# Real-DRL: Teach and Learn in Reality

Yanbing Mao[1,*], Yihao Cai[1,*], Lui Sha[2]

[1]Wayne State University

[2]University of Illinois Urbana-Champaign

*Indicates Equal Contribution

**Presenter:** Yihao Cai

# Preliminary

# Motivation

## Runtime Safety for Deep Reinforcement Learning (DRL)



**Autonomous Vehicles[1]**

**Unmanned Aircraft[2]**

**Quadruped Robots[3]**

**Humanoid Robots[4]**

## How do we ensure runtime safety in safety-critical autonomous systems while DRL agents perform online learning?

*Reference:*
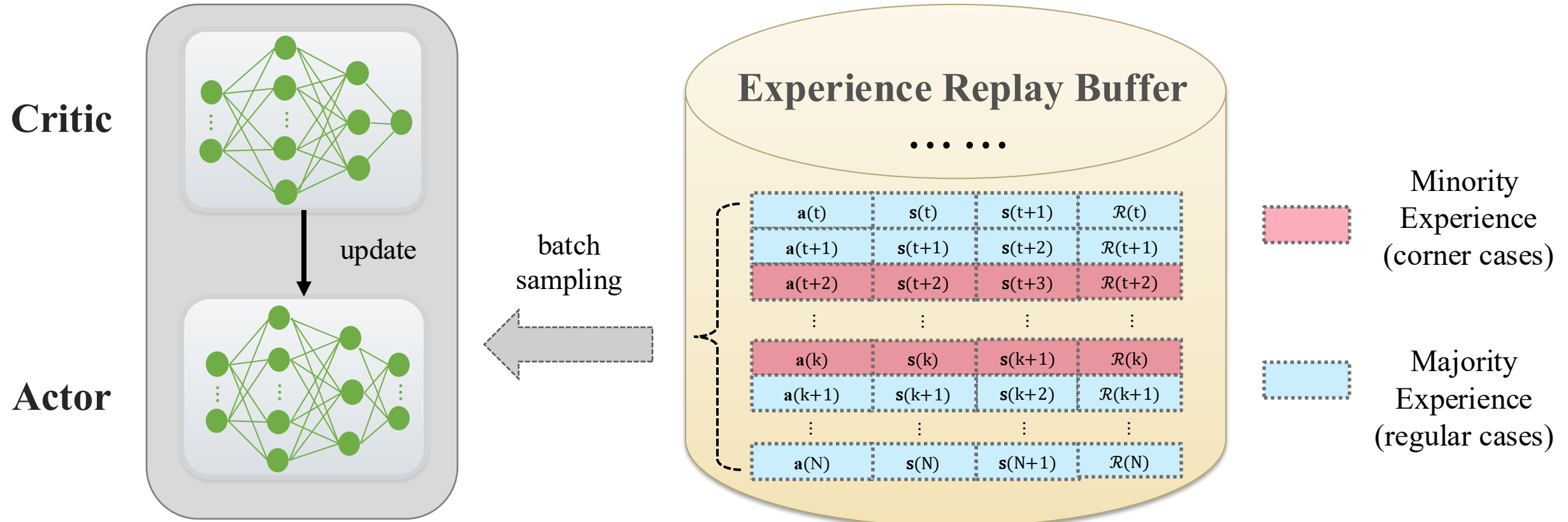*[1] https://www.wired.com/story/dashcam-footage-shows-driverless-cars-cruise-waymo-clogging-san-Francisco/*
*[2] https://flyfrompti.com/unmanned-aircraft-systems-uas-drones/*
*[3] https://droneblocks.io/product/go2-edu-quadruped-robot/?srsltid=AfmBOoqbUHBaaWUpBTC0kkCZOT4tc_DKzTiHbY6uM4-DF36bHmMejDqA*
*[4] https://manlybattery.com/guide-to-leading-humanoid-robots/?srsltid=AfmBOoo1P5Dza-0L1jEdroApnsv2Um_yD2Wxozw_w1V-tYzqF2XObhkJ*

# Motivation

## Data Imbalance Issue from Sampling

**How can the challenges of data imbalance be tackled to achieve more robust and generalizable DRL policies?**

# Challenges

- **Runtime Learning Safety**

  ➤ The risky nature of <u>trial-and-error</u> exploration in DRL

  ➤ Learning in <u>hard-to-predict</u> and <u>hard-to-simulate</u> environments requires timely and adaptive responses

- **Safety-related Data Imbalance Issue**

  ➤ Underrepresentation of <u>rare</u> but <u>crucial</u> data → poor safety at critical moments

  ➤ Leading to training bias and <u>limited generalization</u> capability
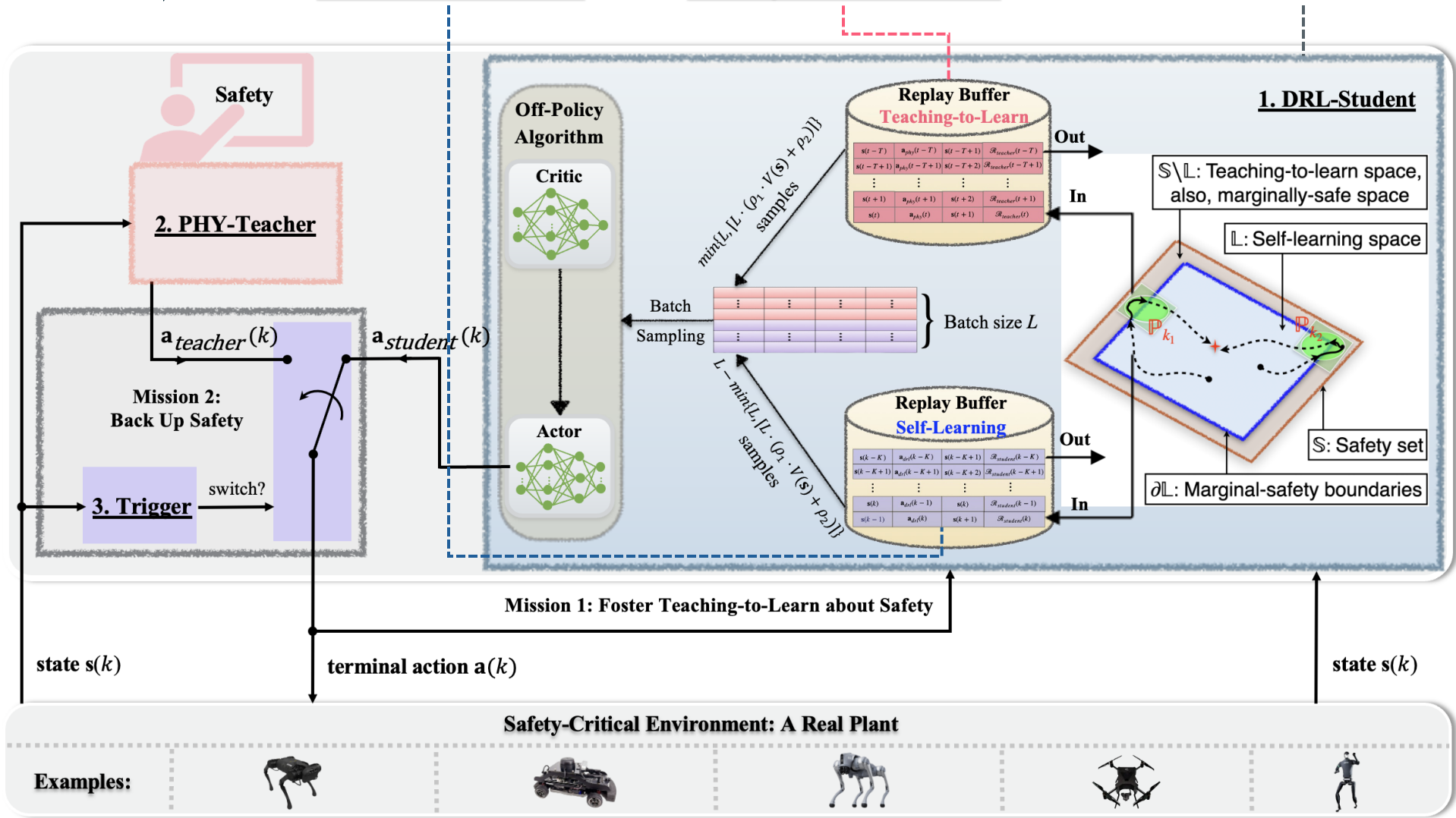
- **Sampling Efficiency**

  ➤ <u>High-quality</u> data fosters efficient and safe learning

  ➤ Inefficient sampling <u>prolongs training</u>, and increases runtime <u>safety risks</u>

# Proposed Framework



Examples:

Component 1: DRL-Student
1. Dual buffer for **self-learning** and **teaching-to-learn** paradigm
2. Safety-informed batch sampling

Component 2: PHY-Teacher
1. Fostering the **teaching-to-learn** mechanism regarding safety
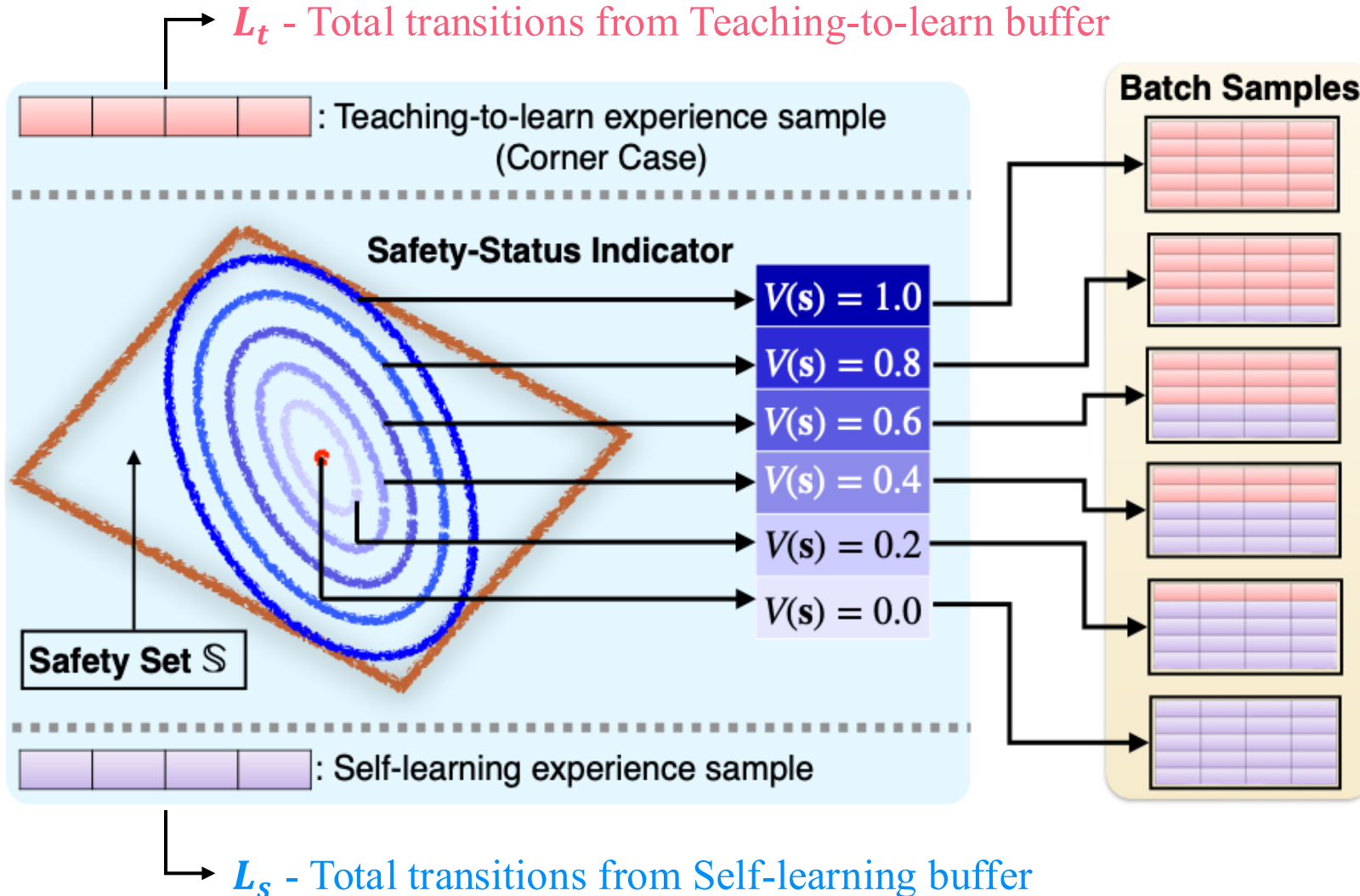2. Safety backup for the real plants

Component 3: Trigger
Monitoring the real-time safety status of the physical plant, and also deciding the terminal action to the plant

# Component-I: DRL-Student

## Safety-informed Batch Sampling

$L_t$ - Total transitions from Teaching-to-learn buffer

: Teaching-to-learn experience sample (Corner Case)

**Safety-Status Indicator**

$V(\mathbf{s}) = 1.0$

$V(\mathbf{s}) = 0.8$

$V(\mathbf{s}) = 0.6$

$V(\mathbf{s}) = 0.4$

$V(\mathbf{s}) = 0.2$

$V(\mathbf{s}) = 0.0$

Safety Set $\mathbb{S}$

: Self-learning experience sample

$L_s$ - Total transitions from Self-learning buffer

**Batch Samples**

**Safety-status indicator**

$$V(\mathbf{s}) \triangleq \mathbf{s}^T \cdot \mathbf{P} \cdot \mathbf{s}$$

$\mathbf{s}$ − real time state $\mathbf{s}(t)$

**Total Sampled Batch Size**

$$L = L_t + L_s$$

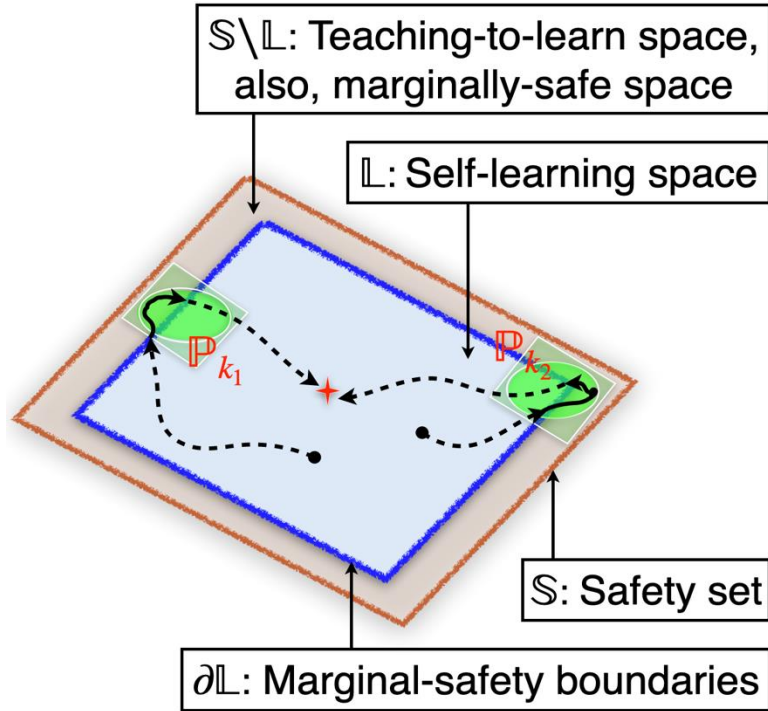$$L_t = \min\{L, \lceil L \cdot (\rho_1 \cdot V(\mathbf{s}(t)) + \rho_2 \rceil\}$$

$$L_s = L - \min\{L, \lceil L \cdot (\rho_1 \cdot V(\mathbf{s}(t)) + \rho_2 \rceil\}$$

$\rho_1, \rho_2$ − hyperparameters

# Component-II: PHY-Teacher

## Real-time Patch Design



$\mathbb{S} \setminus \mathbb{L}$: Teaching-to-learn space, also, marginally-safe space

$\mathbb{L}$: Self-learning space

$\mathbb{P}_{k_1}$  $\mathbb{P}_{k_2}$

$\mathbb{S}$: Safety set

$\partial\mathbb{L}$: Marginal-safety boundaries

**System Dynamics**: $s(t+1) = f(s(t), \mathbf{a}(t))$, $t \in \mathbb{N}$

**Safety Set**: $\mathbb{S} \triangleq \{s \in \mathbb{R}^n | -\mathbf{c} < \mathbf{C} \cdot s < \mathbf{c}\}$

**Self-Learning Space**: $\mathbb{L} \triangleq \{s \in \mathbb{R}^n | -\eta \cdot \mathbf{c} < \mathbf{C} \cdot s < \eta \cdot \mathbf{c}\}$, $0 < \eta < 1$

**Real-time Patch Design $\mathbb{P}_k$**

① **Control Goal** : $s_k^* \triangleq \chi \cdot s(k)$,  $s(k\text{-}1) \in \mathbb{L}$ and $s(k) \in \partial\mathbb{L}$

② **LMI Feasibility** : Construct LMIs and optimize for $\boldsymbol{F}_k$

③ **Action Policy** : $\mathbf{a}_{teacher}(t) = \boldsymbol{F}_k \cdot (s(t) - s_k^*)$,  $t \in \mathbb{T}_k$

# Component-III: Trigger

<u>Triggering Condition</u> $\mathcal{T}$ : $s(k\text{-}1) \in \mathbb{L}$ and $\mathbf{s}(k) \in \partial\mathbb{L}$

<u>PHY-Teacher Action Step</u>: $\mathbb{T}_k \triangleq \{k+1, k+2, \ldots, k+\tau_k\}$

**Switching Law**

$$\mathbf{a}(t) = \begin{cases} \mathbf{a}_{teacher}(t), & \text{if } \mathcal{T} \text{ holds at } k, \text{ and } t \in \mathbb{T}_k \\ \mathbf{a}_{student}(t), & \text{if } s(k) \in \mathbb{L} \end{cases}$$

NEURAL INFORMATION PROCESSING SYSTEMS
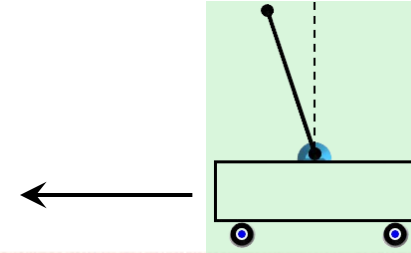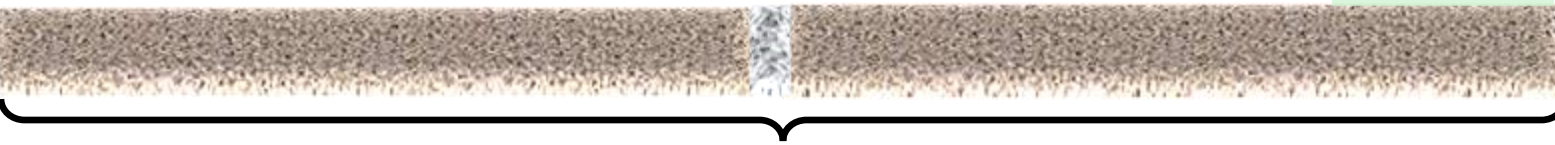
# Experiment

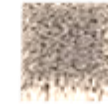# Experiment-I: Cartpole

## Environmental Setup



Testing Environment

Runtime Learning Environment

Friction coefficient of cart-road: 40
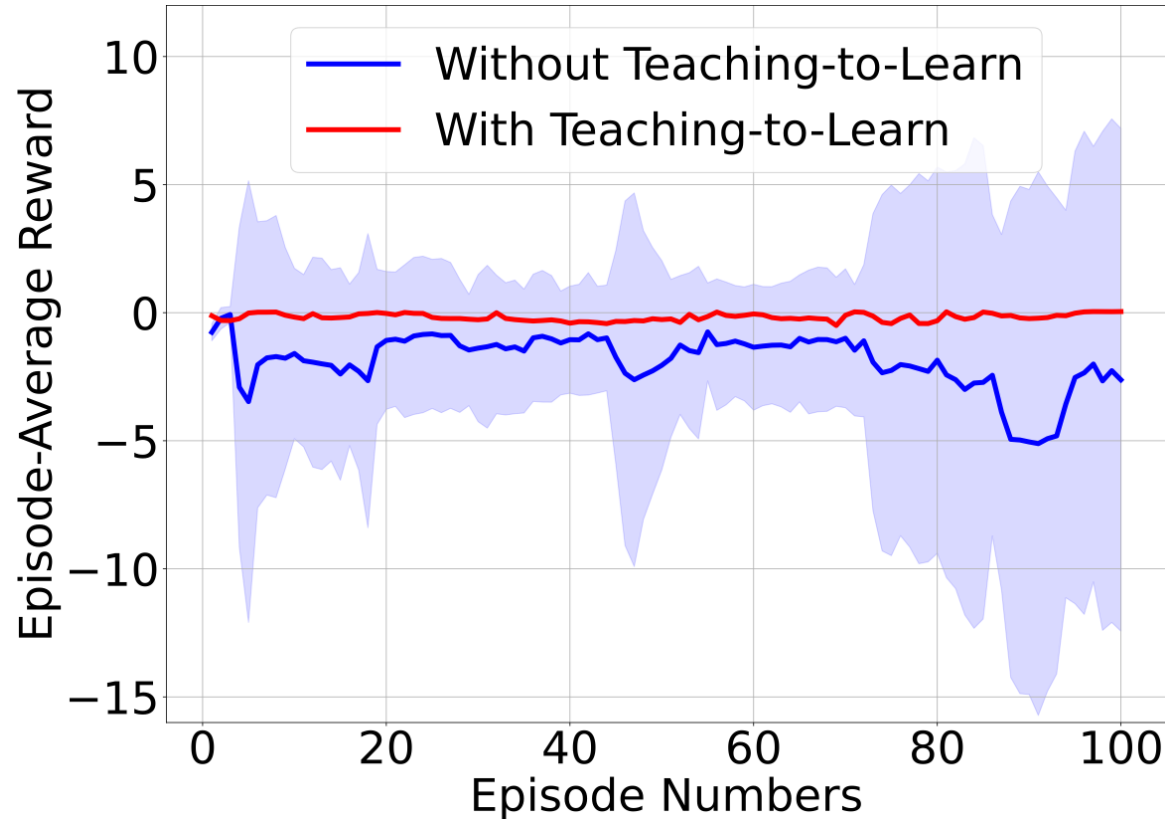
Friction coefficient of cart-road: 10

## Ablation Study – Demonstrating Three Key Features of Real-DRL

**Three Key Features**

- ○ **Feature I:** Teaching-to-Learn Mechanism
- ○ **Feature II:** Safety-informed Batch Sampling
- ○ **Feature III:** Automatic Hierarchical Learning

# Experiment-I: Cartpole

**Feature I:** Teaching-to-Learn mechanism



Episode-Average Reward

**Episode-Average Reward:**

$$\frac{\text{Return (i.e., cumulative reward) in one episode}}{\text{DRL-Student's total activation time in one episode}}$$
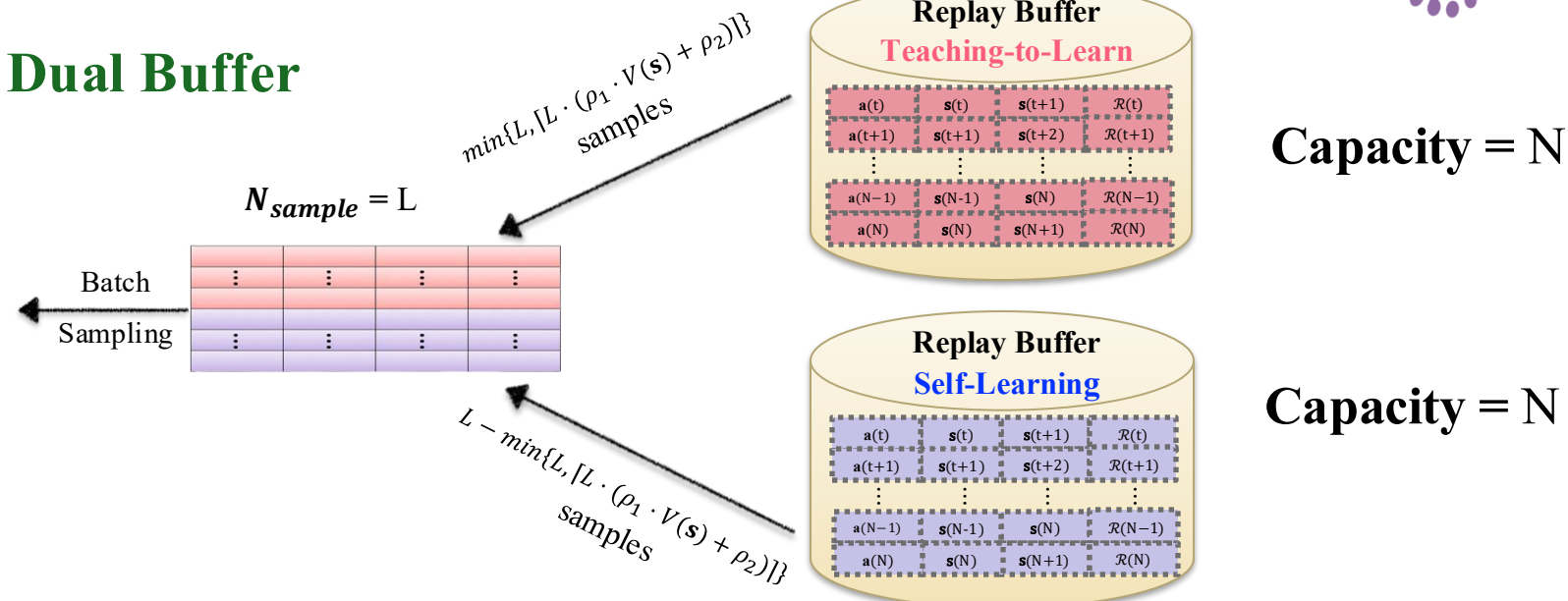
Adopting Teaching-to-Learn paradigm leads to overall improved episode-average reward and stable learning
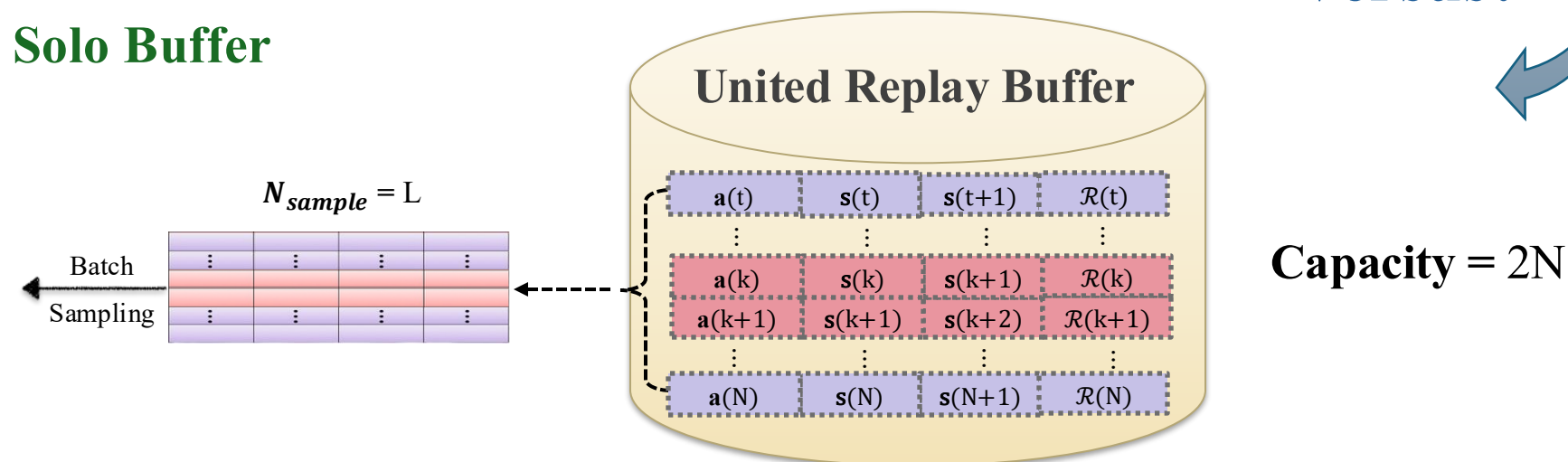
# Experiment-I: Cartpole
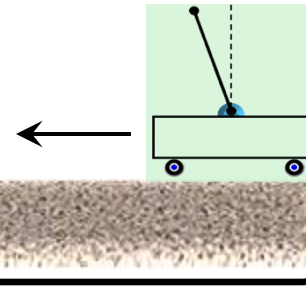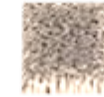
## Feature II: Safety-informed Batch Sampling

# Experiment-I: Cartpole

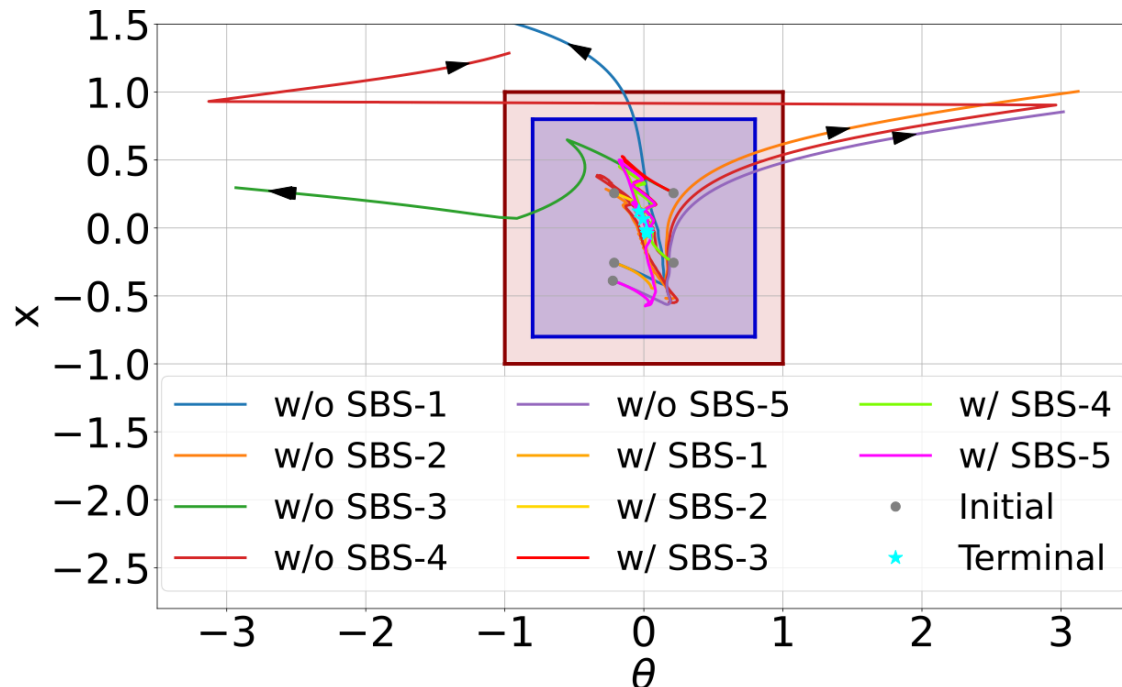## Feature II: Safety-informed Batch Sampling



**Testing Environment**

**Runtime Learning Environment**

Friction coefficient = 40

Friction coefficient = 10



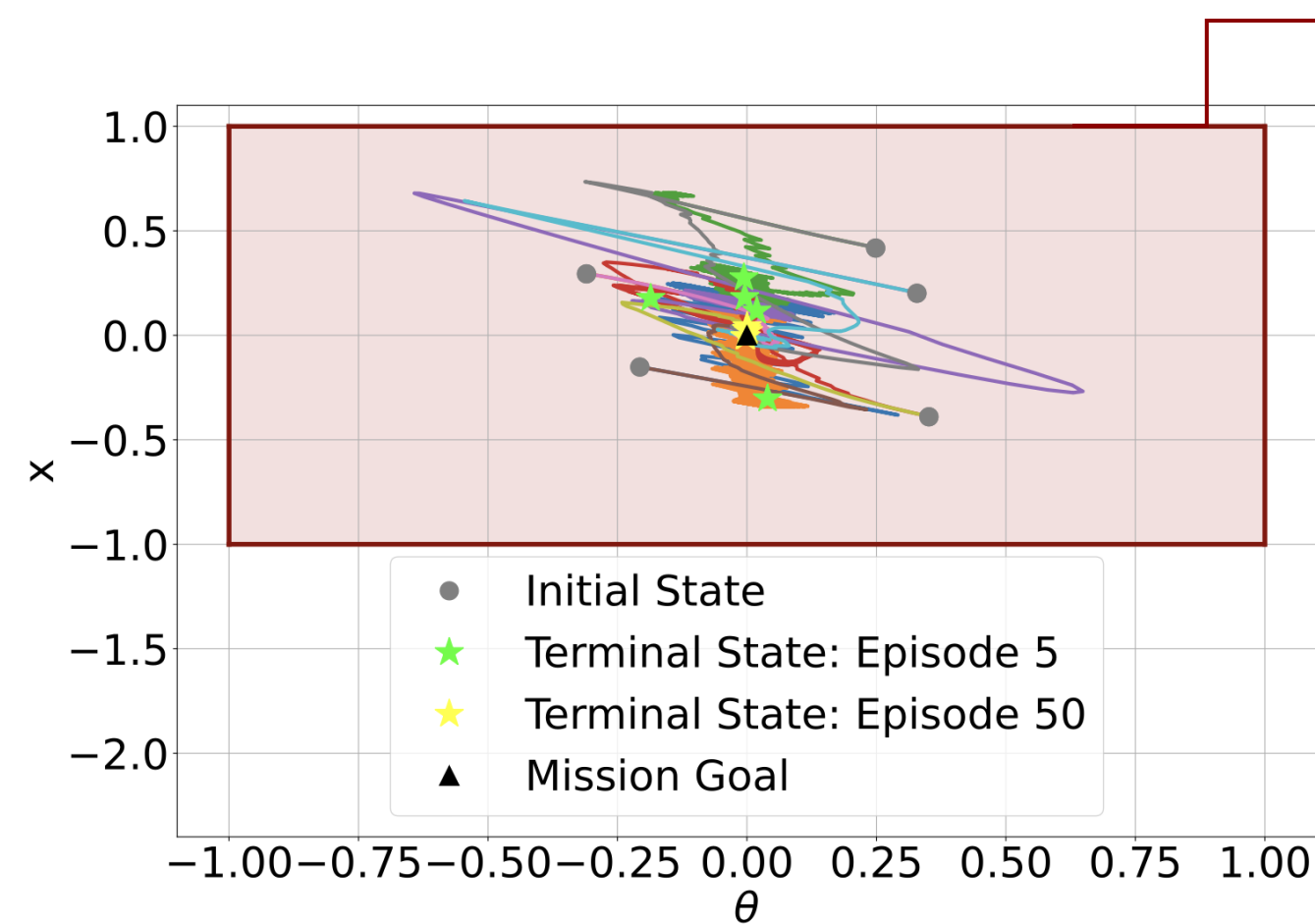Phase Plot (with vs. without safety-informed sampling)

1. Agent built on Real-DRL maintains safety on **both** cases after runtime learning

2. Agent sampling from a united replay buffer maintains safety in the **majority** cases but failed on **corner** cases

# Experiment-I: Cartpole

**Feature III:** Automatic Hierarchical Learning

Safety Set



Agent Trajectory from Different Episodes (Inference)

**Task Goal:** $(\bar{x}^*, \bar{\theta}^*) = (0, 0)$

From the same initial state, after **5** episodes learning by Real-DRL, the trajectory of the agent is within the safety set (**safety-first**); after **20** episodes, the trajectory gets closer to the control goal (**high-performance**)
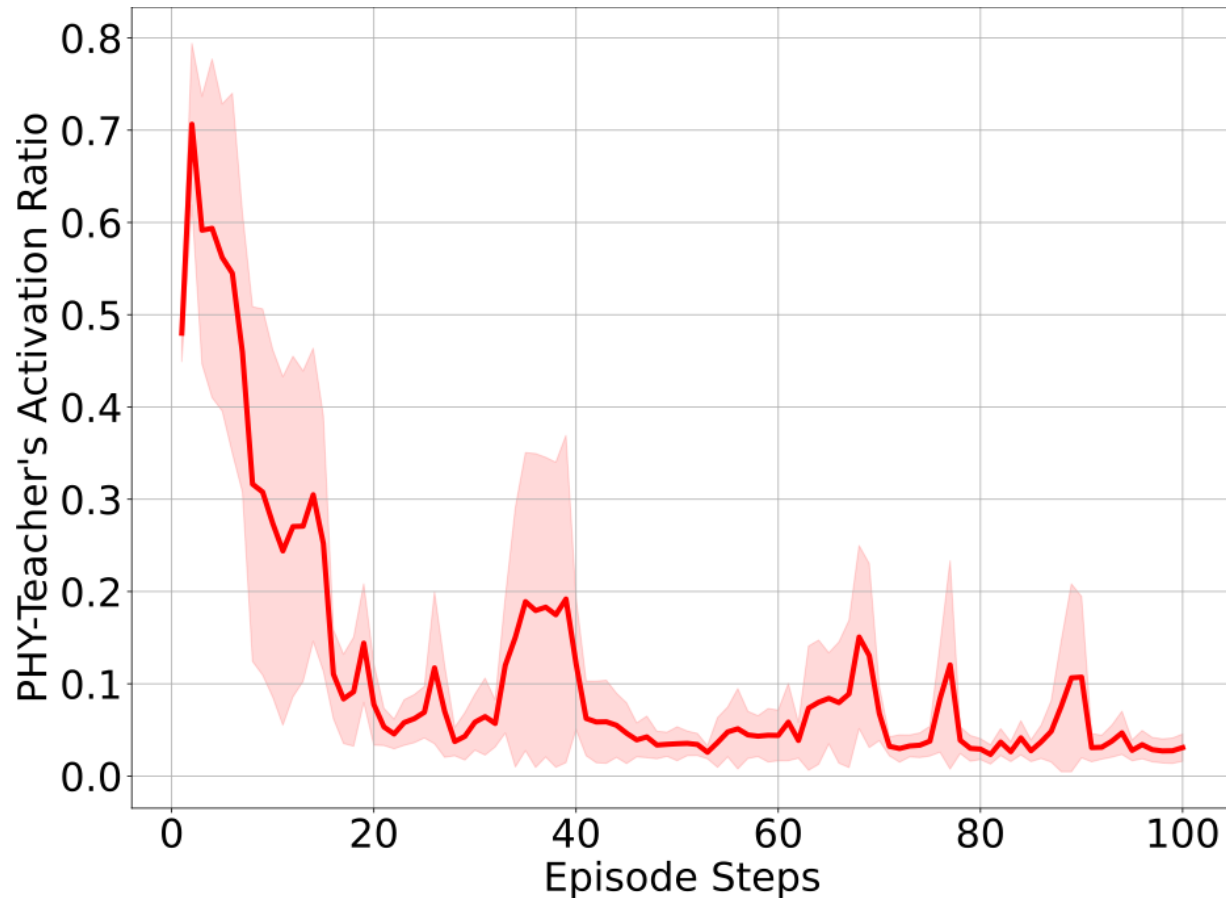
➡ **Automatic Hierarchical Learning:**

**Safety-first** ⤍⤍⤍➤ **High-Performance**

# Experiment-I: Cartpole
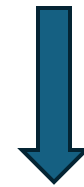
**Feature III:** Automatic Hierarchical Learning



PHY-Teacher Activation Ratio

**PHY-Teacher Activation Ratio:**

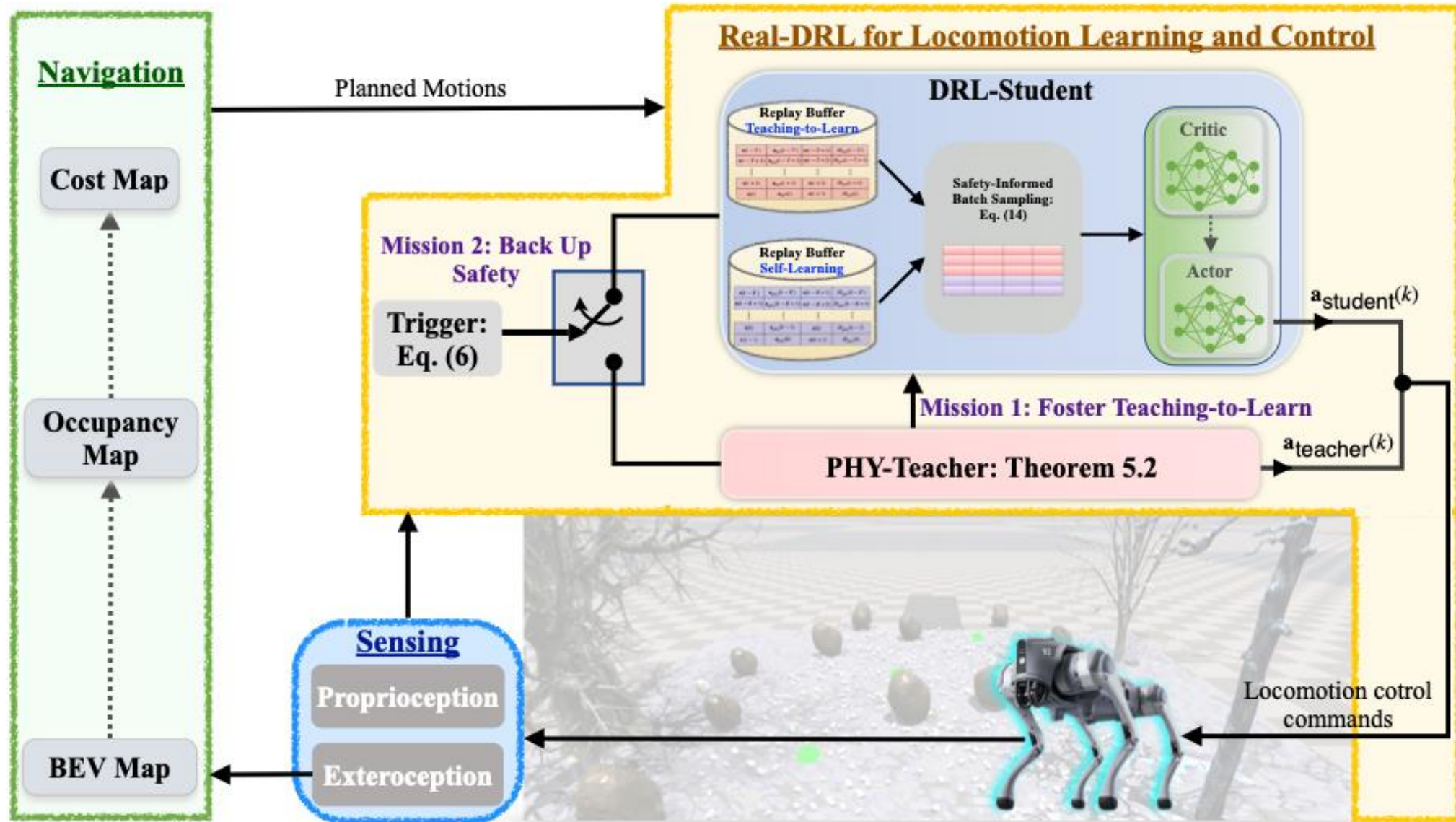$$\frac{\text{PHY-Teacher's total activation times in one episode}}{\text{one episode length}}$$

**The activation ratio of PHY-Teacher within an episode decreases over time**

**DRL-Student becomes independent of PHY-Teacher as learning evolves**

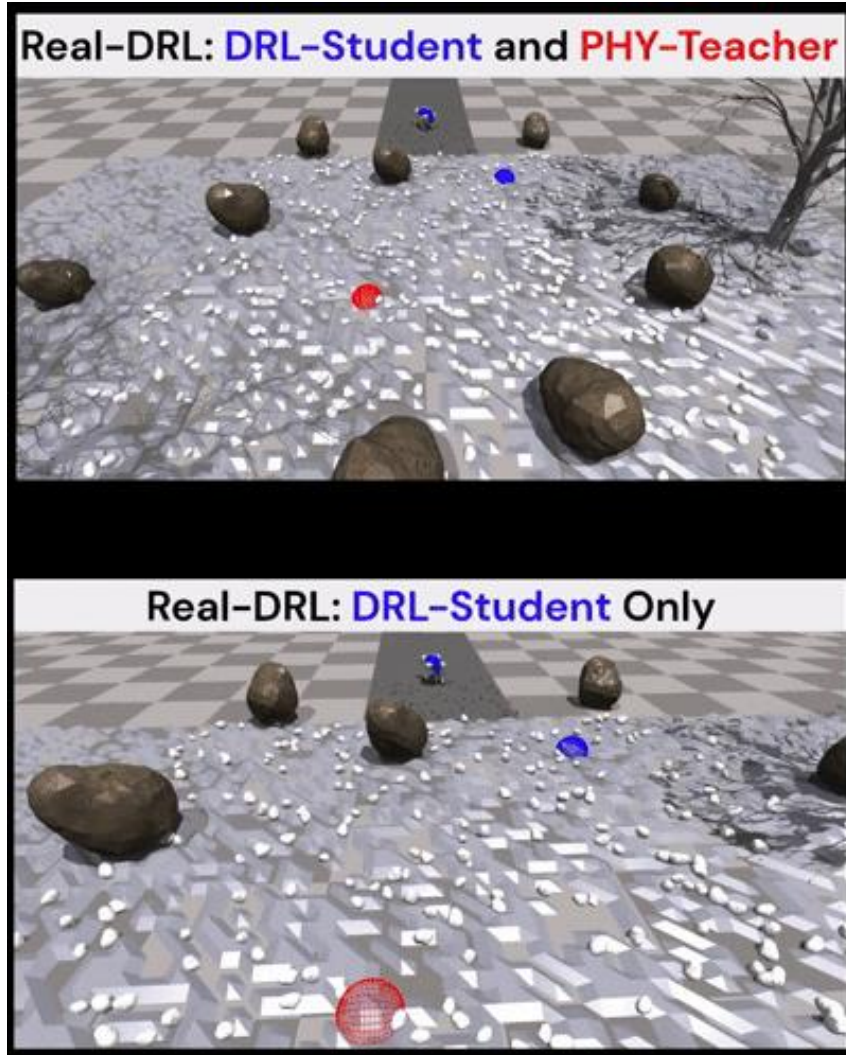# Experiment-II: Go2 in IsaacGym

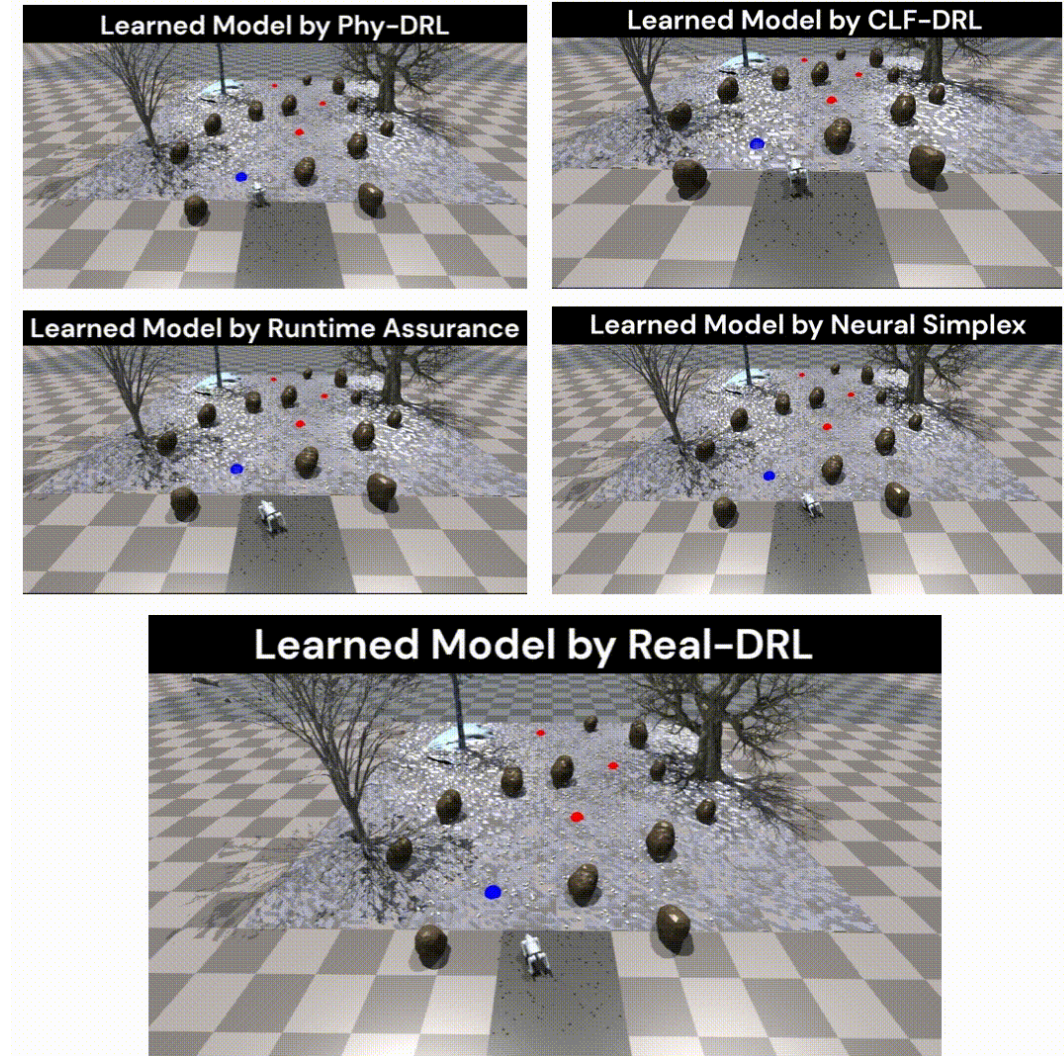## Architecture built on Real-DRL

# Experiment-II: Go2 in IsaacGym

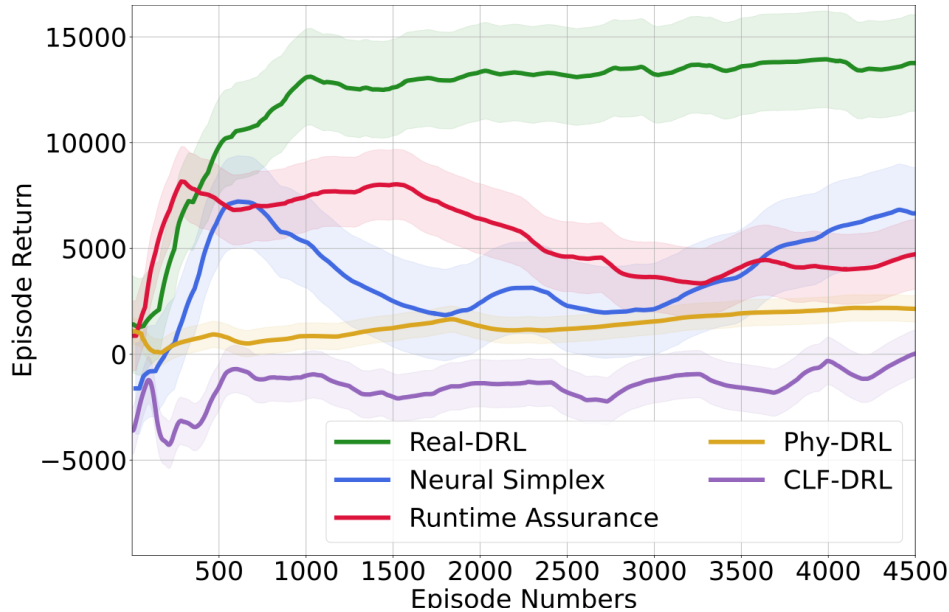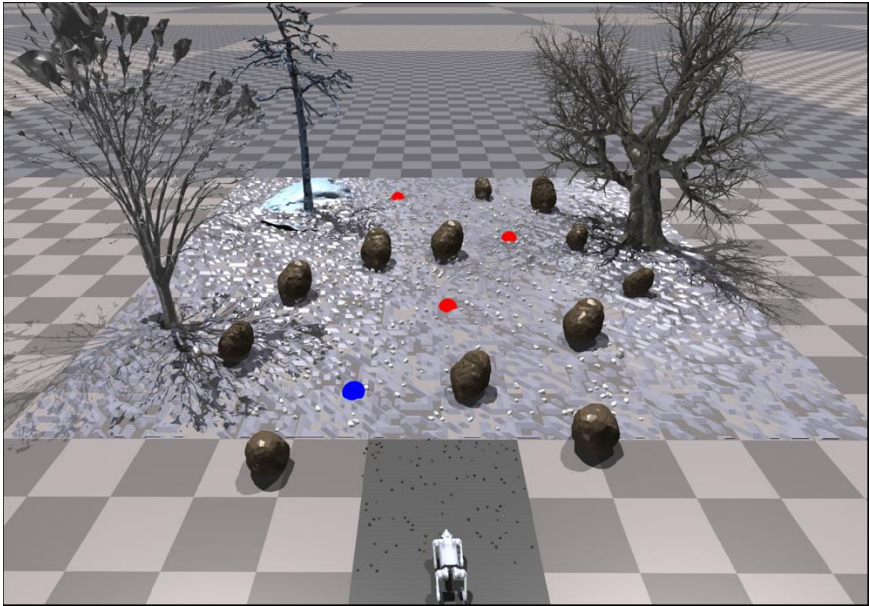## Evaluation Result



Real-DRL in safety guarantee



Real-DRL in learning high-performance policy
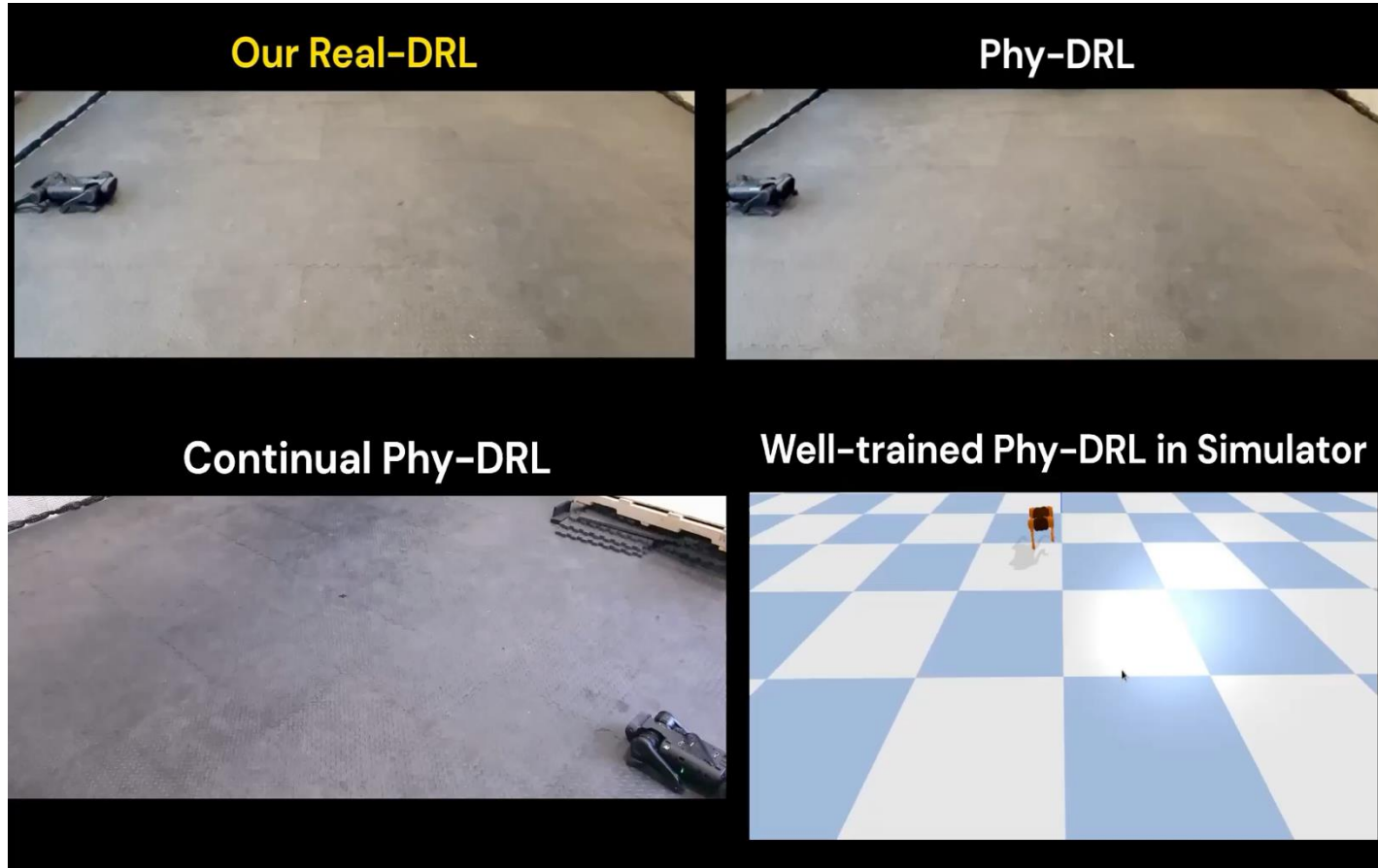
# Experiment-II: Go2 in IsaacGym

## Comparison with SOTA

| Model ID | Navigation Performance | | | | | Energy Efficiency | |
|---|---|---|---|---|---|---|---|
| | **Success** | **Is Fall** | **Collision** | **Num (wp)** | **Travel Time (s)** | **Avg Power (W)** | **Total Energy (J)** |
| CLF-DRL | No | Yes | No | 0 | N/A | N/A | N/A |
| Phy-DRL | No | No | Yes | 1 | ∞ | 507.9441 | ∞ |
| Runtime Assurance | No | Yes | No | 2 | N/A | N/A | N/A |
| Neural Simplex | No | No | Yes | 2 | ∞ | 487.9316 | ∞ |
| PHY-Teacher | **Yes** | No | No | **4** | 55.5327 | 482.8468 | 26817.68 |
| Our Real-DRL | **Yes** | No | No | **4** | **45.3383** | **479.4638** | **21742.42** |

# Experiment-III: A1 in Real World

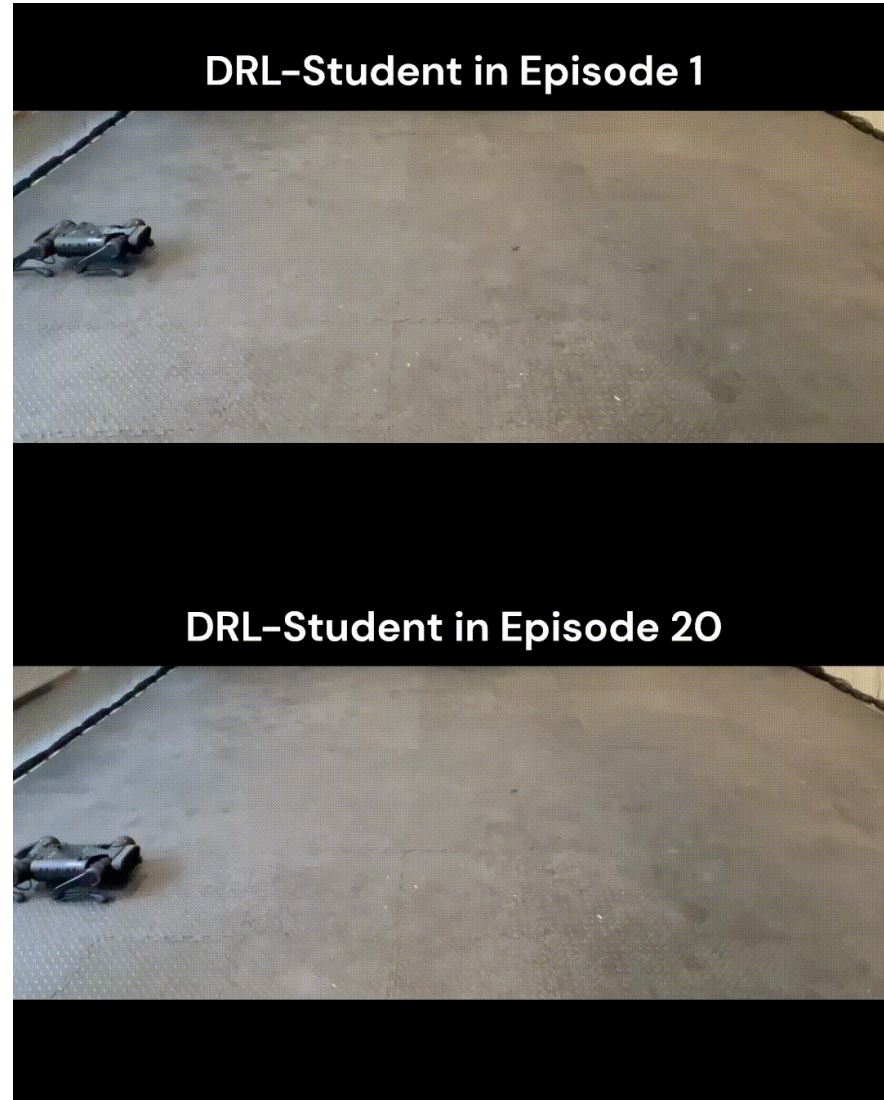## Sim2Real using Real-DRL



Effectiveness of Real-DRL in Sim2Real

Robot Dog Phase Plot

# Experiment-III: A1 in Real World

## Real-DRL Fosters Safety-first Learning



Learning in the first Episode (early stage)

Safe learning with Real-DRL after 20 episodes

# Summary

# Conclusion

➢ *Core Contribution*

- ❑ <u>Real data</u> collection from the hard-to-predict environment

- ❑ <u>Good data</u> (regarding safety) generation from a verifiable PHY-Teacher

- ❑ An <u>innovative RL architecture</u> that supports modular design

➢ *Three Notable Features*

- ❑ <u>Teaching-to-learn Mechanism</u> (e.g., foster safe learning and fast convergence)

- ❑ <u>Automatic Hierarchy Learning</u> (e.g., learn safety first and high-performance policy)

- ❑ <u>Safety-informed batch sampling</u> (e.g., resolve data imbalance caused by corner cases)

➢ *Soundness and Generality*

- ❑ The framework is evaluated across <u>a variety of</u> autonomous systems

- ❑ The experiments incorporate both <u>simulation</u> and <u>real-world</u> evaluations

- ❑ The design of PHY-Teacher provides a <u>theoretical proof</u> of soundness

# Miscellaneous

## ECVXCONE – A Toolbox Towards Real-DRL on Edge Devices

**Cross-Platform** and **Runtime-Efficient** Conic Optimization Toolbox for LMIs

| Hardware Platforms | CPU Configurations | | | Runtime Memory Usage | | LMIs Solve Time | |
|---|---|---|---|---|---|---|---|
| | Arch | Core | Frequency | CVXPY | ECVXCONE | CVXPY | ECVXCONE |
| Dell XPS 8960 Desktop | x86/64 | 32 | 5.4 GHz | 485 MB | 9.87 MB | 49.15 ms | 13.81 ms |
| Intel GEEKOM XT 13 Pro Mini | x86/64 | 20 | 4.7 GHz | 443 MB | 7.32 MB | 61.76 ms | 33.26 ms |
| NVIDIA Jetson AGX Orin | ARM64 | 12 | 2.2 GHz | 423 MB | 8.16 MB | 137.54 ms | 35.73 ms |
| Raspberry Pi 4 Model B | ARM64 | 4 | 1.5 GHz | 436 MB | 8.21 MB | 509.41 ms | 149.87 ms |

Python **CVXPY** vs **C** **ECVXCONE** (Computational Overhead)