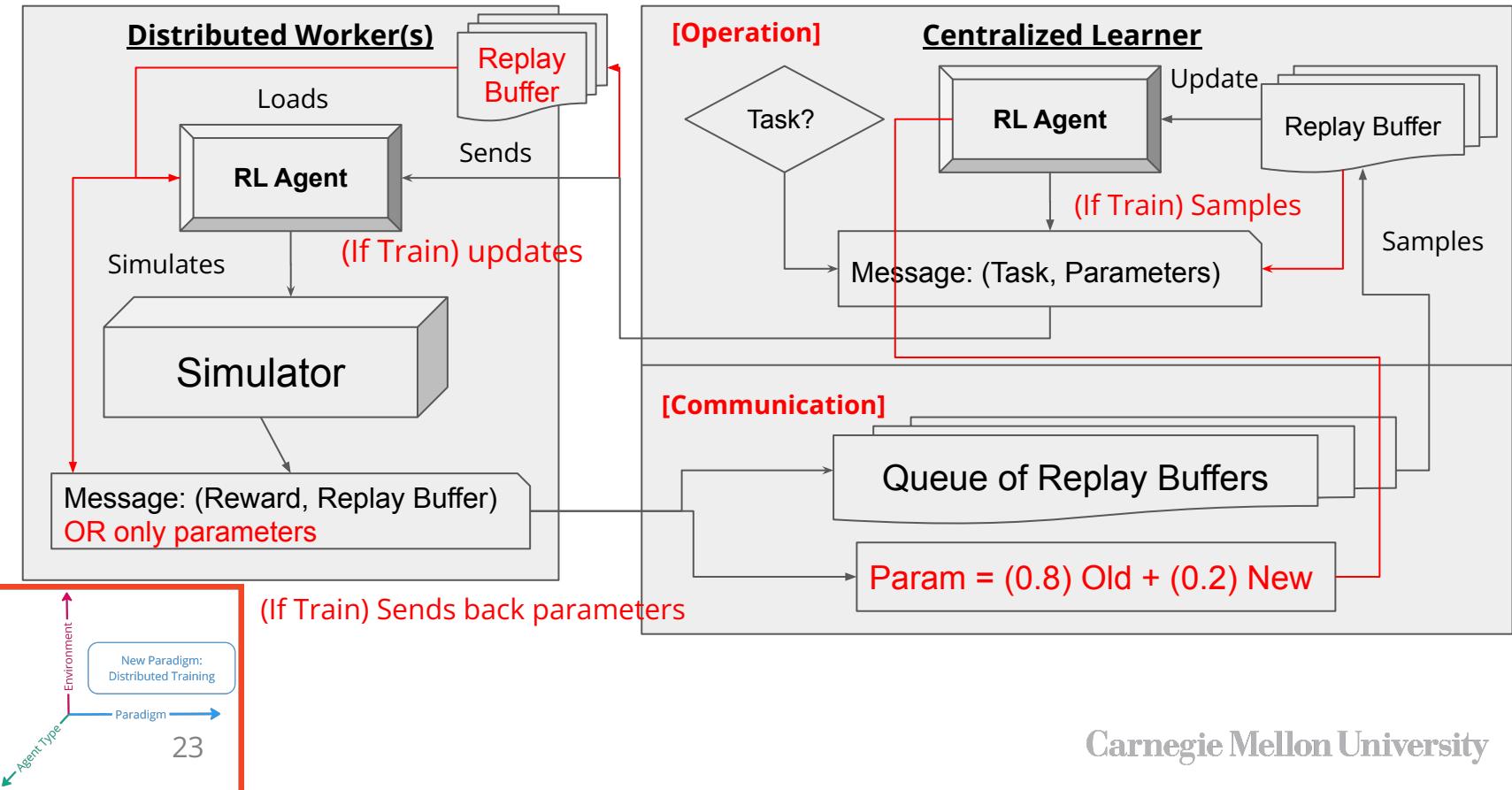


(4) Worker sends back updated parameters without simulation

Paradigm - Distributed Training / Update



Paradigm - Distributed Training / Update



Design - Agent and Environments

RL Agent Design - OpenAI's SAC

- Motivated by non-convergence of bipedal walker from existing approach
- Adopted OpenAI Spinning-Up [1] implementation of SAC agent and replay buffer

```
class MLPActorCritic():
    # Input: Encoded observation and action embedding
    # Output: Recommended action, Log prob. of action w.r.t policy distribution
    self.pi = (
        net_layer = MLP((obs_dim + act_dim) -> 256 -> 256),
        mu_layer = Linear(256 -> act_dim),
        log_std_layer = Linear(256 -> act_dim),
        pi_distribution = Normal(mu, log),
        pi_action = pi_distribution.sample(),
        logp_pi = pi_distribution.log_prob(pi_action)
    )

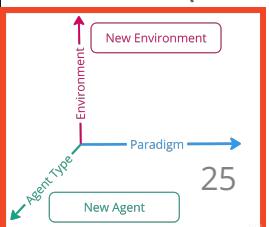
    # Input: Encoded observation and action embedding
    # Output: Scalar Q-value (quality of action in current state)
    self.q1 = self.q2 = MLP((obs_dim + act_dim) -> 256 -> 256 -> 1)
```

Policy-Network: learns to output the probability distribution of actions over current state, for action selection

Q-value Network: learns to output the quality of taking certain actions based on current state

miro

(Pseudocode)



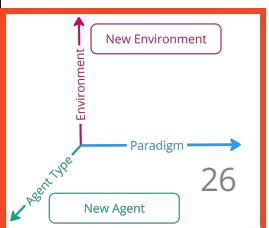
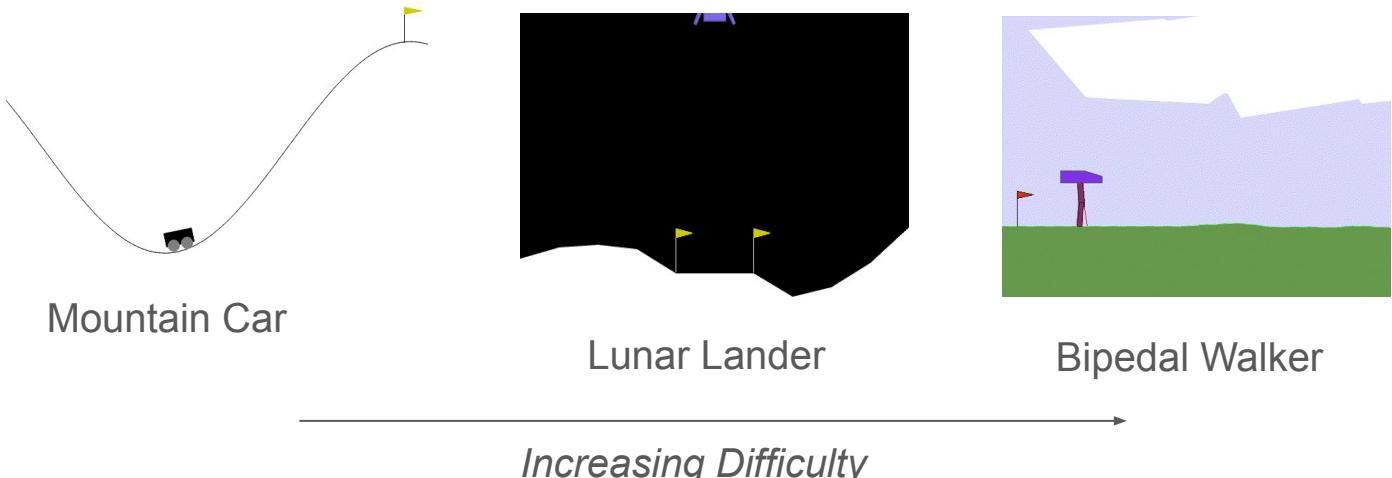
25

Two Q-networks for stability and robustness
(target Q-net periodically synced using weighted update of parameters)

[1] <https://github.com/openai/spinningup>

Environment Design - OpenAI's Gym

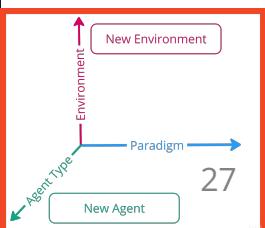
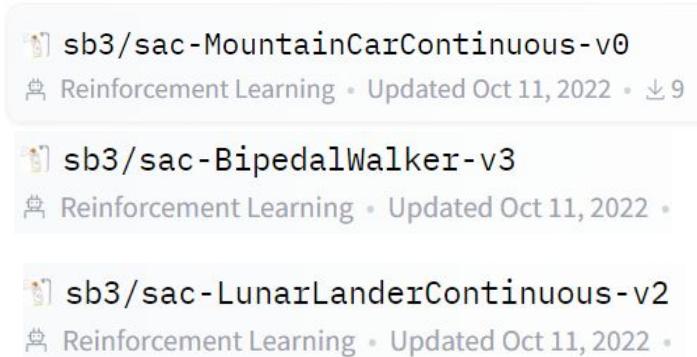
- Compatibility with OpenAI's Box-2d Continuous environments
- Newly added Lunar Lander environment
- Baselines for evaluating distributed paradigm



Parameters Design - Hugging Face

- Reinforcement Learning is prone to **hyperparameter initialization**
- For OpenAI's Gym environment's hyperparameter design
 - Stable-Baselines3 [1]
 - Spinning-Up [2]

Environment Id	Observation Space	Action Space	Reward Range	tStepL	Trials	rThresh
MountainCar-v0	Box(2,)	Discrete(3)	(-inf, inf)	200	100	-110.0
MountainCarContinuous-v0	Box(2,)	Box(1,)	(-inf, inf)	999	100	90.0
Pendulum-v0	Box(3,)	Box(1,)	(-inf, inf)	200	100	None
CartPole-v0	Box(4,)	Discrete(2)	(-inf, inf)	200	100	195.0



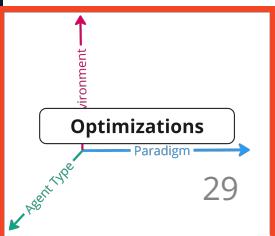
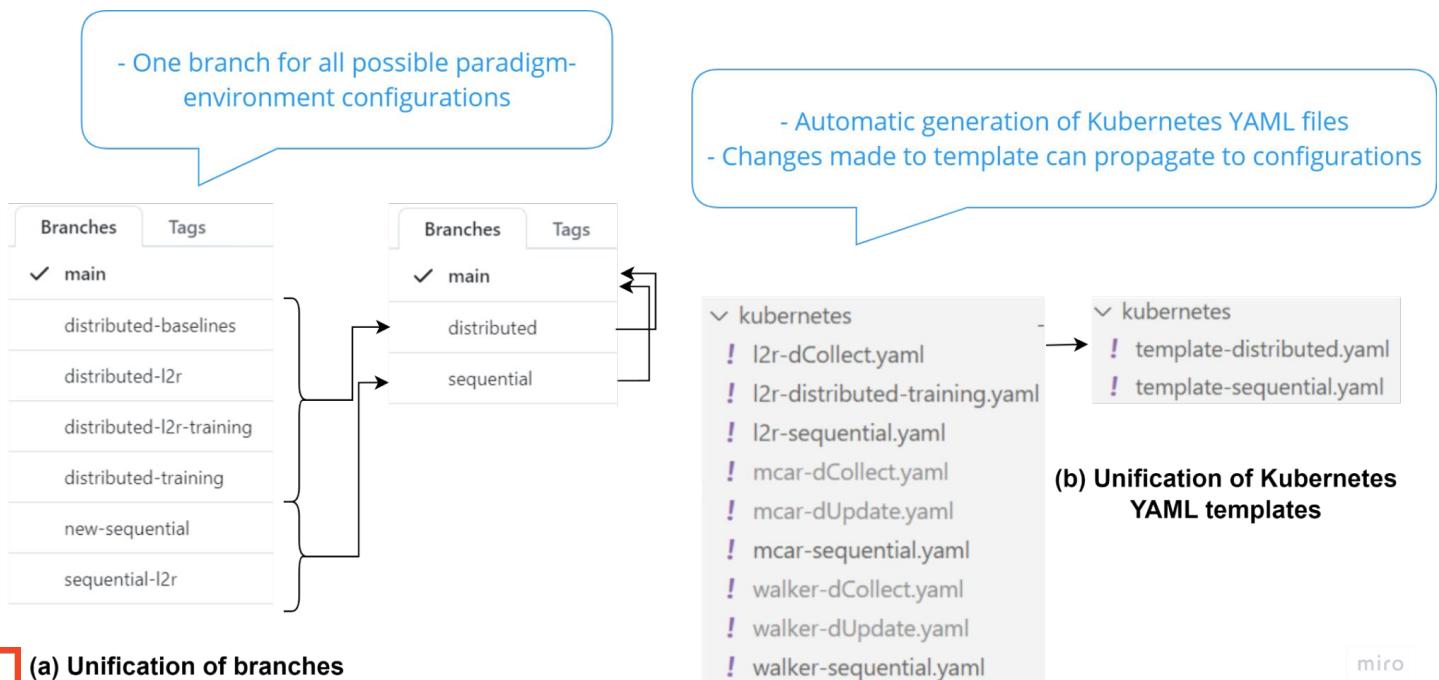
27

[1] <https://github.com/DLR-RM/stable-baselines3>[2] <https://github.com/openai/spinningup>

Design - Optimizations

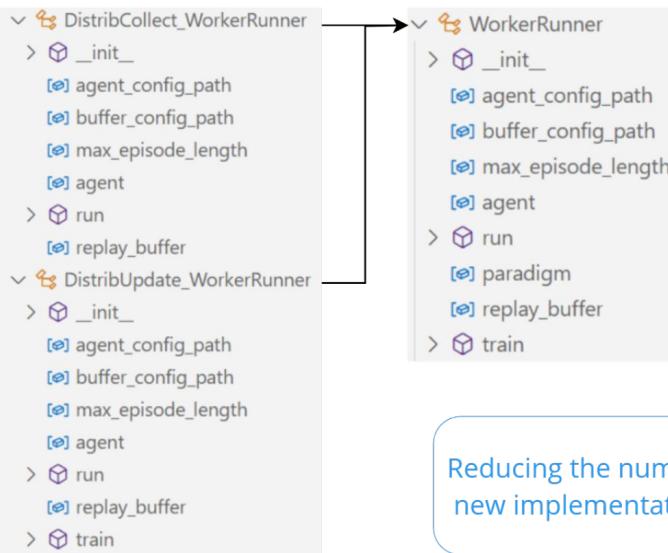
Design

Optimization - Unifications / Encapsulation / Modularity

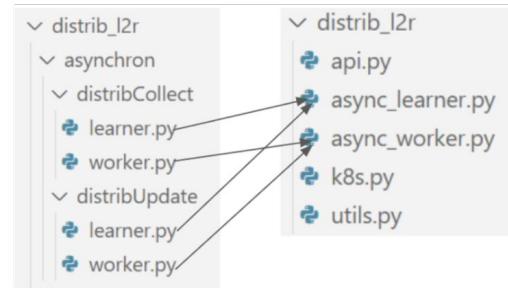


Optimization - Unifications / Encapsulation / Modularity

(c) Unification of Worker Runner class



(d) Unification of learner and worker definitions



Reducing the number of classes and files, allowing users to add new implementations conforming to the base class definitions

miro

Optimizations

Paradigm

Agent Type

30

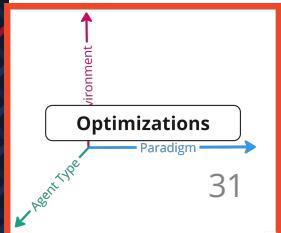
Optimization - Unifications / Encapsulation / Modularity

(e) Unification of worker configuration YAML



One config file per environment, instead of per environment-agent-paradigm configuration

miro



Optimization - Automatic Resource Control

- **Resource launching:** automatically create, store, and launch Kubernetes YAML files based on user's input configurations. Naming of resources descriptive and unique to avoid conflict.

```
sbian@phoebe-login:~/distributed$ python3 launch.py
Select RL environment (l2r/mcar/walker/walker-openai/lander-openai): lander-openai
Select training paradigm (sequential/dCollect/dUpdate): dUpdate
Input number of distributed workers: 3
Input WandB experiment name: lander-openai-dupdate-3workers
-----
RL Environment = [lander-openai] | Training Paradigm = [dUpdate] | Number of Workers = [3] | Experiment Name = [lander-openai-dupdate-3workers] ---
-----
[Warning] Port name exceeded 15 characters, renaming: lander-openai-dupdate -> lander-open-z9x
replicaset.apps/lander-openai-dupdate-workers created
pod/lander-openai-dupdate-learner created
service/lander-openai-dupdate-learner created
```

Resource Launching Script

- **Resource shutdown:** automatically locate and destroy the resources according to user input configuration.

```
sbian@phoebe-login:~/distributed$ python3 shutdown.py
Select RL environment (l2r/mcar/walker/walker-openai/lander-openai): lander-openai
Select training paradigm (sequential/dCollect/dUpdate): dUpdate
---
--- Shutting down [lander-openai] with training paradigm [dUpdate] ---
---
replicaset.apps "lander-openai-dupdate-workers" deleted
service "lander-openai-dupdate-learner" deleted
pod "lander-openai-dupdate-learner" deleted
```

Resource Shutdown

Optimizations

32

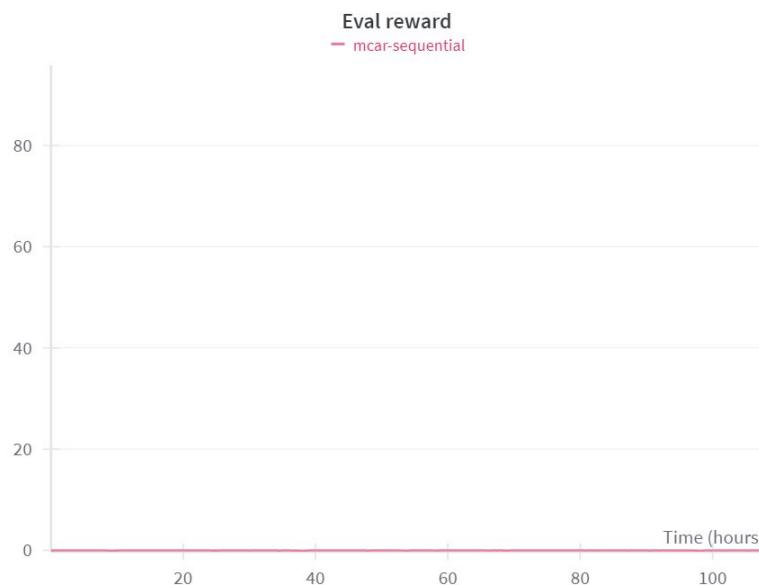
Results

Results - Overview

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car			
Bipedal Walker			
Bipedal Walker (OpenAI)			
Lunar Lander (OpenAI)			
Learn-to-Race			

Mountain Car + Sequential

- Convergence, after prolonged period of time



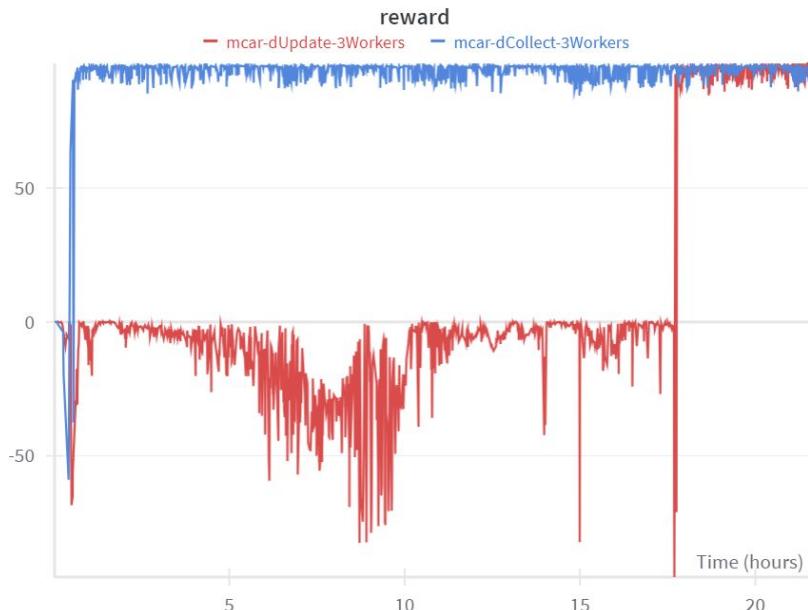
Results - Overview

- Legend: ✓ means convergence, ~ means non-convergence

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (very slow, ~100 h)		
Bipedal Walker			
Bipedal Walker (OpenAI)			
Lunar Lander (OpenAI)			
Learn-to-Race			

Mountain Car + Distributed Paradigms

- Faster convergence, distributed update paradigm works better with less workers (right)



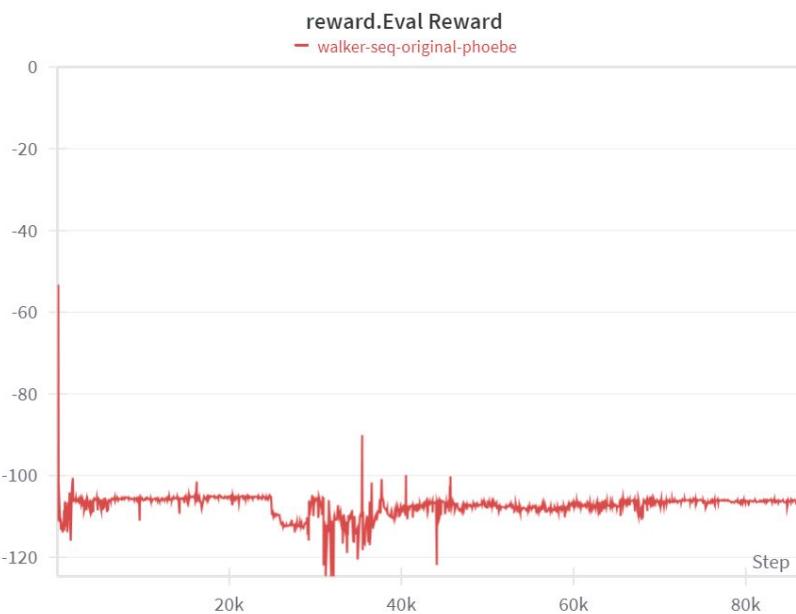
Results - Overview

- Legend: ✓ means convergence, ~ means non-convergence

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)
Bipedal Walker			
Bipedal Walker (OpenAI)			
Lunar Lander (OpenAI)			
Learn-to-Race			

Bipedal Walker + Sequential

- Non convergence (motivation to try out OpenAI's implementation)



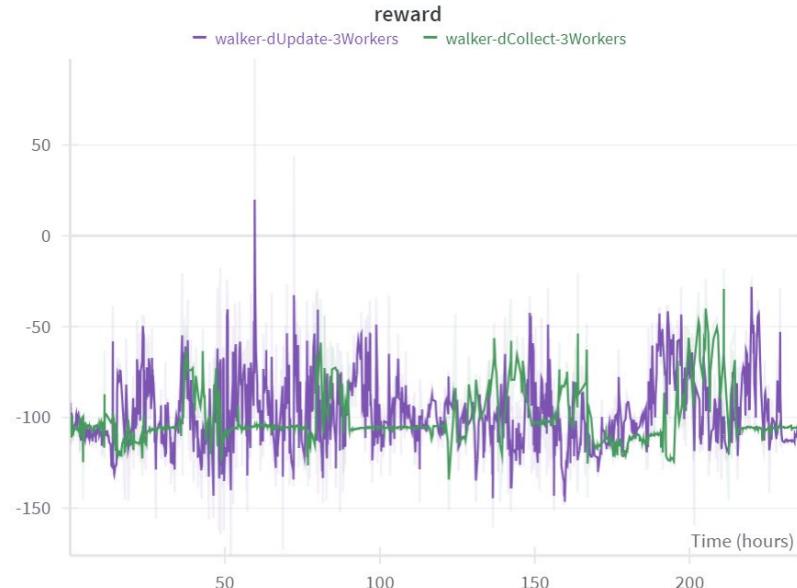
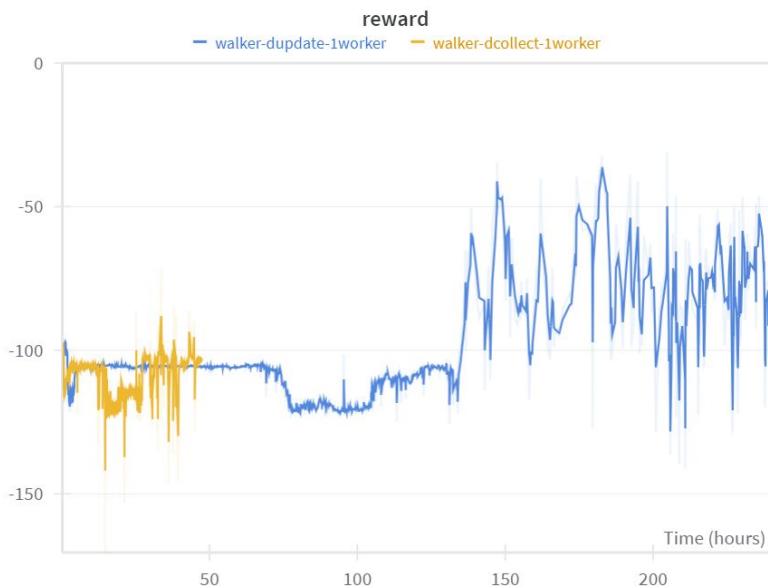
Results - Overview

- Legend: ✓ means convergence, ~ means non-convergence

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)
Bipedal Walker	~ (reward < 0)		
Bipedal Walker (OpenAI)			
Lunar Lander (OpenAI)			
Learn-to-Race			

Bipedal Walker + Distributed Paradigms

- Non convergence, patterns of periodical surges, with distributed update having more fluctuations and crossing 0



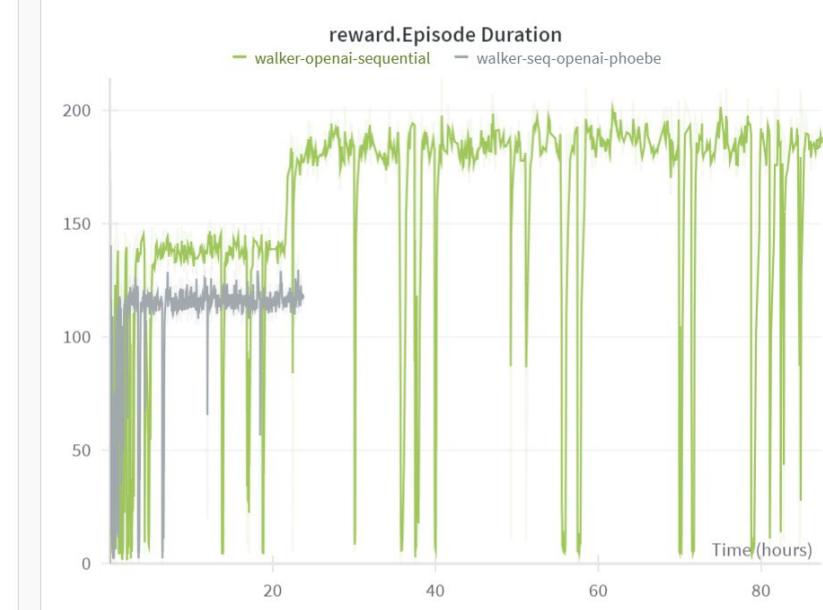
Results - Overview

- Legend: ✓ means convergence, ~ means non-convergence

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)
Bipedal Walker	~ (reward < 0)	~ (periodical surges)	~ (can cross 0, fluctuations)
Bipedal Walker (OpenAI)			
Lunar Lander (OpenAI)			
Learn-to-Race			

Bipedal Walker (OpenAI) + Sequential

- Rough convergence to 0, but not exceeding
- Episode duration has a positive correlation



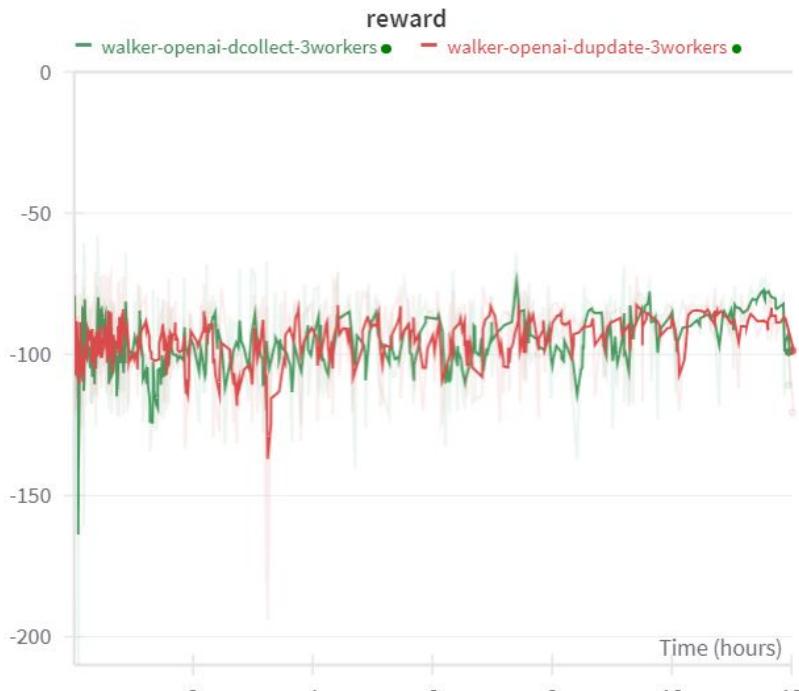
Results - Overview

- Legend: ✓ means convergence, ~ means non-convergence

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)
Bipedal Walker	~ (reward < 0)	~ (periodical surges)	~ (can cross 0, fluctuations)
Bipedal Walker (OpenAI)	~ (rough converge to reward near 0)		
Lunar Lander (OpenAI)			
Learn-to-Race			

Bipedal Walker (OpenAI) + Distributed Paradigms

- No clear convergence



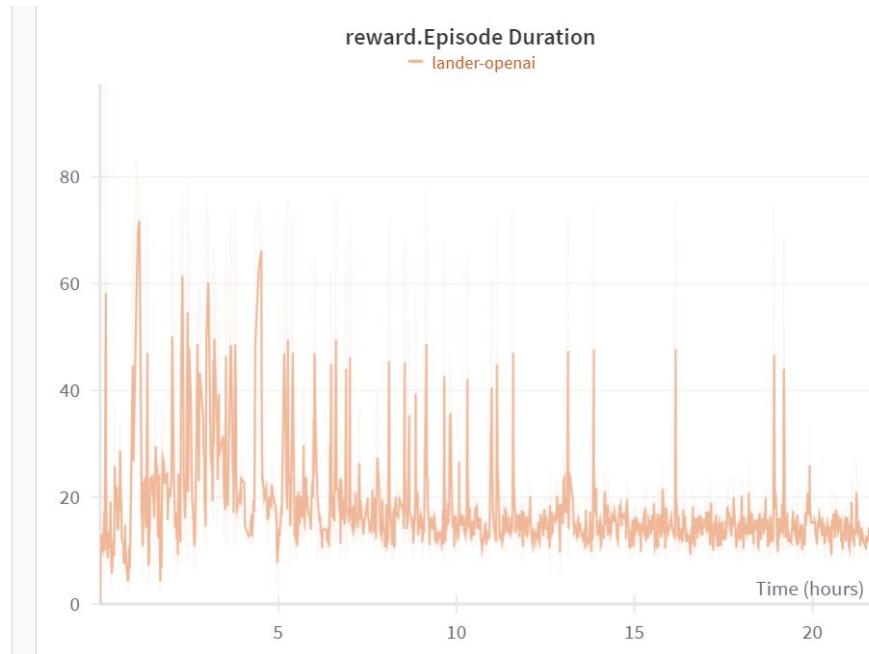
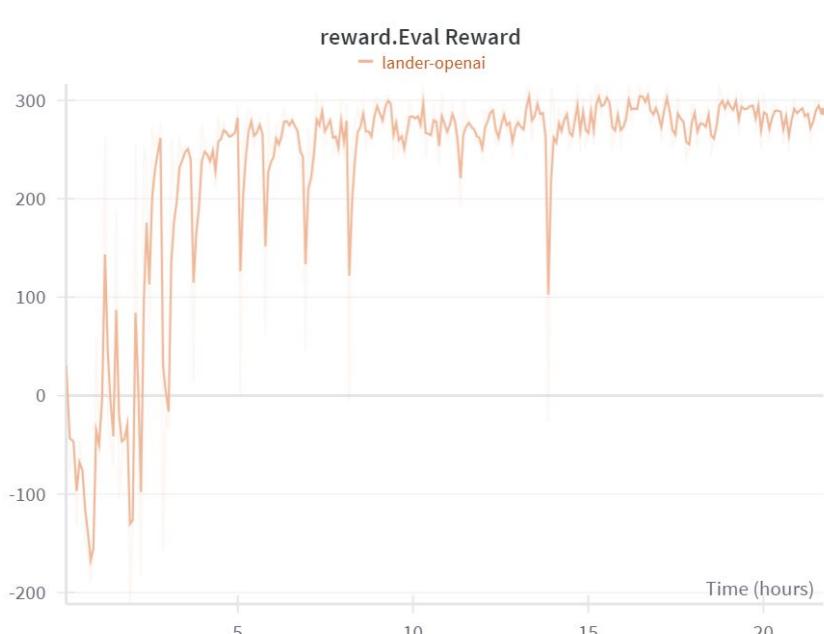
Results - Overview

- Legend: ✓ means convergence, ~ means non-convergence

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)
Bipedal Walker	~ (reward < 0)	~ (periodical surges)	~ (can cross 0, fluctuations)
Bipedal Walker (OpenAI)	~ (rough converge to reward near 0)	~ (no convergence)	~ (no convergence)
Lunar Lander (OpenAI)			
Learn-to-Race			

Lunar Lander (OpenAI) + Sequential

- Convergence, no sign of forgetting, so episode duration stays low



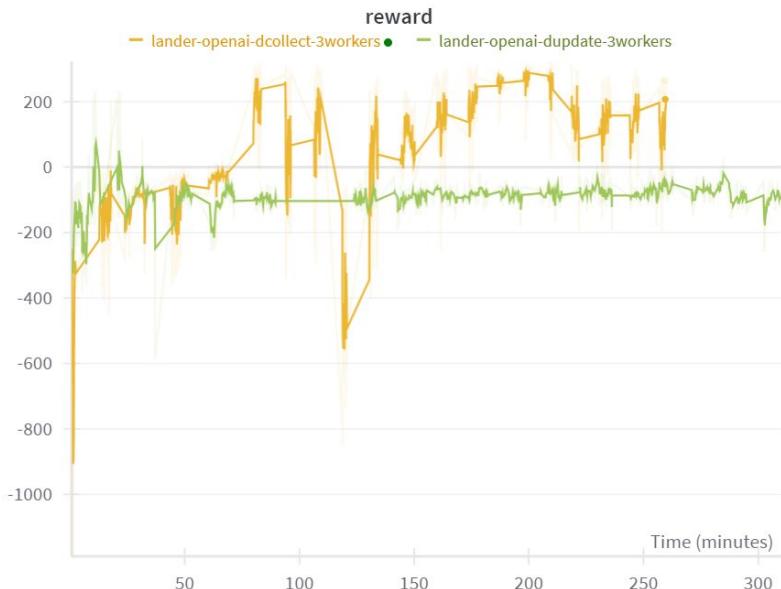
Results - Overview

- Legend: ✓ means convergence, ~ means non-convergence

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)
Bipedal Walker	~ (reward < 0)	~ (periodical surges)	~ (can cross 0, fluctuations)
Bipedal Walker (OpenAI)	~ (rough converge to reward near 0)	~ (no convergence)	~ (no convergence)
Lunar Lander (OpenAI)	✓ (~10 h, no sign of forgetting)		
Learn-to-Race			

Lunar Lander (OpenAI) + Distributed Paradigms

- Able to reach high reward fast (diverse environment?), but then degradation (forgetting?) - may benefit from checkpointing



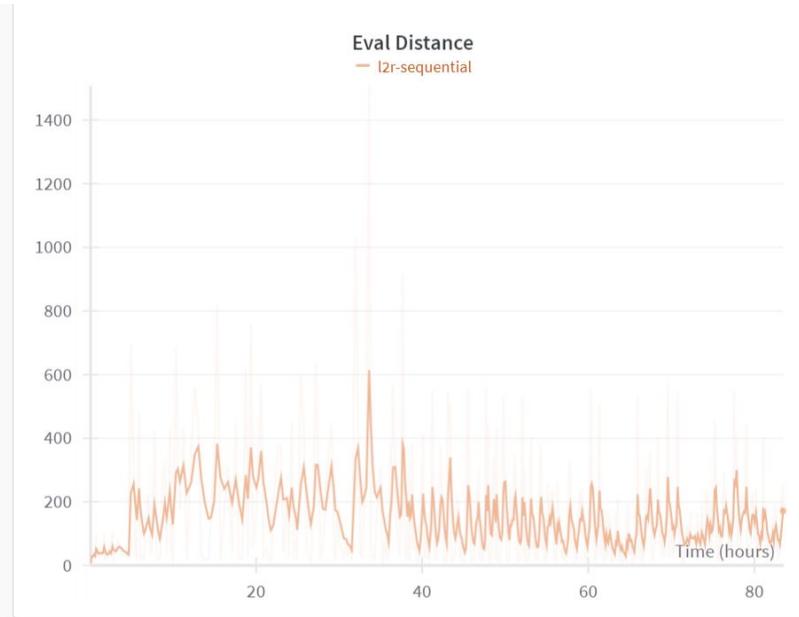
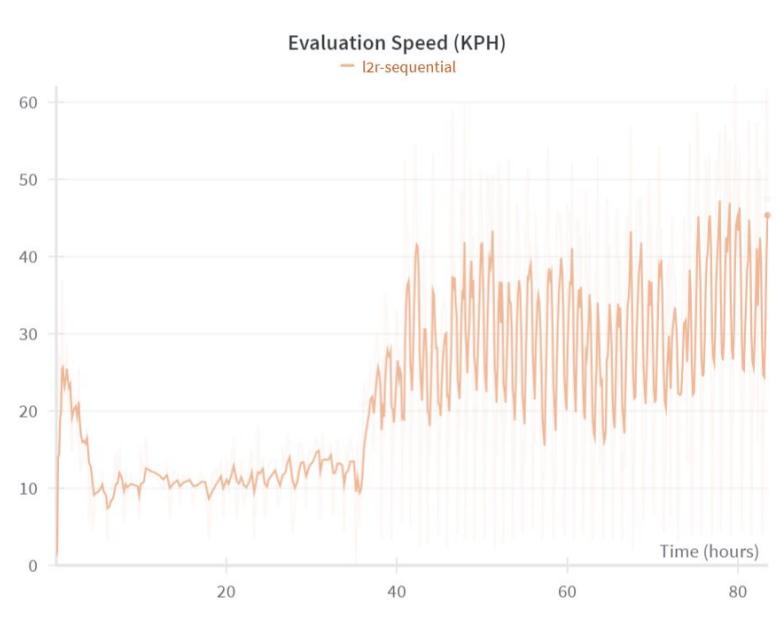
Results - Overview

- Legend: ✓ means convergence, ~ means non-convergence

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)
Bipedal Walker	~ (reward < 0)	~ (periodical surges)	~ (can cross 0, fluctuations)
Bipedal Walker (OpenAI)	~ (rough converge to reward near 0)	~ (no convergence)	~ (no convergence)
Lunar Lander (OpenAI)	✓ (~10 h, no sign of degradation)	✓ (reach high reward <2h, with <i>slight degradation</i>)	~ (reach high reward <1h, but then degradation)
Learn-to-Race			

Learn-to-Race + Sequential

- Using a revised base Docker image
- Slow convergence [1] of car traveling speed



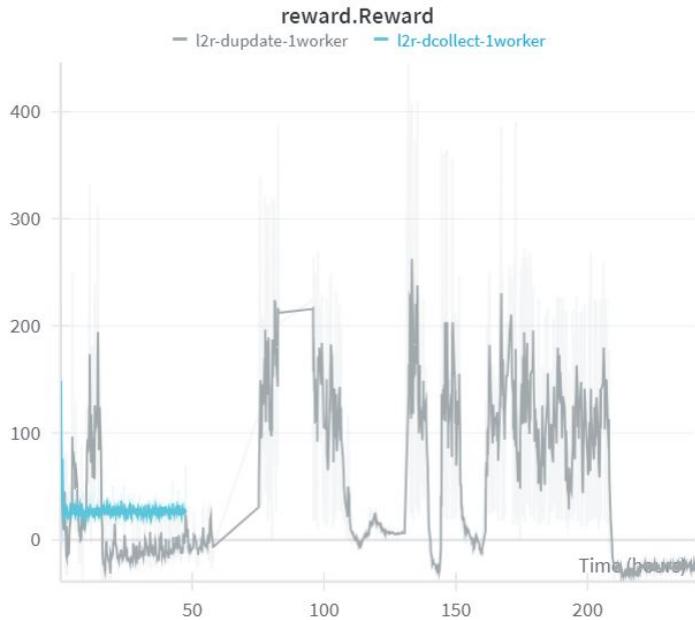
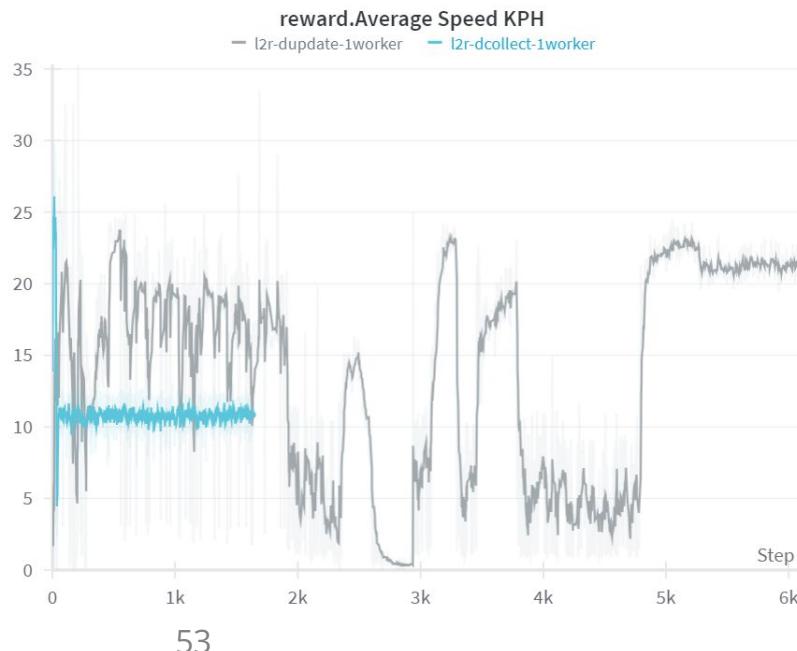
Results - Overview

- Legend: ✓ means convergence, ~ means non-convergence

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)
Bipedal Walker	~ (reward < 0)	~ (periodical surges)	~ (can cross 0, fluctuations)
Bipedal Walker (OpenAI)	~ (rough converge to reward near 0)	~ (no convergence)	~ (no convergence)
Lunar Lander (OpenAI)	✓ (~10 h, no sign of degradation)	✓ (reach high reward <2h, with <i>slight degradation</i>)	~ (reach high reward <1h, but then degradation)
Learn-to-Race	✓ (speed KPH)		

Learn-to-Race + Distributed Paradigms

- Distributed collection paradigm started high but converged to low
- Distributed update paradigm periodically achieves better performance



Results - Overview

- Legend: ✓ means convergence, ~ means non-convergence

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)
Bipedal Walker	~ (reward < 0)	~ (periodical surges)	~ (can cross 0, fluctuations)
Bipedal Walker (OpenAI)	~ (rough converge to reward near 0)	~ (no convergence)	~ (no convergence)
Lunar Lander (OpenAI)	✓ (~10 h, no sign of degradation)	✓ (reach high reward <2h, with <i>slight degradation</i>)	~ (reach high reward <1h, but then degradation)
Learn-to-Race	✓ (speed KPH)	~ (stay at low performance)	~ (periodic improvement)

Analysis

Analysis

- Let's divide-and-conquer this table

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)
Bipedal Walker	~ (reward < 0)	~ (periodical surges)	~ (can cross 0, fluctuations)
Bipedal Walker (OpenAI)	~ (rough converge to reward near 0)	~ (no convergence)	~ (no convergence)
Lunar Lander (OpenAI)	✓ (~10 h, no sign of degradation)	✓ (reach high reward <2h, with <i>slight degradation</i>)	~ (reach high reward <1h, but then degradation)
Learn-to-Race	✓ (speed KPH)	~ (stay at low performance)	~ (periodic improvement)

Analysis

RL Environments	Sequential	Distributed Collection	Distributed Update
Mountain Car 	✓ (slow, ~100 h)	✓ (faster with more workers)	✓ (faster with less workers)

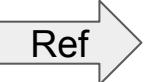
- **Simple environment:** all able to converge, no forgetting/degradation
- **Distributed paradigms:** achieves ~10x faster convergence
- **Distributed collection faster with more workers**
 - Fast training, learner not bottleneck, workers collect diverse experience
- **Distributed update slower with more workers**
 - Network, sharing of parameters, coordinations all become overhead

Analysis

RL Environments	Sequential	Distributed Collection	Distributed Update
Bipedal Walker (OpenAI) 	~ (rough converge to reward near 0)	~ (no convergence)	~ (no convergence)
Lunar Lander (OpenAI) 	✓ (~10 h, no sign of degradation)	~ (reach high reward <2h, but then degradation)	~ (reach high reward <1h, but then degradation)

- **OpenAI Vanilla baselines:** sequential version outperforms original gym baselines, but distributed baseline behavior is not stable
- **Distributed Collection (dCollect) v.s. Distributed Update (dUpdate)**
 - Both performs similarly, with unsatisfactory results for bipedal walker, and the typical “learn-fast, forget-fast” issue for lunar lander

Analysis

RL Environments	Distributed Collection	Distributed Update
Bipedal Walker 	~ (periodical surges)	~ (can cross 0, fluctuations)
Learn-to-Race 	~ (converge to low performance)	~ (periodic improvement)

- **Harder environments:** fluctuations due to diverse environments, potential to reach high performance faster, but suffering from degradation (i.e., learning fast but forgetting fast)
- **Distributed Collection (dCollect) v.s. Distributed Update (dUpdate)**
 - dUpdate induces more fluctuations due to distributed update of parameters (damping factor)
 - dUpdate (due to high fluctuations) generally has more potential to reach better rewards
 - dUpdate can benefit from checkpointing (greedily saving best checkpoint, if environment is static)

Conclusion

Summary of Analysis

RL Environments	Sequential Paradigm	Distributed Paradigms
Easy Gym (mountain car)	✓ stable convergence ✗ very slow	✓ stable convergence, much faster ✗ communication overhead
Hard Gym (bipedal walker)	✓ if convergence, should stay stable ✗ often stuck at worse reward	✓ reach high reward fast (learn fast) due to diverse simulations ✗ suffer from degradation (forget fast)
Vanilla SAC	✓ better convergence performance than non-Vanilla	✗ generally poor behavior (need to investigate further)
Learn-to-Race	✓ observed stable improvement ✗ often stuck at worse reward	✓ observed periodical improvement, faster in comparison ✗ fluctuations, unstable

Conclusion

- **Contribution**
 - o New training paradigm that distribute training tasks to workers
 - o Codebase and resource control optimizations
 - o Experiments and analysis of diverse configuration combinations
- **Observations**
 - o *Handling RL and distributed system design is challenging*
 - o For easy tasks like mountain car, distributed paradigms tend to converge faster
 - o For medium tasks like lunar lander and bipedal walker, distributed paradigms with diverse environments fluctuate and can thus reach high reward pretty quickly, but suffers from forgetting and degradation → may be mitigated through checkpointing
 - o For Learn-to-Race task, distributed Update paradigm has periodical improvements
- **Future Work**
 - o Torch device handling (store replay buffer on CPU, Gym action must be numpy)
 - o Importance sampling and dynamically adjusted task allocations
 - o Should take other thrusts into consideration

Appendix

Terminology

- **Reinforcement learning (RL)**: a machine learning training procedure in which an agent learns how to take actions within an environment based on reward and penalty feedback.
- **Agents**: an intelligent artificial entity can perceive its environment and respond reasonably.
- **State**: the parameters that define the status of the agent.
- **Action**: the action taken by the agent based on its environment and state to maximize the reward function.
- **Environment**: the environmental feedback that the agent is within.
- **Reward function**: a formulated representation indicating the gain of the agent at a certain state and environment through taking a certain action.

References

1. Xingjian SHI, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-kin Wong, and Wang-chun WOO. Convolutional Istm network: A machine learning approach for precipitation nowcasting. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 28. Curran Associates, Inc., 2015.
2. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017
3. Weights and Biases: <https://wandb.ai/>
4. Wang, Limin, et al. "Tdn: Temporal difference networks for efficient action recognition." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.
5. T. Weiss and M. Behl, "DeepRacing: A Framework for Autonomous Racing," 2020 Design, Automation & Test in Europe Conference & Exhibition (DATE), Grenoble, France, 2020, pp. 1163-1168, doi: 10.23919/DATE48585.2020.9116486.

Reflection

- **Documentations**
 - **Vision:** define the project scope, goals, and objectives
 - **Requirement:** identify and prioritize project requirements within our scope
 - **Design:** detailed understanding of modularized dependencies and risks
 - **Plan:** project timelines, milestones, and syncing between thrusts
- **Change of perception**
 - Inter-thrust communication
 - Communication with teaching staff
- **Working as a group**
 - Fixed meeting schedules
 - Fixed allocation of work and responsibilities

Changes To Previous Deliverables

- We have changed to use the new docker image jingyua to substitute l2r2022 since it seemed to provide better performance for the learn-to-race environment. The results, reasoning, and analysis are included in the previous slides.
- We have made significant changes to the previous github repositories to reduce the number of branches to two (sequential and distributed paradigm), which was discussed in the previous slides.