# Runtime Learning Machine

HONGPENG CAO, Technical University of Munich, Germany

YANBING MAO, Wayne State University, USA

YIHAO CAI, Wayne State University, USA

LUI SHA, University of Illinois Urbana-Champaign, USA

MARCO CACCAMO, Technical University of Munich, Germany

This paper proposes the runtime learning machine for safety-critical learning-enabled cyber-physical systems (CPS). The learning machine has three interactive components: a high-performance (HP)-Student, a high-assurance (HA)-Teacher, and a Coordinator. The HP-Student is a high-performance but not fully verified Phy-DRL (physics-regulated deep reinforcement learning) agent that performs safe runtime learning in <u>real</u> plant, using <u>real</u>-time sensor data from <u>real</u>-time physical environments. On the other hand, HA-Teacher is a verified but simplified design, focusing on safety-critical functions. As a complementary, HA-Teacher's novelty lies in real-time patch for two missions: i) correcting unsafe learning of HP-Student, and ii) backing up safety. The Coordinator manages the interaction between HP-Student and HA-Teacher. Powered by the three interactive components, the runtime learning machine notably features i) assuring lifetime safety (i.e., safety guarantee in any runtime learning stage), ii) tolerating unknown unknowns, iii) addressing Sim2Real gap, and iv) automatic hierarchy learning (i.e., safety-first learning, and then high-performance learning). Experiments involving a cart-pole system and two quadruped robots, as well as comparisons with state-of-the-art safe DRL, fault-tolerant DRL, and approaches for addressing Sim2Real gap, demonstrate the learning machine's effectiveness and unique features.

CCS Concepts: • **Computing methodologies** → **Machine learning**; • **Computer systems organization** → **Reliability**.

Additional Key Words and Phrases: Runtime Learning, Safety, Phy-DRL, Real-time Patch, Unknown Unknowns, Sim2Real Gap

## 1 Introduction

Deep reinforcement learning (DRL) has been incorporated into numerous cyber-physical systems (CPS) and has shown significant advancements in making sequential and complex decisions in various CPS fields, such as autonomous driving [28, 30] and robot locomotion [24, 31]. The DRL-enabled CPS has the potential to revolutionize many processes across different industries, leading to tangible economic impacts [46]. However, the public-facing AI Incident database in [1] has revealed that machine learning (ML) techniques, including DRL, can achieve remarkable performance without ensuring safety [52]. For instance, a report by the National Highway Traffic Safety Administration highlighted 351 car crashes related to advanced driver assistance systems from July 2023 to March 2024 in the US alone [35]. Therefore,

ensuring high-performance DRL with verifiable safety is even more crucial for CPS today, aligning well with the market's demand for safe ML techniques.

## 1.1 Safety Challenges and Open Problems

Our considered safety challenges are rooted in the unknown unknowns and the Sim2Real gap, which are detailed below.

**Safety Challenge 1: Unknown Unknowns**. The unknown unknowns generally refer to outcomes, events, circumstances, or consequences that are not known in advance and cannot be predicted in time and distributions [4]. The dynamics of many safety-critical CPS (e.g., autonomous vehicles [39], airplanes [41], and quadrupedal robots [7]) are governed by a combination of known knowns (e.g., Newton's laws of motion), known unknowns (e.g., Gaussian noise without knowing to mean and variance), and unknown unknowns. The unknown unknowns are due to, for example, unforeseen operating environments and DNN's colossal parameter space, intractable activation, and hard-to-verify. The safety assurance also requires resilience to unknown unknowns, which is very challenging. The reasons stem from characteristics of unknown unknowns: there is almost zero historical data, unpredictable timing and distributions, resulting in the unavailability of models for scientific discoveries and understanding.

**Safety Challenge 2: Sim2Real Gap**. The prevalent DRL involves training a policy within a simulator using synthetic data and deploying it to real CPS platforms. However, the difference between the simulated environment and the real CPS world creates a gap known as the Sim2Real gap. This gap causes a drop in performance when using pre-trained DRL in real CPS environments. Numerous approaches have been developed to address the Sim2Real gap [13, 19, 25, 34, 36, 45, 47, 49, 51]. These methods aim to improve the realism of the simulator and can mitigate the Sim2Real and domain gaps to varying degrees. Nevertheless, undisclosed gaps and missing dynamics continue to hinder the safety assurance of real CPS.

To address Safety Challenges 1 and 2, the most appealing solution is the *Prospect: Runtime learning for a high-performance action policy in _real_ plant – using _real_-time sensor data generated from _real_-time physical environments while prioritizing safety!* However, two open problems arise about bringing the prospect into reality.

**Problem 1**: *If the DRL agent's actions lead to a safety violation, how can we correct his unsafe learning and back up the safety of real plants in a timely manner?*

**Problem 2**: *How to tolerate and also teach the DRL agent to tolerate unknown unknowns and Sim2Real gap for assuring safety of real plants?*

## 1.2 Related Work

Significant efforts have been devoted to enhance DRL safety for CPS by developing safe DRL and fault-tolerant DRL, which are summarized below.

**Safe DRL**. One research focus of safe DRL is the safety-embedded reward, as a DRL agent must learn a high-performance action policy with verifiable safety. The control Lyapunov function (CLF) proposed in [6, 14, 37, 54] is a candidate. Meanwhile, seminal work in [48] revealed that a CLF-like reward could enable DRL with verifiable stability. At the same time, enabling verifiable safety is achievable by extending CLF-like rewards with given safety regulations. However, systematic guidance for constructing such CLF-like rewards remains open. The residual action policy is another shift in safe DRL, which integrates data-driven action policy and physics-model-based action policy. The existing residual diagrams focus on stability guarantee [17, 26, 32, 40], with the exception being [16] on safety guarantee. However, the physics models considered are nonlinear and intractable, which thwarts delivering a verifiable safety guarantee or assurance, if not impossible. The recently developed Phy-DRL (physics-regulated DRL) framework

[11, 12] can satisfactorily address the open problems of safe DRL. Summarily, Phy-DRL permits simplifying the model of nonlinear dynamics to an analyzable and tractable linear one. This linear model can then be a model-based guide for constructing the safety-embedded (CLF-like) reward and residual action policy. Meanwhile, the Phy-DRL exhibits verifiable safety. However, it is only mathematically or theoretically possible due to the underlying assumptions of manageable Sim2Real gap and unknown unknowns. In other words, Phy-DRL cannot offer verifiable safety for real plants in the face of unknown unknowns and the Sim2Real gap.

***Fault-tolerant DRL.*** This is another direction for DRL safety in real plants. Recent approaches include neural Simplex [38], runtime assurance [9, 15, 44], and model predictive shielding [3, 5]. They treat the DRL agent as a high-performance module (HPM) but a black box that runs in parallel with a verified high-assurance module (HAM). Normally, HPM controls the real plants. HAM takes over once safety violation occurs. These architectures can ensure the safe running of DRL in real plants under the assumption that Challenges 1 and 2 do not cause HAM to fail, which is not practical for systems whose operating environments are dynamic and unpredictable. Furthermore, they are not solutions to Problems 1 and 2. Specifically, in all these architectures, HAM and HPM are independent, that is, HPM cannot learn from HAM, and HAM cannot teach HPM how to be safe. Meanwhile, HAM is the static model-based controller, and its action will be unreliable if the real-time unknown unknowns and Sim2Real gap create a significant model mismatch.

## 1.3   Contribution: Runtime Learning Machine: From Theory To Implementation

To address Safety Challenges 1 and 2 and answer Problems 1 and 2, we propose the runtime learning machine, whose framework is shown in Figure 1. The machine constitutes high-performance (HP)-Student, high-assurance (HA)-Teacher, and Coordinator. HP-Student is a Phy-DRL agent that can be pre-trained and continue to learn in real plants that operate in real-time physical environments. HA-Teacher is a verified and physics-based design, with its functionality being reduced to a safety-critical level. Coordinator manages interactions between HP-Student and HA-Teacher. *As a metaphor,*



Fig. 1.  Runtime learning machine framework.

*HP-Student's runtime learning in our machine is like a student's journey. First, he learns from teachers in middle school, high school, college, etc., who have verified domain knowledge in subjects like physics and mathematics, to gain essential knowledge. Then, he delves deeper into specific areas during graduate studies to acquire expertise in those fields.* Summarily, our runtime machine learning has following three distinct characteristics.

***Characteristic 1: Automatic Hierarchy Learning Mechanism***. HP-Student's growth in our runtime learning machine is an automatic hierarchical learning mechanism that respects safety-first principles for safety-critical CPS without compromising mission performance. As depicted in Figure 2, HP-Student undergoes a two-stage learning process:
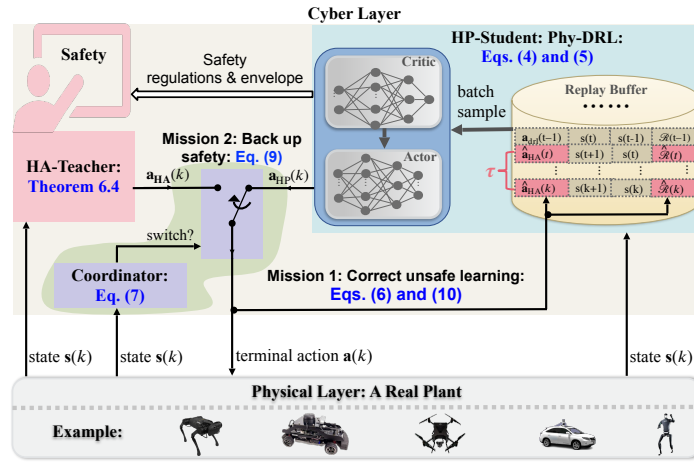
- Stage 1: Safety-first Learning. HP-Student first learns from HA-Teacher how to be safe (i.e., constraining the system states of real plants into a safety set). Meanwhile, Figure 2 illustrates that prioritizing safety does not compromise mission performance. In other words, violating safety protocols results in decreased mission performance.

- Stage 2: Self High-performance Learning. After HP-Student has learned how to control system states within safety envelopes, Coordinator rarely



Fig. 2. HP-Student's two learning stages.

activates HA-Teacher. Consequently, HP-Student engages in self-learning within the safety envelope for a high-performance action policy, such as the car closely following the planned blue path in stage 2 in Figure 2.
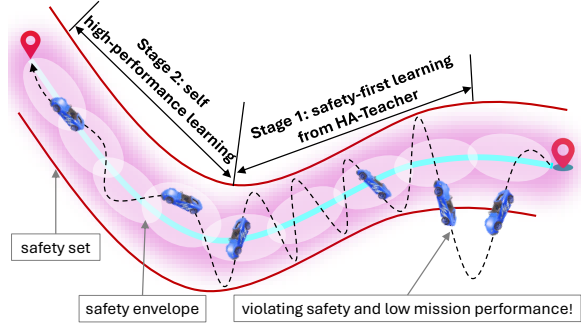
***Characteristic 2: Assuring Safety by Tolerating Unknown Unknowns and Sim2Real Gap.*** HA-Teacher's real-time patch, enabled by a real-time model, real-time action mission, and real-time model-based policy computation, aim to ensure lifetime safety. This means guaranteeing the safety of real plants during any runtime learning stages, regardless of HP-Student's failures, and in the face of real-time unknown unknowns and the Sim2Real gap.

***Characteristic 3: Highly Interactive HP-Student and HA-Teacher.*** The interactions between HP-Student and HA-Teacher in the runtime learning machine occur in two dimensions:

- HP-Student $\longrightarrow$ HA-Teacher: HP-Student shares his safety regulations and envelope with HA-Teacher for his real-time patch design.

- HP-Student $\longleftarrow$ HA-Teacher: Showing in Figure 1, HA-Teacher has two missions: i) correct unsafe learning of HP-Student and ii) back up the safety of the real plants, in the face of unknown unknowns and Sim2Real gap.

| Notations throughout Paper | |
|---|---|
| $\mathbb{R}^n$ | Set of $n$-dimensional real vectors |
| $\mathbb{N}$ | Set of natural numbers |
| $[\mathbf{x}]_i$ | The $i$-th entry of vector $\mathbf{x}$ |
| $[\mathbf{W}]_{i,:}$ | The $i$-th row of matrix $\mathbf{W}$ |
| $[\mathbf{W}]_{:,j}$ | The $j$-th column of matrix $\mathbf{W}$ |
| $[\mathbf{W}]_{i,j}$ | Matrix $\mathbf{W}$'s element at row $i$ and column $j$ |
| $\mathbf{P} \succ (\prec)\ 0$: | Matrix $\mathbf{P}$ is positive (negative) definite |
| $\top$ | Transposition of a matrix or vector |
| $\mathbf{P}^{-1}$ | Inverse of matrix $\mathbf{P}$ |
| $\lceil \cdot \rceil$ | Ceiling function |
| $\ln(a)$ | Natural logarithm of the number $a > 0$ |

## 2 Preliminaries: Definitions of Safety and High Performance

We introduce the dynamics model of a DRL-enabled CPS:

$$\mathbf{s}(k+1) = \mathbf{A}(\mathbf{s}(k)) \cdot \mathbf{s}(k) + \mathbf{B}(\mathbf{s}(k)) \cdot \mathbf{a}(k) + \mathbf{f}(\mathbf{s}(k)),\ k \in \mathbb{N} \qquad (1)$$

whose equilibrium point is $\mathbf{s}^* = \mathbf{0}$. In Equation (1), $\mathbf{f}(\mathbf{s}(k)) \in \mathbb{R}^n$ is the model mismatch, $\mathbf{A}(\mathbf{s}(k)) \in \mathbb{R}^{n \times n}$ and $\mathbf{B}(\mathbf{s}(k)) \in \mathbb{R}^{n \times m}$ denote system matrix and control structure matrix, respectively, $\mathbf{s}(k) \in \mathbb{R}^n$ is real plant's state in real-time, $\mathbf{a}(k) \in \mathbb{R}^m$ is the action command in real-time.

The safety issue arises from the system's state $\mathbf{s}(k)$ and the associated safety regulations, defining the permissible state space for the system:

$$\text{Safety set: } \mathbb{X} \triangleq \left\{ \mathbf{s} \in \mathbb{R}^n \middle| \underline{\mathbf{v}} \leq \mathbf{D} \cdot \mathbf{s} \leq \overline{\mathbf{v}}, \quad \text{with } \mathbf{D} \in \mathbb{R}^{h \times n}, \overline{\mathbf{v}}, \underline{\mathbf{v}} \in \mathbb{R}^h \right\}. \tag{2}$$

where $\mathbf{D}$, $\overline{\mathbf{v}}$ and $\underline{\mathbf{v}}$ are given in advance for formulating $h \in \mathbb{N}$ safety regulations. Inequalities in Equation (2) are generic, as they can cover many safety regulations, such as speed regulation, collision avoidance, lane keeping, and tracking for autonomous vehicles. Building on the safety set, the lifetime safety is formally defined below.

**Definition 2.1** (Lifetime Safety). Consider the safety set $\mathbb{X}$ in Equation (2). The CPS in Equation (1) is said to have lifetime safety, if given any $\mathbf{s}(1) \in \mathbb{X}$, the $\mathbf{s}(k) \in \mathbb{X}$ holds at any time $k \in \mathbb{N}$, regardless of HP-Student's failure.

**Definition 2.2.** 'High Performance' in this paper has two-dimensional definition: 1) mission performance (measured by, for example, tracking errors in the lane-tracking and path-following tasks) and 2) operation performance (measured by, for example, jerky movements for customers' comfort in autonomous vehicles). In the learning machine, HP-Student's reward encodes safety regulations, mission, and operation for learning a high-performance action policy with a safety guarantee. On the other hand, HA-Teacher's function is reduced to be safety-critical only, and his performance consideration is only about the operation regulations.

## 3 Design Overview

Our proposed runtime learning machine aims to address Safety Challenges 1 and 2 and answer Problems 1 and 2. To do so, showing in Figure 1, it is designed to have three interactive components:

- <u>HP-Student</u> builds on Phy-DRL agent, which can be pre-trained in a simulator or another domain but performs runtime learning in a real CPS to tolerate unknown unknowns and address the Sim2real gap.
- <u>HA-Teacher</u> is a verified safety-only design whose novelty lies in real-time patches with two missions: timely correcting unsafe learning of HP-Student and backing up safety of real CPS.
- <u>Coordinator</u> is responsible for monitoring the real-time safety status and facilitating interactions between HP-Student and HA-Teacher. Specifically, when the real-time safety status of the plant being controlled by HP-Student approaches the safety boundary, Coordinator prompts HA-Teacher to intervene and assure the safety of real plant, and correct unsafe learning of HP-Student. When the real-time states return to a safe region, Coordinator triggers the switch back to HP-Student and terminates the learning correction.

Next, we will describe the designs of the three interactive components in Sections 4 to 6, respectively.

## 4 Runtime Learning Machine: HP-Student Component

### 4.1 HP-Student Candidates

We acknowledge that DRL is unable to directly embed high-dimensional or many safety regulations into the reward function due to the reward being a one-dimensional real value, creating a dimension gap. To bridge this dimension gap, the literature [11, 12] introduces the concept of a safety envelope, which has the one-dimensional condition and can be

designed as a subset of the safety set $\mathbb{X}$ (see Figures 2 and 3 for visualization).

$$\text{Safety envelope: } \Omega \triangleq \left\{ \mathbf{s} \in \mathbb{R}^n \,\middle|\, \mathbf{s}^\top \cdot \mathbf{P} \cdot \mathbf{s} \le 1, \ \mathbf{P} \succ 0 \right\}. \tag{3}$$

So, our HP-Student candidates are those whose rewards can successfully embed the safety envelope in Equation (3), and such safety-embedded rewards can be shared with HA-Teacher for his real-time patch design. Along with this direction, DRL with CLF-like reward proposed in [48] and Phy-DRL (physics-regulated DRL) proposed in [11, 12] are two preferred candidates, as they can successfully embed the safety envelope into their rewards. Finally, HP-Student adopts Phy-DRL because Phy-DRL also features fast training theoretically and experimentally, which is desirable for runtime learning in real plants. Next, we will review the HP-Student design.

### 4.2 HP-Student: Phy-DRL: Residual Action Policy and Safety-embedded Reward

Recalling Phy-DRL in [11, 12], HP-Student has residual action policy formula:

$$\mathbf{a}_{\text{HP}}(k) = \underbrace{\mathbf{a}_{\text{drl}}(k)}_{\text{data-driven}} + \underbrace{\mathbf{a}_{\text{phy}}(k) \ (= \mathbf{F} \cdot \mathbf{s}(k))}_{\text{model-based}}, \tag{4}$$

where $\mathbf{a}_{\text{drl}}(k)$ denotes a date-driven action from DRL, while $\mathbf{a}_{\text{phy}}(k)$ is a model-based action. Referring to safety envelope in Equation (3), HP-Student's safety-embedded reward is

$$\mathcal{R}(\mathbf{s}(k), \mathbf{a}_{\text{drl}}(k)) = \underbrace{\mathbf{s}^\top(k) \cdot \mathbf{P} \cdot \mathbf{s}(k) - \mathbf{s}^\top(k+1) \cdot \mathbf{P} \cdot \mathbf{s}(k+1)}_{\triangleq \ r(\mathbf{s}(k), \ \mathbf{s}(k+1))} + w(\mathbf{s}(k), \mathbf{a}_{\text{HP}}(k)), \tag{5}$$

where the sub-reward $r(\mathbf{s}(k), \mathbf{s}(k+1))$ is safety-embedded, while the sub-reward $w(\mathbf{s}(k), \mathbf{a}(k))$ aims at high operation performance (e.g., minimizing energy consumption of resource-limited robots and avoiding jerks for customers' comfort in autonomous vehicles). The matrices $\mathbf{F}$ in Equation (4) and $\mathbf{P}$ in Equation (3) and Equation (5) are the design variables. Their automatic computation by the CVXPY toolbox is detailed in [12].

*Remark* 4.1 (**Safety- And Also Mission-Embedded**). The equilibrium $\mathbf{s}^* = \mathbf{0}$ means that the system described in Equation (1) can be interpreted as the dynamics of mission-tracking error. For instance, in a path-following task, the path represents the mission goal, while $\mathbf{s}(k)$ denotes the real-time tracking error of the path. Additionally, as indicated in Equation (3), the center of the safety envelope is the $\mathbf{s}^* = \mathbf{0}$. Based on this, we can conclude that the sub-reward $r(\mathbf{s}(k), \mathbf{s}(k+1))$ defined in Equation (5) encompasses both safety and mission considerations, and HP-Student's learning encourages actions that increase $r(\mathbf{s}(k), \mathbf{s}(k+1))$ over time. Furthermore, an increase in $r(\mathbf{s}(k), \mathbf{s}(k+1))$ signifies progress towards both the envelope center and the mission goal. This also explains Figure 2, where prioritizing safety does not compromise mission performance (violating safety protocols results in decreased mission performance).

### 4.3 HP-Student: Correction of Unsafe Runtime Learning

HP-Student can be pre-trained in a simulator or another domain, and then he performs runtime learning in real plants within a real-time physical environment. HP-Student utilizes the actor-critic architecture-based DRL such as those outlined in [33] and [23] for runtime learning, in order to learn a safe data-driven policy that maximizes the expected return. HP-Student consists of an action policy and an action-value function.

Sampling efficiency is crucial for runtime learning. Experience replay (ER) [2] enables off-policy algorithms to reuse past experiences, significantly improving sampling efficiency and preventing forgetting of learned knowledge [29].

ER also helps break the correlation between adjacent transitions to avoid sampling bias for a stable learning process, which is important when online data is limited due to the expensive interaction with physical systems. During online inference, we continuously store real transitions resulting from the actions of HP-Student and corrected unsafe actions by HA-Teacher in the replay buffer. As shown in Figure 1, if the action $\mathbf{a}_{\mathrm{HP}}(k)$ from HP-Student leads to unsafe behavior of a real plant, HA-Teacher takes control to ensure safety of real plant, and corrects the unsafe data-driven action to $\widehat{\mathbf{a}}_{\mathrm{HA}}(k)$ and the corresponding reward to $\widehat{\mathcal{R}}(k)$, according to

$$\mathbf{a}_{\mathrm{drl}}(k) \leftarrow \widehat{\mathbf{a}}_{\mathrm{HA}}(k) \triangleq \mathbf{a}_{\mathrm{HA}}(k) - \mathbf{a}_{\mathrm{phy}}(k), \qquad \mathcal{R}(\mathbf{s}(k), \mathbf{a}_{\mathrm{drl}}(k)) \leftarrow \widehat{\mathcal{R}}(k) \triangleq \mathcal{R}(\mathbf{s}(k), \widehat{\mathbf{a}}_{\mathrm{HA}}(k)), \tag{6}$$

where $\mathbf{a}_{\mathrm{phy}}(k)$ is HP-Student's model-based action in Equation (4), and $\mathbf{a}_{\mathrm{HA}}(k)$ is the action from HA-Teacher, whose design is presented in Section 6. In the meantime, during runtime learning, a minibatch of transitions is uniformly sampled for training or learning [21].

*Remark* 4.2. Equation (6) states that according to HP-Student's residual action policy in Equation (4), the action correction is only applied to the data-driven $\mathbf{a}_{\mathrm{drl}}(k)$, as the model-based action policy $\mathbf{a}_{\mathrm{phy}}(k) = \mathbf{F} \cdot \mathbf{s}(k)$ is invariant.

## 5 Runtime Learning Machine: Coordinator Component

Coordinator manages interactions between HP-Student and HA-Teacher according to

$$\text{Triggering condition: } \mathbf{s}^{\top}(k-1) \cdot \mathbf{P} \cdot \mathbf{s}(k-1) \leq 1 \text{ and } \mathbf{s}^{\top}(k) \cdot \mathbf{P} \cdot \mathbf{s}(k) > 1, \tag{7}$$

coupled with which, we introduce the active time phase of HA-Teacher:

$$\text{HA-Teacher's active phase: } \mathbb{T}_{\sigma(k)} \triangleq \{k,\ k+1,\ \ldots,\ k+\tau\},\ \tau \in \mathbb{N} \tag{8}$$

where $\sigma(k)$ represents a piece-wise signal for notation. For instance, $\sigma(k) = i$ for $k \in \mathbb{T}_{\sigma(k)}$ signifies the $i$-th time that HA-Teacher is triggered, and its active phase this time is $\mathbb{T}_i$. The switching logic of actions applied to a real plant for backing up safety is as follows:

$$\mathbf{a}(t) \leftarrow \begin{cases} \mathbf{a}_{\mathrm{HA}}(t), & \text{if triggering condition (7) holds at } k \text{ and } t \in \mathbb{T}_{\sigma(k)} \\ \mathbf{a}_{\mathrm{HP}}(t), & \text{otherwise} \end{cases} \tag{9}$$

synchronizing with which is the correcting logic of HP-Student's unsafe action and reward:

$$\mathbf{a}_{\mathrm{drl}}(t) \leftarrow \begin{cases} \widehat{\mathbf{a}}_{\mathrm{HA}}(t), & \text{if triggering condition (7) holds at } k \text{ and } t \in \mathbb{T}_{\sigma(k)} \\ \mathbf{a}_{\mathrm{drl}}(k), & \text{otherwise} \end{cases} \tag{10a}$$

$$\mathcal{R}(t) \leftarrow \begin{cases} \widehat{\mathcal{R}}(t), & \text{if triggering condition (7) holds at } k \text{ and } t \in \mathbb{T}_{\sigma(k)} \\ \mathcal{R}(\mathbf{s}(t), \mathbf{a}_{\mathrm{drl}}(t)), & \text{otherwise} \end{cases} \tag{10b}$$

where $\widehat{\mathbf{a}}_{\mathrm{HA}}(t)$ and $\widehat{\mathcal{R}}(t)$ are the corrected action and reward by HA-Teacher, defined in Equation (6).

*Remark* 5.1 (**Enabling Automatic Hierarchy Learning**). Operating within the safety envelope, Coordinator activates HA-Teacher, if the real-time states of the CPS move outside the safety envelope. Once the active phase ends, control transitions back to HP-Student, and HA-Teacher's correction of unsafe learning concludes. If condition (7) is no longer met, HP-Student will have successfully learned to control the CPS within safety envelope, and continual runtime learning will then focus on achieving high mission and operation performance.

*Remark* 5.2 (**Active Phase**). Referring to Equations (8) to (10), the symbol $\tau$ represents the correction horizon for unsafe action and reward of HP-Student and the dwell time of HA-Teacher. Its allowable minimum value is one. However, if the value of $\tau$ is very small, the patch center may not sufficiently attract system states to the envelope inside, and HA-Teacher will dominate the learning machine, only ensuring safety. Corollary C.1 in Appendix C guides determining the appropriate value for $\tau$.

## 6  Runtime Learning Machine: HA-Teacher Component

Enabling runtime learning in real plants is straightforward in addressing the Sim2Real gap, but not so for unknown unknowns, because unknown unknowns lack historical data and cannot be predicted in time and distribution. When an unknown unknown creates safety issues in a time-critical environment, it is crucial to update the dynamics models, action plans, and mission goals promptly to ensure safe and effective responses in real time. The insight inspires us to develop the real-time patch as the HA-Teacher. Its model knowledge, action policy, and mission goal are dynamic and real-time. The mathematical formula for a real-time patch is

$$\Psi_{\sigma(k)} \triangleq \{ \mathbf{s} \mid (\mathbf{s} - \chi \cdot \widehat{\mathbf{s}}_{\sigma(k)})^{\top} \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot (\mathbf{s} - \chi \cdot \widehat{\mathbf{s}}_{\sigma(k)}) \le (1-\chi)^2 \cdot \widehat{\mathbf{s}}_{\sigma(k)}^{\top} \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \widehat{\mathbf{s}}_{\sigma(t)} \}, \tag{11}$$

coupled with which is the real-time action policy:

$$\mathbf{a}_{\text{HA}}(k) = \widehat{\mathbf{F}}_{\sigma(k)} \cdot (\mathbf{s}(k) - \chi \cdot \widehat{\mathbf{s}}_{\sigma(k)}), \text{ with } \chi \in (0,1) \text{ such that } \chi^2 \cdot \widehat{\mathbf{s}}_{\sigma(k)}^{\top} \cdot \mathbf{P} \cdot \widehat{\mathbf{s}}_{\sigma(k)} < 1, \tag{12}$$

where $\widehat{\mathbf{P}}_{\sigma(k)} \succ 0$, the $\chi \cdot \widehat{\mathbf{s}}_{\sigma(k)}$ represents the patch center (i.e., the yellow dots in Figure 3), and the $\widehat{\mathbf{s}}_{\sigma(k)}$ denotes the real-time state that triggers HA-Teacher and remains constant for defining patch center during HA-Teacher's active phase $\mathbb{T}_{\sigma(k)}$ (defined in Equation (8)), i.e.,

$$\widehat{\mathbf{s}}_{\sigma(t)} = \mathbf{s}(k) \text{ for } t \in \mathbb{T}_{\sigma(k)}, \text{ with } \mathbf{s}(k) \text{ satisfying triggering condition (7)}. \tag{13}$$

*Remark* 6.1 (**Why named patch?**). In today's world, there are two approaches to achieving the same control task of CPS: a high-dimensional data-driven DRL and a low-dimensional physics-model-based controller. The data-driven DRL provides superior performance but is challenging to verify (due to DNN's huge parameter, nonlinear activation, etc.). On the other hand, the physics-model-based approach offers analyzable and verifiable behavior but has limited performance (due to model mismatch). This explains why the set in Equation (11) follows a very similar safety envelope formula in Equation (3), but it is referred to as a patch: the envelope represents a DRL design, while the patch represents a physics-model-based design with a small verifiable-safety region.

When a plant under the control of HP-Student experiences a safety violation at time $k$ (as indicated by the condition in Equation (7)), Coordinator activates HA-Teacher. HA-Teacher then utilizes real-time sensor data $\widehat{\mathbf{s}}_{\sigma(k)}$ to update the physics-model knowledge $(\mathbf{A}(\widehat{\mathbf{s}}_{\sigma(k)}), \mathbf{B}(\widehat{\mathbf{s}}_{\sigma(k)}))$. This update is used to compute the real-time patch in Equation (11) and the coupled action policy in Equation (12). The real-time patch and its coupled action policy will empower HA-Teacher to achieve backing up safety and correcting unsafe learning of HP-Student. However, to deliver the targeted capabilities, real-time patch must meet following three requirements.

Fig. 3. System behavior.

Requirement 1: Attracting Toward Safety Envelope. The center of the patch must be within the safety envelope. If it's not, as shown by the patch $\Psi_4$ in Figure 3, the system's state can get stuck in the
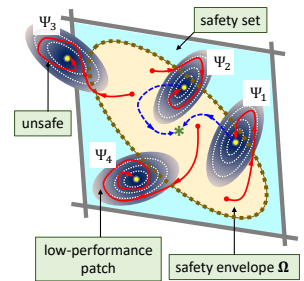
patch. This can lead to HA-Teacher dominating the machine during runtime learning, and HP-Student being unable to self-learn for a high-performance action policy.

Requirement 2: Conformity with Safety Regulations. The real-time patches must be subsets of the safety set in Equation (2). If not, the patch will not be able to ensure safety, as shown by patch $\Psi_3$ in Figure 3, where the system states leave the safety set.

Requirement 3: Conformity with Operation Regulations. It is necessary to confine the real-time action $\mathbf{a}_{\text{HA}}(k)$ within a physically-feasible bounded action space:

$$\mathbb{A} \triangleq \left\{ \mathbf{a}_{\text{HA}} \in \mathbb{R}^m \,\middle|\, \underline{\mathbf{z}} \leq \mathbf{C} \cdot \mathbf{a}_{\text{HA}} \leq \overline{\mathbf{z}}, \text{ with } \mathbf{C} \in \mathbb{R}^{g \times m}, \, \underline{\mathbf{z}}, \, \overline{\mathbf{z}} \in \mathbb{R}^m \right\}, \tag{14}$$

where $\mathbf{C}$, $\overline{\mathbf{z}}$ and $\underline{\mathbf{z}}$ are given in advance for formulating operation regulations.

The $\widehat{\mathbf{F}}_{\sigma(k)}$ and $\widehat{\mathbf{P}}_{\sigma(k)}$ in Equation (12) and Equation (11) are our design variables for delivering the real-time patch and coupled action policy. Our design focus is on how $\widehat{\mathbf{F}}_{\sigma(k)}$ and $\widehat{\mathbf{P}}_{\sigma(k)}$ can meet Requirements 1–3. We observe from Equations (11) and (12) that the patch center $\chi \cdot \widehat{\mathbf{s}}_{\sigma(k)}$ meets Requirement 1 because it is located inside the safety envelope (due to $\chi^2 \cdot \widehat{\mathbf{s}}_{\sigma(k)}^\top \cdot \mathbf{P} \cdot \widehat{\mathbf{s}}_{\sigma(k)} < 1$) to attract systems toward the envelope. So, the remaining task is to follow Requirements 2 and 3 to design $\widehat{\mathbf{F}}_{\sigma(k)}$ and $\widehat{\mathbf{P}}_{\sigma(k)}$, which relies on a tracking-error dynamics model obtained from Equation (1):

$$\mathbf{e}(k+1) = \mathbf{A}(\widehat{\mathbf{s}}_{\sigma(k)}) \cdot \mathbf{e}(k) + \mathbf{B}(\widehat{\mathbf{s}}_{\sigma(k)}) \cdot \mathbf{a}_{\text{HA}}(k) + \mathbf{h}(\mathbf{e}(k)), \text{ with } \mathbf{e}(k) \triangleq \mathbf{s}(k) - \chi \cdot \widehat{\mathbf{s}}_{\sigma(k)}. \tag{15}$$

Next, we present a practical and common assumptions regarding the model mismatch for the design.

**Assumption 6.2.** The model mismatch in $\mathbf{h}(\cdot)$ in Equation (15) is locally Lipschitz in $\Psi_{\sigma(k)}$, i.e.,

$$(\mathbf{h}(\mathbf{e}_1) - \mathbf{h}(\mathbf{e}_2))^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot (\mathbf{h}(\mathbf{e}_1) - \mathbf{h}(\mathbf{e}_2)) \leq \kappa \cdot (\mathbf{e}_1 - \mathbf{e}_2)^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot (\mathbf{e}_1 - \mathbf{e}_2), \, \forall \mathbf{e}_1, \mathbf{e}_2 \in \Psi_{\sigma(k)}.$$

We also assume the computing hardware, mechanical components, sensors, and operating systems function correctly.

The following Theorem 6.3 presents the design, meeting Requirements 2 and 3; its proof is in Appendix B.

**Theorem 6.3 (Real-time Patch Design).** *Consider the HA-Teacher's action policy in Equation (12), the patch $\Psi_{\sigma(k)}$ in Equation (11), and the action space $\mathbb{A}$ in Equation (14), where the matrices $\widehat{\mathbf{F}}_{\sigma(k)}$ and $\widehat{\mathbf{P}}_{\sigma(k)}$ are computed according to*

$$\widehat{\mathbf{F}}_{\sigma(k)} = \widehat{\mathbf{R}}_{\sigma(k)} \cdot \widehat{\mathbf{Q}}_{\sigma(k)}^{-1}, \qquad \widehat{\mathbf{P}}_{\sigma(k)} = \widehat{\mathbf{Q}}_{\sigma(k)}^{-1}, \tag{16}$$

*with $\widehat{\mathbf{R}}_{\sigma(k)}$ and $\widehat{\mathbf{Q}}_{\sigma(k)}$ satisfying the conditions in Equations (23) and (26), and*

$$\widehat{\mathbf{Q}}_{\sigma(k)} - \mu \cdot \mathbf{P}^{-1} \succ 0, \text{ with } \mu > 0 \tag{17}$$

$$(1 - \chi \cdot \gamma_1) \cdot \mu \geq 1 - 2 \cdot \chi + \frac{\chi}{\gamma_1} > 0, \tag{18}$$

$$\begin{bmatrix} \widehat{\mathbf{Q}}_{\sigma(k)} & \widehat{\mathbf{R}}_{\sigma(k)}^\top \\ \widehat{\mathbf{R}}_{\sigma(k)} & \mathbf{T} \end{bmatrix} \succ 0, \tag{19}$$

$$\begin{bmatrix} \left(\alpha - \kappa \cdot \left(1 + \frac{1}{\gamma_2}\right)\right) \cdot \widehat{\mathbf{Q}}_{\sigma(k)} & \widehat{\mathbf{Q}}_{\sigma(k)} \cdot \mathbf{A}^\top(\widehat{\mathbf{s}}_{\sigma(k)}) + \widehat{\mathbf{R}}_{\sigma(k)}^\top \cdot \mathbf{B}^\top(\widehat{\mathbf{s}}_{\sigma(k)}) \\ \mathbf{A}(\widehat{\mathbf{s}}_{\sigma(k)}) \cdot \widehat{\mathbf{Q}}_{\sigma(k)} + \mathbf{B}(\widehat{\mathbf{s}}_{\sigma(k)}) \cdot \widehat{\mathbf{R}}_{\sigma(k)} & \frac{\widehat{\mathbf{Q}}_{\sigma(k)}}{1+\gamma_2} \end{bmatrix} \succ 0, \tag{20}$$

*where $0 < \chi < 1$, $\gamma_1 > 0$, $\gamma_2 > 0$, $0 < \alpha < 1$, and $\mathbf{P}$ is given in Equation (3). Under Assumption 6.2, we have the following properties:*

(1) The real-time patch $\Psi_{\sigma(k)} \subseteq \mathbb{X}$ holds for any time $k$.

(2) The $\mathbf{e}^\top(t+1) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t+1) \le \alpha \cdot \mathbf{e}^\top(t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t)$ holds for any time $t \in \mathbb{T}_{\sigma(k)}$ defined in Equation (8).

(3) The HA-Teacher's real-time action satisfies $\mathbf{a}_{HA}(t) \in \mathbb{A}$ for any time $t \in \mathbb{T}_{\sigma(k)}$ defined in Equation (8).

*Remark* 6.4 (**Safety Knowledge from HP-Student**). The safety regulations and envelope provided by HP-Student are applied in Equations (17) and (23) for the patch design. The resulting properties in Items 1 and 3 of Theorem 6.3 show that the designed patch meets Requirements 2 and 3. The property in Item 2 is used to develop guidance (i.e., Corollary C.1 in Appendix C) for determining $\tau$, which is the dwell time and correction horizon of HA-Teacher.

*Remark* 6.5. The $\widehat{\mathbf{F}}_{\sigma(k)}$ and $\widehat{\mathbf{P}}_{\sigma(k)}$ are automatically computed from Equations (16) to (20), (23) and (26), using the LMI toolbox [8, 22]. The computation time is quite short (0.01–0.04 seconds) and can be disregarded.

## 7 Experiment

The experiment involved comprehensive comparisons, a cart-pole system, and a real A1 quadruped robot. The open-source codes of the experiment are available at Github: https://github.com/Charlescai123/Runtime-Learning-Machine. Appendix G summarizes computation resources for implementing the runtime learning machine in the real A1 robot.

### 7.1 Cart-Pole System

We pre-train HP-Student using the OpenAI Gym [10]. The pre-training process includes domain randomization [34, 42] to bridge the Sim2Real gap, through introducing random force disturbances and randomizing the friction force. We use the simulator to mimic the real plant. The Sim2Real gap is intentionally created by inducing a friction force that is out of the distribution of those in pre-training. Unknown unknowns are disturbances applied to HP-Student's action commands, generated by a randomized Beta distribution. Appendix D explains why the randomized Beta distribution can be one kind of unknown unknown.

The system's state consists of the pendulum angle $\theta$, the cart position $x$, and their respective velocities $\omega = \dot{\theta}$ and $v = \dot{x}$. The goal of HP-Student is to stabilize system at the equilibrium $\mathbf{s}^* = [0, 0, 0, 0]^\top$, while keeping the system states within the safety set $\mathbb{X} = \{\, \mathbf{s} \mid |x| \le 1, |\theta| < 0.8 \,\}$. The action space of HA-Teacher is $\mathbb{A} = \{\, \mathbf{a}_{HA} \in \mathbb{R} \mid |\mathbf{a}_{HA}| \le 40 \,\}$. Additionally, Appendix E presents the pre-training and runtime learning configurations, and the design information of HA-Teacher and HP-Student.



(a) Initial condition 1               (b) Initial condition 2               (c) Rewards (5 seeds, %95 CI).
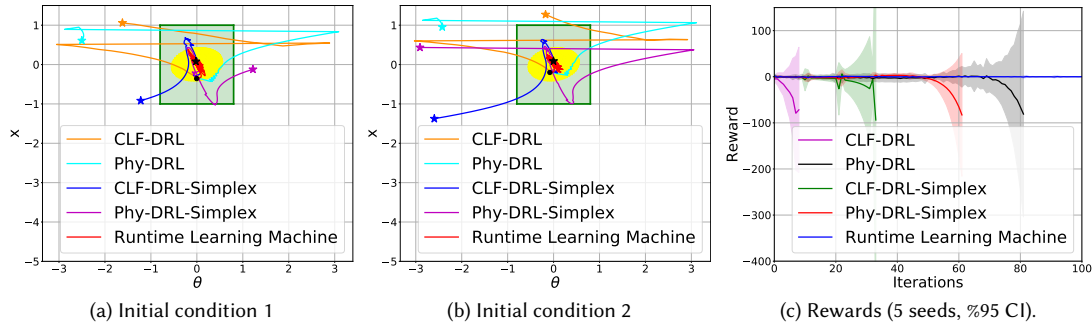
Fig. 4. (a) and (b): Phase plots in Episode 1 under two initial conditions, where the black dot and star denote the initial condition and final location, respectively. Green and yellow areas denote safety set and envelope, respectively. (c): Reward trajectories over five random seeds.

When we disable HA-Teacher's real-time patch and unsafe learning correction, our runtime learning machine degrades to the recently developed fault-tolerant DRL: runtime assurance [15] and neural Simplex [38]. Since runtime assurance is an extension of Simplex [43], we refer to the two compared models as 'CLF-DRL-Simplex' and 'Phy-DRL-Simplex', with their high-performance components being the newly developed Phy-DRL [12] and CLF-DRL [48], respectively. When HA-Teacher is completely disabled, they further degrade to pure Phy-DRL and CLF-DRL. Therefore, we now have five models for comparison. The phase plots of position and angle, as well as the trajectories of learning reward, are presented in Figure 4. It is shown in Figure 4 (a) and (b) that our runtime learning machine can assure lifetime safety in the face of unknown unknowns and the Sim2Real gap, as system states (magenta curves) never leave the safety set (green area) in any learning stage. In contrast, current fault-tolerant DRL and safe DRL cannot achieve this. Meanwhile, as seen in Figure 4 (c), our runtime learning machine provides remarkably stable and fast agent learning.

We next showcase the automatic hierarchy learning mechanism. We disable HA-Teacher in episodes 5 and 20 and observe the system's behavior under the control of the sole HP-Student. The phase plots with ten random initial conditions (each runs for 2000 steps) are displayed in Figure 5. Upon observing Figure 5 (a), we can conclude that HP-Student has successfully learned from the HA-Teacher how to ensure safety in episode 5: his action policy can confine the system states to the safety set (green area). HP-Student will automatically become independent of HA-Teacher and self-learn for a



Fig. 5. Automatic hierarchy learning.

high-performance action policy. This is evident in Figure 5 (b), where in episode 20, HP-Student consistently confines the system within her safety envelope (yellow area), and HA-Teacher is seldom triggered by the condition in Equation (7). Additionally, the action policy of HP-Student in episode 20 demonstrates higher mission performance: faster clustering and much closer proximity to the mission goal, as observed in Figure 5 (b).

We next define the HA-Teacher's activation ratio = $\frac{\text{HA-Teacher's total dwell/activation time in one episode}}{\text{one episode length}} \in (0, 1)$, where a ratio of 0 means HA-Teacher is never activated throughout the entire episode of learning, while a ratio of 1 means HA-Teacher completely dominates HP-Student for the entire episode. Figure 6 (a) illustrates the activation ratio trajectories over the episode steps during runtime learning for three different episode lengths and five random seeds. From the graph, we can conclude that HA-Teacher is rarely activated to correct the unsafe learning of HP-Student and support the safety of real plants after 15 episode-steps of runtime learning. This indicates that HP-Student has learned sufficient safety from HA-Teacher.

Finally, the trajectories of HP-Student's learning reward are shown in Figure 6 (b) and (c) for the two learning machines: one with unsafe-learning correction and one without. These trajectories were generated using the same two random initial conditions and ten seeds. Figure 6 (b) and (c) emphasizes the important role of HA-Teacher's unsafe learning correction in contributing to HP-Student's fast and stable learning, with larger reward values.

## 7.2 Real Unitree A1 Quadruped Robot

The robot's learning is to control its CoM height, CoM x-velocity, and other states to track commands $r_{v_x}$, $r_h$, and zeros, constraining system states to a safety set $\mathbb{X} = \left\{ \mathbf{s} \mid \left| \text{CoM x-velocity} - r_{v_x} \right| \leq 0.3 \text{ m/s}, \left| \text{CoM z-height} - r_h \right| \leq 0.15 \text{ m} \right\}$. HA-Teacher's action space is $\mathbb{A} = \left\{ \mathbf{a}_{\text{HA}} \in \mathbb{R}^6 \mid \left| \mathbf{a}_{\text{HA}} \right| \leq [30, 30, 30, 60, 60, 60]^\top \right\}$. The design information of HP-Student

(a) Activation ratio        (b) Initial Condition 1        (c) Initial Condition 2
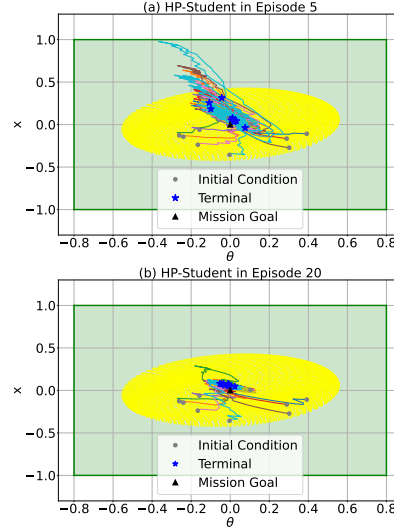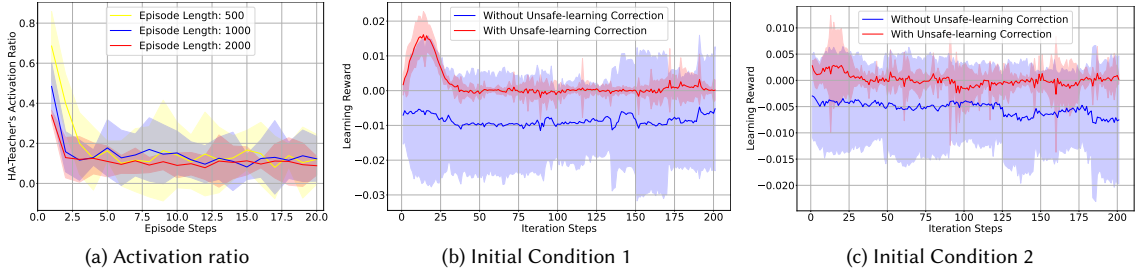
Fig. 6. (a): HA-Teacher's activation ratio over 20 episodes, five random seeds (%95 CI). (b) and (c): HP-Student's learning rewards for two runtime learning machines: one with HA-Teacher's unsafe learning correction and one without: two random initial conditions and ten seeds (%95 CI).

and HA-Teacher, and the learning configurations are presented in Appendix F. During pre-training in the simulator, we set $r_{v_x}$ = 0.6 m/s and $r_h$ = 0.24 m. To better demonstrate the runtime learning machine, the real robot's velocity command is $r_{v_x}$ = 0.35 m/s, which is different from the one in simulator. For the runtime learning, one episode is defined as "*running robot for 15 seconds*."

We compared the runtime learning machine with existing approaches to address the Sim2Real gap in training HP-Student in the simulator. The approach we compared is called 'delay + domain,' which involves concurrent delay randomization [25] and domain randomization [42] (by randomizing friction force). This approach resulted in two comparison models. 1) 'Continual Phy-DRL: delay + domain,' represents a well-trained Phy-DRL using the 'delay + domain' approach in the simulator, which performed continual learning in the real robot to fine-tune the action policy. 2) 'Phy-DRL: delay + domain,' represents the well-trained Phy-DRL policy directly deployed to the real robot. The comparison video for episode 1 is available at comparison video link [anonymous hosting and browsing] and the trajectories of the robot's CoM height and CoM-x velocity in episode 1 are shown in Figure 7. After watching the comparison video and observing Figure 7, we concluded that a well-trained Phy-DRL in the simulator cannot guarantee the safety of the real robot due to the Sim2Real gap and unknown unknowns that the delay randomization and force randomization failed to capture. In contrast, our runtime learning machine can provide the safety guarantee.

We continue the comparison with 'Continual Phy-DRL: delay + domain.' It is a well-trained Phy-DRL in the simulator and performs continual learning in the real robot for 20 episodes. Figure 8 presents the trajectories of learning



Fig. 7. Trajectories.

reward in terms of iteration steps and the episode-average reward. This demonstrates that the runtime learning machine features stable, fast, and safe learning in real plants. This notable feature is attributed to HA-Teacher's real-time patch for correcting unsafe learning and backing up safety. In addition, HA-Teacher enables HP-Student's safety-first learning from him in the learning machine. To verify this, we deactivate HA-Teacher in episodes 1 and 20, and compare system
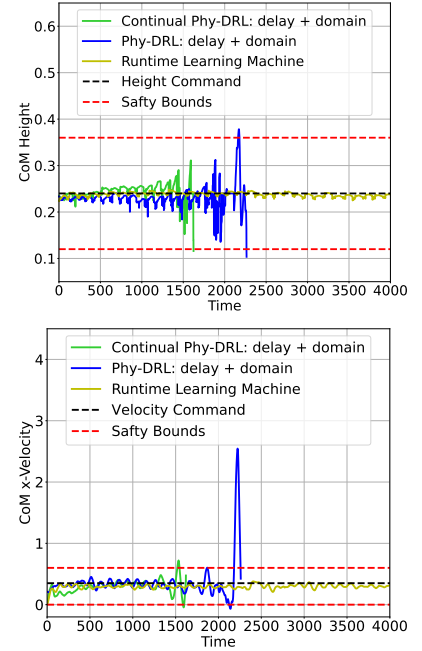
behavior. The demonstration video is available at safety-first-learning video link [anonymous hosting and browsing], which illustrates that HP-Student quickly learned from HA-Teacher to be safe, within 20 episodes (i.e., 300 seconds).

Finally, we showcase the learning machine's ability to tolerate various unknown unknowns. In addition to inherent unknowns, our experiment includes five unknown unknowns that have never occurred in HP-Student's historical training and learning. They are 1) **Beta**: Disturbances injected into HP-Student's actions, generated by a randomized Beta distribution (see Appendix D for an explanation of its representation of unknown unknowns); 2) **PD**: Random and sudden payload (around 4 lbs) drops on the robot's back; 3) **Kick**: Random and sudden kick by a human; 4) **DoS**: A real denial-of-service fault of the platform, which can be caused by task scheduling latency, communication delay, communication block, etc., but is unknown to us; and 5) **SP**: A sudden side push. We consider three combinations of these unknown unknowns applied to the runtime learning stage: i) '**Beta + PD**,' ii) '**Beta + DoS + Kick**,' and iii) '**Beta + SP**.' The demonstration video is available at unknown-unknown video link [anonymous hosting and browsing], which demonstrate that our learning machine successfully ensures the safety of the real plant by tolerating such complex combined unknowns.

### 7.3 Unitree Go2 Quadruped Robot: NVIDIA Isaac Sim

Many safety-critical CPS, such as quadruped robots, drones, UAVs, and autonomous vehicles, interact dynamically with their environments. For instance, the movement dynamics on a sandy road will be different from those on a surface covered in freezing rain. As a result, the operating environment plays a crucial role in introducing real-time unknown unknowns, Sim2Real gap, and domain gap. So, this subsection's experiment aims to demonstrate the safety assurance of our runtime learning

Fig. 8. Rewards.

machine in challenging environments. To do so, we initially pre-trained HP-Student for the A1 robot in the PyBullet simulator, using a flat terrain environment, as the same one in Section 7.2. After this pre-training, we directly deployed HP-Student to the Go2 robot. We utilized NVIDIA Isaac Sim to create an operating environment for showcasing the Go2 robot's runtime-learning capabilities. This environment transitions from flat terrain to unstructured and uneven ground, further complicated by ice from unforseen freezing rain. We here can conclude that Go2 robot's operating environment are non-stationary and unforeseen, and never occur in the pre-training stage. Besides, A1 and Go2 robots are very different in their motors, weights, heights, mass, etc. For the Go2 robot, its safety set is

$$\mathbb{X} = \left\{ \mathbf{s} \mid \left| \text{CoM x-velocity} - r_{v_x} \right| \le 0.4 \text{ m/s}, \ \left| \text{CoM z-height} - r_h \right| \le 0.15 \text{ m} \right\}. \tag{21}$$

All other designs are the same as those of A1 quadruped robot, presented in Appendix F.

In the challenging real-time operating environments, our first mission command sent to the robot is *walking forward at velocity 0.7 m/s (i.e., $r_{v_x}$ = 0.7 m/s) and maintaining CoM height at $r_h$ = 0.3 m, while constraining them to the safety set in Equation* (21). When we disable HA-Teacher's real-time patch and unsafe learning correction, our runtime learning machine degrades to the recently runtime assurance [9, 15, 44], which is also proposed to support runtime learning in real plants. When HA-Teacher is completely disabled for backing up safety, our runtime learning machine further degrade to pure Phy-DRL [11, 12].
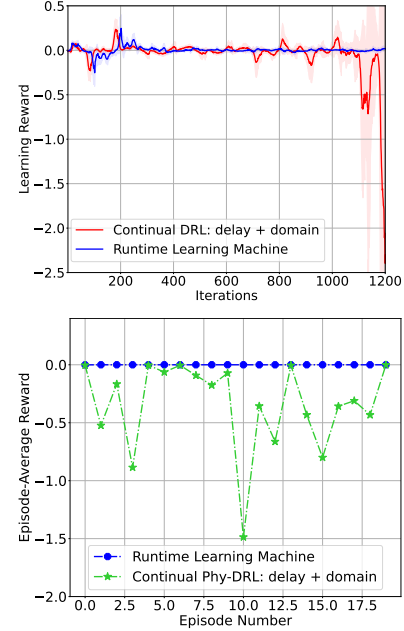
The demonstration video of the well-pretrained HP-Student (Phy-DRL) in PyBullet, along with the execution of mission command by our runtime learning machine and runtime assurance, is available at Go2 Forward [anonymous hosting and browsing]. We also set the second mission command: *walking backward at velocity 0.7 m/s (i.e., $r_{v_x}$ = -0.7 m/s) while maintaining CoM height at $r_h$ = 0.3 m, while constraining them to the safety set in Equation* (21). The demonstration video is available at Go2 Backward [anonymous hosting and browsing]. In addition to the non-stationary, unstructured, and uneven operating environments, we inject unknown-unknown noise into state samplings and randomly kick the Go2 robot to demonstrate the machine's capability of assuring safety. Following the method of inducing action noise in Section 7, the unknown-unknown noise for state samplings are generated by a randomized Beta distribution. Appendix D explains why the randomized Beta distribution generate one kind of unknown unknown. In presence of the two additional unknown unknowns, the demonstration video of robot's execution of the two mission commands (given in Section 7.3) by our runtime learning machine is available at Go2: Kick–Sensor [anonymous hosting and browsing]. All these videos well demonstrates the significantly enhanced safety assurance by our runtime learning machine in the complex environments.

## 8 Conclusion and Discussion

This paper presents a runtime learning machine designed for safety-critical CPS. The learning machine consists of the interactive HP-Student, HA-Teacher, and Coordinator. The machine's goal is to facilitate runtime learning for a high-performance action policy with verified safety in real plants, using real-time sensor data from real-time physical environments. The learning machine ensures lifetime safety by accommodating unknown unknowns and addressing the Sim2Real gap. The runtime learning machine also serves as an automatic hierarchy learning mechanism for HP-Student. Hierarchically, HP-Student first learns from the HA-Teacher to prioritize safety. After mastering safety-first learning, HP-Student autonomously self-learns to develop a high-performance action policy with a safety guarantee. Our runtime learning machine has shown outstanding features compared to state-of-the-art safe DRL and fault-tolerant DRL, with approaches to addressing the Sim2Real gap. These were demonstrated through comprehensive experiments on a cart-pole system and two quadruped robots.

## 9 Acknowledgments

## References

[1] [n. d.]. AI INCIDENT DATABASE. https://incidentdatabase.ai/entities/.

[2] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. 2017. Hindsight experience replay. *Advances in neural information processing systems* 30 (2017).

[3] Arko Banerjee, Kia Rahmani, Joydeep Biswas, and Isil Dillig. 2024. Dynamic Model Predictive Shielding for Provably Safe Reinforcement Learning. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*. https://openreview.net/forum?id=x2zY4hZcmg

[4] Thomas Bartz-Beielstein. 2019. Why we need an AI-resilient society. *arXiv:1912.08786* (2019). https://arxiv.org/pdf/1912.08786.pdf

[5] Osbert Bastani. 2021. Safe reinforcement learning with nonlinear dynamics via model predictive shielding. In *2021 American control conference*. IEEE, 3488–3494.

[6] Felix Berkenkamp, Matteo Turchetta, Angela Schoellig, and Andreas Krause. 2017. Safe model-based reinforcement learning with stability guarantees. *Advances in Neural Information Processing Systems* 30 (2017).

[7] Gerardo Bledt, Matthew J Powell, Benjamin Katz, Jared Di Carlo, Patrick M Wensing, and Sangbae Kim. 2018. MIT Cheetah 3: Design and control of a robust, dynamic quadruped robot. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2245–2252.

[8]   Stephen Boyd, Laurent El Ghaoui, Eric Feron, and Venkataramanan Balakrishnan. 1994. *Linear matrix inequalities in system and control theory.*
        SIAM.

[9]   Guillaume Brat and Ganeshmadhav Pai. 2023. Runtime assurance of aeronautical products: preliminary recommendations. *NTRS - NASA Technical
        Reports Server* (2023).

[10]  Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. OpenAI Gym.
        arXiv:arXiv:1606.01540

[11]  Hongpeng Cao, Yanbing Mao, Lui Sha, and Marco Caccamo. 2023. Physics-Model-Regulated Deep Reinforcement Learning towards Safety &
        Stability Guarantees. In *62nd IEEE Conference on Decision and Control.* 8300–8305.

[12]  Hongpeng Cao, Yanbing Mao, Lui Sha, and Marco Caccamo. 2024. Physics-Regulated Deep Reinforcement Learning: Invariant Embeddings. In *The
        Twelfth International Conference on Learning Representations.* https://openreview.net/forum?id=5Dwqu5urzs

[13]  Hongpeng Cao, Mirco Theile, Federico G. Wyrwal, and Marco Caccamo. 2022. Cloud-Edge Training Architecture for Sim-to-Real Deep Reinforcement
        Learning. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS).* 9363–9370. doi:10.1109/IROS47612.2022.9981565

[14]  Ya-Chien Chang and Sicun Gao. 2021. Stabilizing neural control using self-learned almost Lyapunov critics. In *2021 IEEE International Conference on
        Robotics and Automation.* IEEE, 1803–1809.

[15]  Shengduo Chen, Yaowei Sun, Dachuan Li, Qiang Wang, Qi Hao, and Joseph Sifakis. 2022. Runtime safety assurance for learning-enabled control of
        autonomous driving vehicles. In *2022 International Conference on Robotics and Automation (ICRA).* IEEE, 8978–8984.

[16]  Richard Cheng, Gábor Orosz, Richard M Murray, and Joel W Burdick. 2019. End-to-end safe reinforcement learning through barrier functions for
        safety-critical continuous control tasks. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 3387–3395.

[17]  Richard Cheng, Abhinav Verma, Gabor Orosz, Swarat Chaudhuri, Yisong Yue, and Joel Burdick. 2019. Control regularization for reduced variance
        reinforcement learning. In *International Conference on Machine Learning.* 1141–1150.

[18]  Jared Di Carlo, Patrick M Wensing, Benjamin Katz, Gerardo Bledt, and Sangbae Kim. 2018. Dynamic locomotion in the mit cheetah 3 through
        convex model-predictive control. In *2018 IEEE/RSJ international conference on intelligent robots and systems (IROS).* IEEE, 1–9.

[19]  Yuqing Du, Olivia Watkins, Trevor Darrell, Pieter Abbeel, and Deepak Pathak. 2021. Auto-Tuned Sim-to-Real Transfer. arXiv:2104.07662 [cs.RO]

[20]  Răzvan Florian. 2005. Correct equations for the dynamics of the cart-pole system. (08 2005).

[21]  Scott Fujimoto, Herke Hoof, and David Meger. 2018. Addressing function approximation error in actor-critic methods. In *International conference on
        machine learning.* PMLR, 1587–1596.

[22]  Pascal Gahinet, Arkadii Nemirovskii, Alan J Laub, and Mahmoud Chilali. 1994. The LMI control toolbox. In *Proceedings of 1994 33rd IEEE conference
        on decision and control*, Vol. 3. IEEE, 2038–2041.

[23]  Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement
        Learning with a Stochastic Actor. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research,
        Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, 1861–1870. https://proceedings.mlr.press/v80/haarnoja18b.html

[24]  Julian Ibarz, Jie Tan, Chelsea Finn, Mrinal Kalakrishnan, Peter Pastor, and Sergey Levine. 2021. How to train your robot with deep reinforcement
        learning: lessons we have learned. *The International Journal of Robotics Research* 40, 4-5 (2021), 698–721.

[25]  Chieko Sarah Imai, Minghao Zhang, Yuchen Zhang, Marcin Kierebiński, Ruihan Yang, Yuzhe Qin, and Xiaolong Wang. 2022. Vision-guided
        quadrupedal locomotion in the wild with multi-modal delay randomization. In *2022 IEEE/RSJ international conference on intelligent robots and
        systems (IROS).* IEEE, 5556–5563.

[26]  Tobias Johannink, Shikhar Bahl, Ashvin Nair, Jianlan Luo, Avinash Kumar, Matthias Loskyll, Juan Aparicio Ojea, Eugen Solowjow, and Sergey Levine.
        2019. Residual reinforcement learning for robot control. In *2019 International Conference on Robotics and Automation (ICRA).* IEEE, 6023–6029.

[27]  Norman L Johnson, Samuel Kotz, and Narayanaswamy Balakrishnan. 1995. *Continuous univariate distributions, volume 2.* Vol. 289. John wiley &
        sons.

[28]  Alex Kendall, Jeffrey Hawke, David Janz, Przemyslaw Mazur, Daniele Reda, John-Mark Allen, Vinh-Dieu Lam, Alex Bewley, and Amar Shah. 2019.
        Learning to drive in a day. In *2019 International Conference on Robotics and Automation.* IEEE, 8248–8254.

[29]  Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup. 2022. Towards continual reinforcement learning: A review and perspectives.
        *Journal of Artificial Intelligence Research* 75 (2022), 1401–1476.

[30]  B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. 2021. Deep reinforcement
        learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems* 23, 6 (2021), 4909–4926.

[31]  Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. 2016. End-to-end training of deep visuomotor policies. *The Journal of Machine
        Learning Research* 17, 1 (2016), 1334–1373.

[32]  Tongxin Li, Ruixiao Yang, Guannan Qu, Yiheng Lin, Steven Low, and Adam Wierman. [n. d.]. Equipping Black-Box Policies with Model-Based
        Advice for Stable Nonlinear Control. *arXiv preprint* https://arxiv.org/pdf/2206.01341.pdf.

[33]  Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. 2016. Continuous
        control with deep reinforcement learning. In *4th International Conference on Learning Representations, ICLR.*

[34]  Anusha Nagabandi, Ignasi Clavera, Simin Liu, Ronald S. Fearing, Pieter Abbeel, Sergey Levine, and Chelsea Finn. 2019. Learning to Adapt in
        Dynamic, Real-World Environments Through Meta-Reinforcement Learning. arXiv:1803.11347 [cs.LG]

[35]  NHTSA. 2022. Summary Report: Standing General Order on Crash Reporting for Level 2 Advanced Driver Assistance Systems. *National Highway
        Traffic Safety Administration* (2022). https://www.nhtsa.gov/sites/nhtsa.gov/files/2022-06/ADAS-L2-SGO-Report-June-2022.pdf

[36] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. 2018. Sim-to-Real Transfer of Robotic Control with Dynamics Randomization. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE. doi:10.1109/icra.2018.8460528

[37] Theodore J Perkins and Andrew G Barto. 2002. Lyapunov design for safe reinforcement learning. *Journal of Machine Learning Research* 3, Dec (2002), 803–832.

[38] Dung T Phan, Radu Grosu, Nils Jansen, Nicola Paoletti, Scott A Smolka, and Scott D Stoller. 2020. Neural Simplex architecture. In *NASA Formal Methods Symposium*. Springer, 97–114.

[39] Rajesh Rajamani. 2011. *Vehicle dynamics and control.* Springer Science & Business Media.

[40] Krishan Rana, Vibhavari Dasagi, Jesse Haviland, Ben Talbot, Michael Milford, and Niko Sünderhauf. [n. d.]. Bayesian controller fusion: Leveraging control priors in deep reinforcement learning for robotics. *arXiv preprint* https://arxiv.org/pdf/2107.09822.pdf.

[41] Jan Roskam. 1995. *Airplane flight dynamics and automatic flight controls.* DARcorporation.

[42] Fereshteh Sadeghi and Sergey Levine. 2017. CAD2RL: Real Single-Image Flight without a Single Real Image. arXiv:1611.04201 [cs.LG]

[43] Lui Sha et al. 2001. Using simplicity to control complexity. *IEEE Software* 18, 4 (2001), 20–28.

[44] Joseph Sifakis and David Harel. 2023. Trustworthy autonomous system development. *ACM Transactions on Embedded Computing Systems* 22, 3 (2023), 1–24.

[45] Jie Tan, Tingnan Zhang, Erwin Coumans, Atil Iscen, Yunfei Bai, Danijar Hafner, Steven Bohez, and Vincent Vanhoucke. 2018. Sim-to-Real: Learning Agile Locomotion For Quadruped Robots. *Robotics: Science and Systems* (2018).

[46] Tian Tolentino. 2019. Autonomous aircraft market worth USD 23.7bn by 2030. https://www.traveldailymedia.com/autonomous-aircraft-market-research/.

[47] Quan Vuong, Sharad Vikram, Hao Su, Sicun Gao, and Henrik I. Christensen. 2019. How to pick the domain randomization parameters for sim-to-real transfer of reinforcement learning policies? arXiv:1903.11774 [cs.LG]

[48] Tyler Westenbroek, Fernando Castaneda, Ayush Agrawal, Shankar Sastry, and Koushil Sreenath. 2022. Lyapunov Design for Robust and Efficient Robotic Reinforcement Learning. *arXiv:2208.06721* (2022). https://arxiv.org/pdf/2208.06721.pdf

[49] Tsung-Yen Yang, Tingnan Zhang, Linda Luu, Sehoon Ha, Jie Tan, and Wenhao Yu. 2022. Safe reinforcement learning for legged locomotion. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2454–2461.

[50] Yuxiang Yang. [n. d.]. GitHub: Quadruped Robot Simulator. https://github.com/yxyang/locomotion_simulation

[51] Wenhao Yu, Jie Tan, C Karen Liu, and Greg Turk. 2017. Preparing for the unknown: Learning a universal policy with online system identification. *Robotics: Science and Systems* (2017).

[52] Arnold Zachary and Toner Helen. 2021. AI Accidents: An Emerging Threat. *Center for Security and Emerging Technology* (2021). https://doi.org/10.51593/20200072

[53] Fuzhen Zhang. 2006. *The Schur complement and its applications.* Vol. 4. Springer Science & Business Media.

[54] Liqun Zhao, Konstantinos Gatsis, and Antonis Papachristodoulou. 2023. Stable and Safe Reinforcement Learning via a Barrier-Lyapunov Actor-Critic Approach. In *62nd IEEE Conference on Decision and Control*. IEEE, 1320–1325.

## A  Auxiliary Lemmas

This section introduces the auxiliary lemmas used to establish the theoretical framework for our proposed runtime learning machine.

**Lemma A.1** (Schur Complement [53]). *For any symmetric matrix* $\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{C} \end{bmatrix}$, *then* $\mathbf{M} \succ 0$ *holds if and only if* $\mathbf{C} \succ 0$ *and* $\mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^\top \succ 0$.

**Lemma A.2** ([12]). *Consider the safety set* $\mathbb{X}$ *defined in Equation* (2) *and define a set*

$$\Omega_{\sigma(k)} \triangleq \{\, \mathbf{s} \mid \mathbf{s}^\top \cdot \widehat{\mathbf{Q}}_{\sigma(k)}^{-1} \cdot \mathbf{s} \le 1, \ \widehat{\mathbf{Q}}_{\sigma(k)} \succ 0 \,\}. \tag{22}$$

*We have* $\Omega_{\sigma(k)} \subseteq \mathbb{X}$ *if*

$$[\underline{\mathbf{D}}]_{i,:} \cdot \widehat{\mathbf{Q}}_{\sigma(k)} \cdot [\underline{\mathbf{D}}^\top]_{:,i} = \begin{cases} \ge 1, & [\mathbf{d}]_i = 1 \\ \le 1, & [\mathbf{d}]_i = -1 \end{cases}, \ and \ [\overline{\mathbf{D}}]_{i,:} \cdot \widehat{\mathbf{Q}}_{\sigma(k)} \cdot [\overline{\mathbf{D}}^\top]_{:,i} \le 1, \ i \in \{1, \ldots, h\} \tag{23}$$

where $\overline{\mathbf{D}} = \frac{\mathbf{D}}{\Lambda}$, $\underline{\mathbf{D}} = \frac{\mathbf{D}}{\underline{\Lambda}}$, and for $i, j \in \{1, \ldots, h\}$, we define:

$$
[\mathbf{d}]_i \triangleq \begin{cases} 1, & [\underline{\mathbf{v}}]_i > 0 \\ 1, & [\overline{\mathbf{v}}]_i < 0 \\ -1, & \text{otherwise} \end{cases}, \quad [\overline{\Lambda}]_{i,j} \triangleq \begin{cases} 0, & i \neq j \\ [\overline{\mathbf{v}}]_i, & [\underline{\mathbf{v}}]_i > 0 \\ [\underline{\mathbf{v}}]_i, & [\overline{\mathbf{v}}]_i < 0 \\ [\overline{\mathbf{v}}]_i, & \text{otherwise} \end{cases}, \quad [\underline{\Lambda}]_{i,j} \triangleq \begin{cases} 0, & i \neq j \\ [\underline{\mathbf{v}}]_i, & [\underline{\mathbf{v}}]_i > 0 \\ [\overline{\mathbf{v}}]_i, & [\overline{\mathbf{v}}]_i < 0 \\ -[\underline{\mathbf{v}}]_i, & \text{otherwise} \end{cases}. \tag{24}
$$

**Lemma A.3.** *Consider the action set $\mathbb{A}$ defined in Equation (14), and*

$$
\Xi \triangleq \left\{ \mathbf{a}_{HA} \in \mathbb{R}^m \mid \mathbf{a}_{HA}^\top \cdot \mathbf{T}^{-1} \cdot \mathbf{a}_{HA} \leq 1, \ \mathbf{T} \succ 0 \right\}. \tag{25}
$$

*We have $\Xi \subseteq \mathbb{A}$, if*

$$
[\underline{\mathbf{C}}]_{i,:} \cdot \mathbf{T} \cdot [\underline{\mathbf{C}}^\top]_{:,i} = \begin{cases} \geq 1, & [\mathbf{c}]_i = 1 \\ \leq 1, & [\mathbf{c}]_i = -1 \end{cases}, \text{ and } [\overline{\mathbf{C}}]_{i,:} \cdot \mathbf{T} \cdot [\overline{\mathbf{C}}^\top]_{:,i} \leq 1, \ i \in \{1, \ldots, m\} \tag{26}
$$

*where $\overline{\mathbf{C}} = \frac{\mathbf{C}}{\Lambda}$ and $\underline{\mathbf{C}} = \frac{\mathbf{C}}{\underline{\Lambda}}$, and for $i, j \in \{1, \ldots, m\}$, we define:*

$$
[\mathbf{c}]_i \triangleq \begin{cases} 1, & [\underline{\mathbf{z}}]_i > 0 \\ 1, & [\overline{\mathbf{z}}]_i < 0 \\ -1, & \text{otherwise} \end{cases}, \quad [\overline{\Lambda}]_{i,j} \triangleq \begin{cases} 0, & i \neq j \\ [\overline{\mathbf{z}}]_i, & [\underline{\mathbf{z}}]_i > 0 \\ [\underline{\mathbf{z}}]_i, & [\overline{\mathbf{z}}]_i < 0 \\ [\overline{\mathbf{z}}]_i, & \text{otherwise} \end{cases}, \quad [\underline{\Lambda}]_{i,j} \triangleq \begin{cases} 0, & i \neq j \\ [\underline{\mathbf{z}}]_i, & [\underline{\mathbf{z}}]_i > 0 \\ [\overline{\mathbf{z}}]_i, & [\overline{\mathbf{z}}]_i < 0 \\ -[\underline{\mathbf{z}}]_i, & \text{otherwise} \end{cases}. \tag{27}
$$

PROOF. Lemma A.3's proof path is exactly the same as the proof of Lemma B.2 in [12], so it is omitted here. □

**Lemma A.4.** *For two vectors $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^n$, and a matrix $\mathbf{P} \succ 0$, we have*

$$
2 \cdot \mathbf{x}^\top \cdot \mathbf{P} \cdot \mathbf{y} \leq \gamma \cdot \mathbf{x}^\top \cdot \mathbf{P} \cdot \mathbf{x} + \frac{1}{\gamma} \cdot \mathbf{y}^\top \cdot \mathbf{P} \cdot \mathbf{y}, \text{ with } \gamma > 0.
$$

PROOF. The proof is straightforward when we consider $\mathbf{P} \succ 0$ and recall the following inequality:

$$
(\sqrt{\gamma} \cdot \mathbf{x} - \frac{1}{\sqrt{\gamma}} \cdot \mathbf{y})^\top \cdot \mathbf{P} \cdot (\sqrt{\gamma} \cdot \mathbf{x} - \frac{1}{\sqrt{\gamma}} \cdot \mathbf{y}) = \gamma \cdot \mathbf{x}^\top \cdot \mathbf{P} \cdot \mathbf{x} + \frac{1}{\gamma} \cdot \mathbf{y}^\top \cdot \mathbf{P} \cdot \mathbf{y} - 2 \cdot \mathbf{x}^\top \cdot \mathbf{P} \cdot \mathbf{y} \geq 0.
$$

□

## B  Proof of Theorem 6.3

The three statements in Theorem 6.3 are proved separately.

*B.0.1 Proof of Statement in Item 1.* The envelope patch in Equation (11) can be equivalently rewritten as:

$$
\Psi_{\sigma(k)} = \Big\{ \mathbf{s} \mid \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s} \leq (1 - \chi)^2 \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k) + 2 \cdot \chi \cdot \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \widehat{\mathbf{s}}_{\sigma(k)}
$$
$$
- \chi^2 \cdot \widehat{\mathbf{s}}_{\sigma(k)}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \widehat{\mathbf{s}}_{\sigma(k)} \Big\}, \tag{28}
$$

which, in light of Equation (13), equivalently transforms to

$$\Psi_{\sigma(k)} = \left\{ \mathbf{s} \mid \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s} \le (1 - 2 \cdot \chi) \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k) + 2 \cdot \chi \cdot \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \widehat{\mathbf{s}}(k) \right\}. \tag{29}$$

In light of Lemma A.4 in Appendix A, we have

$$2 \cdot \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k) \le \gamma_1 \cdot \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s} + \frac{1}{\gamma_1} \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k), \text{ with } \gamma_1 > 0$$

substituting which into the inequality in Equation (29) and considering $0 < \chi < 1$ yields

$$\mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s} \le (1 - 2 \cdot \chi) \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k) + 2 \cdot \chi \cdot \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \widehat{\mathbf{s}}(k)$$

$$\le (1 - 2 \cdot \chi + \frac{\chi}{\gamma_1}) \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k) + \chi \cdot \gamma_1 \cdot \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}, \tag{30}$$

which leads to

$$(1 - \chi \cdot \gamma_1) \cdot \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s} \le (1 - 2 \cdot \chi + \frac{\chi}{\gamma_1}) \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k). \tag{31}$$

We conclude from Equations (29) to (31) that if the inequality for defining the envelope patch $\Psi_{\sigma(k)}$ in Equation (29) holds, the inequality in Equation (31) holds as well. Therefore, we can define the first auxiliary set:

$$\Theta_1 = \left\{ \mathbf{s} \mid (1 - \chi \cdot \gamma_1) \cdot \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s} \le (1 - 2 \cdot \chi + \frac{\chi}{\gamma_1}) \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k) \right\}, \tag{32}$$

and it satisfies

$$\Psi_{\sigma(k)} \subseteq \Theta_1. \tag{33}$$

Considering $\mu > 0$, we can conclude from Equation (18) that $1 - \chi \cdot \gamma > 0$. Therefore, the set in Equation (32) can be equivalently transformed to

$$\Theta_1 = \left\{ \mathbf{s} \mid \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s} \le \frac{1 - 2 \cdot \chi + \frac{\chi}{\gamma_1}}{1 - \chi \cdot \gamma_1} \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k) \right\}. \tag{34}$$

Considering $\widehat{\mathbf{P}}_{\sigma(k)} = \widehat{\mathbf{Q}}_{\sigma(k)}^{-1}$ and $\mu > 0$, the condition in Equation (17) is equivalent to

$$\frac{1}{\mu} \cdot \mathbf{P} \succ \widehat{\mathbf{P}}_{\sigma(k)},$$

substituting which into the inequality in Equation (34) results in

$$\mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s} \le \frac{1 - 2 \cdot \chi + \frac{\chi}{\gamma_1}}{1 - \chi \cdot \gamma_1} \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k) \le \frac{1 - 2 \cdot \chi + \frac{\chi}{\gamma_1}}{1 - \chi \cdot \gamma_1} \cdot \frac{1}{\mu} \cdot \mathbf{s}^\top(k) \cdot \mathbf{P} \cdot \mathbf{s}(k) = \frac{1 - 2 \cdot \chi + \frac{\chi}{\gamma_1}}{1 - \chi \cdot \gamma_1} \cdot \frac{1}{\mu}, \tag{35}$$

where last equality is obtained because $\mathbf{s}(k)$ approaches the boundary of the safety envelope, i.e., $\mathbf{s}^\top(k) \cdot \mathbf{P} \cdot \mathbf{s}(k) = 1$. From this, we can conclude that if the inequality defining the set $\Theta_1$ in Equation (34) holds, the inequality in Equation (35) holds as well. Therefore, we can define the second auxiliary set as:

$$\Theta_2 = \left\{ \mathbf{s} \mid \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s} \le \frac{1 - 2 \cdot \chi + \frac{\chi}{\gamma_1}}{1 - \chi \cdot \gamma_1} \cdot \frac{1}{\mu} \right\}, \tag{36}$$

and it satisfies

$$\Theta_1 \subseteq \Theta_2. \tag{37}$$

Moving forward, we note that the condition in Equation (18) is equivalent to $0 < \frac{1-2\chi+\frac{\chi}{\gamma_1}}{1-\chi\cdot\gamma_1} \cdot \frac{1}{\mu} < 1$. Therefore, we can define the third auxiliary set as

$$\Theta_3 = \left\{ \mathbf{s} \mid \mathbf{s}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s} \leq 1, \ \widehat{\mathbf{P}}_{\sigma(k)} \succ 0 \right\}, \tag{38}$$

and referring to Equation (36), it satisfies

$$\Theta_2 \subseteq \Theta_3. \tag{39}$$

At the moment, we can draw conclusions from Equations (33), (37) and (40):

$$\Psi_{\sigma(k)} \subseteq \Theta_1 \subseteq \Theta_2 \subseteq \Theta_3. \tag{40}$$

Observing Equations (22) and (38), we have $\Theta_3 = \Omega_{\sigma(k)}$. Then, applying Lemma A.2 in Appendix A, we have $\Theta_3 = \Omega_{\sigma(k)} \subseteq \mathbb{X}$, which, in light of Equation (40), results in $\Psi_{\sigma(k)} \subseteq \mathbb{X}$. We thus complete the proof of the statement in Item 1.

*B.0.2 Proof of Statement in Item 2.* We define a Lyapunov candidate for the tracking-error dynamics described in Equation (15) as:

$$V(t) = \mathbf{e}^\top (t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e} (t), \tag{41}$$

which, combined with the dynamics in Equation (15) and the action policy in Equation (12), leads to

$$
\begin{aligned}
V(t+1) &- \alpha \cdot V(t) \\
&= \mathbf{e}^\top (t+1) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e} (t+1) - \alpha \cdot \mathbf{e}^\top (t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e} (t) \\
&= \mathbf{e}^\top (t) \cdot \left( \overline{\mathbf{A}}^\top (\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \overline{\mathbf{A}}(\widehat{\mathbf{s}}_{\sigma(t)}) - \alpha \cdot \widehat{\mathbf{P}}_{\sigma(t)} \right) \cdot \mathbf{e} (t) + \mathbf{h}^\top (\mathbf{e} (t)) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{h} (\mathbf{e} (t)) \\
&\qquad\qquad\qquad\qquad\qquad\qquad + 2 \cdot \mathbf{e}^\top (t) \cdot \left( \overline{\mathbf{A}}(\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \right) \cdot \mathbf{h} (\mathbf{e} (t)),
\end{aligned}
\tag{42}
$$

where we define:

$$\overline{\mathbf{A}}(\widehat{\mathbf{s}}_{\sigma(t)}) \stackrel{\Delta}{=} \mathbf{A}(\widehat{\mathbf{s}}_{\sigma(t)}) + \mathbf{B}(\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{F}}_{\sigma(t)}. \tag{43}$$

After applying Lemma A.4 in Appendix A, we have:

$$
\begin{aligned}
2\mathbf{e}^\top (t) \cdot \left( \overline{\mathbf{A}}(\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \right) \cdot \mathbf{h} (\mathbf{e} (t)) &\leq \gamma_2 \cdot \mathbf{e}^\top (t) \cdot \overline{\mathbf{A}}^\top (\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \overline{\mathbf{A}}(\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \mathbf{e} (t) \\
&\qquad\qquad + \frac{1}{\gamma_2} \cdot \mathbf{h}^\top (\mathbf{e} (t)) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{h} (\mathbf{e} (t)),
\end{aligned}
\tag{44}
$$

where $\gamma_2 > 0$.

We note that Assumption 6.2 implies:

$$\mathbf{h}^\top (\mathbf{e} (t)) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{h}(\mathbf{e} (t)) \leq \kappa \cdot \mathbf{e}^\top (t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e} (t). \tag{45}$$

Substituting inequalities in Equations (44) and (45) into Equation (42) yields:

$$V(t+1) - \alpha \cdot V(t) \leq \mathbf{e}^\top (t) \cdot \left( (1+\gamma_2) \cdot \overline{\mathbf{A}}^\top (\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \overline{\mathbf{A}}(\widehat{\mathbf{s}}_{\sigma(t)}) - (\alpha - \kappa \cdot (1 + \frac{1}{\gamma_2})) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \right) \cdot \mathbf{e} (t). \tag{46}$$

Recalling the Schur Complement in Lemma A.1 of Appendix A and considering $\widehat{\mathbf{P}}_{\sigma(t)} \succ 0$, we conclude that the inequality in Equation (20) is equivalent to

$$(\alpha - \kappa \cdot (1 + \frac{1}{\gamma_2})) \cdot \widehat{\mathbf{Q}}_{\sigma(t)} - (1 + \gamma_2) \cdot (\widehat{\mathbf{Q}}_{\sigma(t)} \cdot \mathbf{A}^\top(\widehat{\mathbf{s}}_{\sigma(t)}) + \widehat{\mathbf{R}}^\top_{\sigma(t)} \cdot \mathbf{B}^\top(\widehat{\mathbf{s}}_{\sigma(t)}) ) \cdot \widehat{\mathbf{Q}}^{-1}_{\sigma(t)} \cdot (\mathbf{A}(\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{Q}}_{\sigma(t)} + \mathbf{B}(\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{R}}_{\sigma(t)}) \succ 0,$$

multiplying both the left-hand side and the right-hand side of which by $\widehat{\mathbf{Q}}^{-1}$ yields:

$$(\alpha - \kappa \cdot (1 + \frac{1}{\gamma_2})) \cdot \widehat{\mathbf{Q}}^{-1}_{\sigma(t)} - (1 + \gamma_2) \cdot (\mathbf{A}^\top(\widehat{\mathbf{s}}_{\sigma(t)}) + \widehat{\mathbf{Q}}^{-1}_{\sigma(t)} \cdot \widehat{\mathbf{R}}^\top_{\sigma(t)} \cdot \mathbf{B}^\top(\widehat{\mathbf{s}}_{\sigma(t)})) \cdot \widehat{\mathbf{Q}}^{-1}_{\sigma(t)}$$
$$\cdot (\mathbf{A}(\widehat{\mathbf{s}}_{\sigma(k)}) + \mathbf{B}(\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{R}}_{\sigma(t)} \cdot \mathbf{Q}^{-1}_{\sigma(t)}) \succ 0,$$

Substituting the definitions in Equation (16) into which, we arrive at

$$(\alpha - \kappa \cdot (1 + \frac{1}{\gamma_2})) \cdot \widehat{\mathbf{P}}_{\sigma(t)} - (1 + \gamma_2) \cdot (\mathbf{A}^\top(\widehat{\mathbf{s}}_{\sigma(t)}) + \widehat{\mathbf{F}}^\top_{\sigma(t)} \cdot \mathbf{B}^\top(\widehat{\mathbf{s}}_{\sigma(t)})) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot (\mathbf{A}(\widehat{\mathbf{s}}_{\sigma(t)}) + \mathbf{B}(\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{F}}_{\sigma(t)}) \succ 0. \quad (47)$$

Recalling Equation (43), the inequality in Equation (47) is equivalent to the following:

$$(1 + \gamma_2) \cdot \overline{\mathbf{A}}^\top(\widehat{\mathbf{s}}_{\sigma(t)}) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \overline{\mathbf{A}}(\widehat{\mathbf{s}}_{\sigma(t)}) - (\alpha - \kappa \cdot (1 + \frac{1}{\gamma_2})) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \prec 0,$$

which, in conjunction with Equation (46), leads to $V(t+1) - \alpha \cdot V(t) \leq 0$, i.e., $\mathbf{e}^\top(t+1) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t+1) \leq \alpha \cdot \mathbf{e}^\top(t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t)$, we thus complete the proof of the statement in Item 2.

*B.0.3  Proof of Statement in Item 3.* With the consideration of $\mathbf{T}^{-1} = \mathbf{V}$, according to Lemma A.1, the condition in Equation (19) implies:

$$\widehat{\mathbf{Q}}_{\sigma(t)} - \widehat{\mathbf{R}}^\top_{\sigma(t)} \cdot \mathbf{T}^{-1} \cdot \widehat{\mathbf{R}}_{\sigma(t)} = \widehat{\mathbf{Q}}_{\sigma(t)} - \widehat{\mathbf{R}}^\top_{\sigma(t)} \cdot \mathbf{V} \cdot \widehat{\mathbf{R}}_{\sigma(t)} \succ 0. \quad (48)$$

Substituting $\widehat{\mathbf{F}}_{\sigma(t)} \cdot \widehat{\mathbf{Q}}_{\sigma(t)} = \widehat{\mathbf{R}}_{\sigma(t)}$ into Equation (48) leads to

$$\widehat{\mathbf{Q}}_{\sigma(t)} - (\widehat{\mathbf{F}}_{\sigma(t)} \cdot \widehat{\mathbf{Q}}_{\sigma(t)})^\top \cdot \mathbf{V} \cdot (\widehat{\mathbf{F}}_{\sigma(t)} \cdot \widehat{\mathbf{Q}}_{\sigma(t)}) \succ 0. \quad (49)$$

multiplying both left-hand and right-hand sides of which by $\widehat{\mathbf{Q}}^{-1}_{\sigma(k)}$ yields:

$$\widehat{\mathbf{Q}}^{-1}_{\sigma(t)} - \widehat{\mathbf{F}}^\top_{\sigma(t)} \cdot \mathbf{V} \cdot \widehat{\mathbf{Q}}_{\sigma(t)} \succ 0,$$

from which we thus have

$$\mathbf{e}^\top(t) \cdot \widehat{\mathbf{Q}}^{-1}_{\sigma(t)} \cdot \mathbf{e}(t) - \mathbf{e}^\top(t) \cdot \widehat{\mathbf{F}}^\top_{\sigma(t)} \cdot \mathbf{V} \cdot \widehat{\mathbf{F}}_{\sigma(t)} \cdot \mathbf{e}(t) = \mathbf{e}^\top(t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t) - \mathbf{a}^\top_{\mathrm{HA}}(t) \cdot \mathbf{V} \cdot \mathbf{a}_{\mathrm{HA}}(t) > 0, \quad (50)$$

which is obtained via considering $\widehat{\mathbf{P}}_{\sigma(t)} = \widehat{\mathbf{Q}}^{-1}_{\sigma(t)}$, and Equation (12) with $\mathbf{e}(t) = \mathbf{s}(t) - \chi \cdot \widehat{\mathbf{s}}_{\sigma(t)}$.

We let $\mathbf{e} = \mathbf{s} - \chi \cdot \widehat{\mathbf{s}}_{\sigma(k)}$. The patch definition in Equation (11) can re-expressed as

$$\Psi_{\sigma(k)} \triangleq \{ \mathbf{e} \mid \mathbf{e}^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{e} \leq (1 - \chi)^2 \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k), \text{ with } \mathbf{s}(k) \text{ subject to Equation (7), and } \widehat{\mathbf{P}}_{\sigma(k)} \succ 0 \}. \quad (51)$$

The inequality in Equation (50) can be expressed as $\mathbf{e}^\top(t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t) > \mathbf{a}^\top_{\mathrm{HA}}(t) \cdot \mathbf{V} \cdot \mathbf{a}_{\mathrm{HA}}(t)$. Based on Equation (38), we can conclude that if $\mathbf{e}(t) \in \Theta_3$, meaning it satisfies $\mathbf{e}^\top(t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t) < 1$, then $\mathbf{a}^\top_{\mathrm{HA}}(t) \cdot \mathbf{V} \cdot \mathbf{a}_{\mathrm{HA}}(t) < 1$. Additionally, considering Equation (40) and Equation (51), if $\mathbf{e}(t) \in \Psi_{\sigma(k)}$, then $\mathbf{a}^\top_{\mathrm{HA}}(t) \cdot \mathbf{V} \cdot \mathbf{a}_{\mathrm{HA}}(t) < 1$. It's important to note that $t \in \{k, \ldots, k + \tau\} = \mathbb{T}_{\sigma(k)}$, and $k$ represents the triggering time of HA-Teacher.

Upon verification from Equation (13), it becomes evident that $\mathbf{e}(k) = \mathbf{s}(k) - \chi \cdot \widehat{\mathbf{s}}_{\sigma(k)} = \mathbf{s}(k) - \chi \cdot \widehat{\mathbf{s}}(k)$, where $\mathbf{e}(k)$ lies on the boundary of the patch: $\mathbf{e}(k)^\top \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{e}(k) = (1 - \chi)^2 \cdot \mathbf{s}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{s}(k)$. Furthermore, as per the

second statement in Item 2, i.e., $\mathbf{e}^\top(t+1) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t+1) \le \alpha \cdot \mathbf{e}^\top(t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t)$ for time $t \in \mathbb{T}_{\sigma(k)}$, we can infer that $\mathbf{e}(t)$ never exits the patch during the active time of HA-Teacher initiated at time $k$. Hence, we can conclude that $\mathbf{a}_{\mathrm{HA}}^\top(t) \cdot \mathbf{V} \cdot \mathbf{a}_{\mathrm{HA}}(t) < 1$ holds for any time $t \in \mathbb{T}_{\sigma(k)}$.

Finally, taking into account Equation (26) and Lemma A.3 in Appendix A, we can establish that $\mathbf{a}_{\mathrm{HA}}(t) \in \mathbb{A}$ for any time $t \in \mathbb{T}_{\sigma(k)}$, thus completing the proof.

## C Guidance for Correction Horizon and Dwell Time

Upon reviewing Figure 1 and Equations (8) to (10), we can conclude that $\tau$ serves as both the correction horizon for the unsafe actions of HP-Student and the dwell time for HA-Teacher. The value of $\tau$ significantly influences HP-Student's runtime learning in achieving a high-performance policy. If $\tau$ is small, the patch center fails to attract the system states to the envelope inside, resulting in HA-Teacher dominating the learning process, solely ensuring safety. Conversely, if $\tau$ is very large, HP-Student is unable to effectively and swiftly self-learn to achieve his goal. Thus, these considerations should guide the selection of $\tau$. This guidance is based on the results from Theorem 6.3, as presented in the following corollary.

**Corollary C.1.** *If the correction horizon and the dwell time of HA-Teacher, denoted as $\tau$, satisfy:*

$$\tau = \left\lceil \frac{\ln(\delta \cdot \mu) - \ln(\mathbf{e}^\top(k) \cdot \mathbf{P} \cdot \mathbf{e}(k))}{\ln \alpha} \right\rceil, \tag{52}$$

*we have $\mathbf{e}^\top(k+\tau) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{e}(k+\tau) \le \delta$, where $k$ denotes the triggering time of HA-Teacher.*

Proof. We obtain from Item 2 that

$$\mathbf{e}^\top(t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t) \le \alpha^{t-k} \cdot \mathbf{e}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{e}(k), \ t \in \mathbb{T}_{\sigma(k)}. \tag{53}$$

Considering $0 < \alpha < 1$, we can verify from Equation (53) that $\alpha^{t-k} \cdot \mathbf{e}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{e}(k) \le \delta$ is equivalent to

$$\tau = t - k \ge \frac{\ln \delta - \ln(\mathbf{e}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{e}(k))}{\ln \alpha}. \tag{54}$$

In addition, considering $\widehat{\mathbf{P}}_{\sigma(k)} = \widehat{\mathbf{Q}}_{\sigma(k)}^{-1}$ and $\mu > 0$, the condition in Equation (17) used for designing real-time patch in Theorem 6.3 is equivalent to

$$\frac{1}{\mu} \cdot \mathbf{P} \succ \widehat{\mathbf{P}}_{\sigma(k)},$$

which, in conjunction with $0 < \alpha < 1$, leads to

$$\frac{\ln \delta - \ln(\mathbf{e}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{e}(k))}{\ln \alpha} \le \frac{\ln \delta - \ln(\frac{1}{\mu} \cdot \mathbf{e}^\top(k) \cdot \mathbf{P} \cdot \mathbf{e}(k))}{\ln \alpha} = \frac{\ln(\delta \cdot \mu) - \ln(\mathbf{e}^\top(k) \cdot \mathbf{P} \cdot \mathbf{e}(k))}{\ln \alpha}$$
$$\le \left\lceil \frac{\ln(\delta \cdot \mu) - \ln(\mathbf{e}^\top(k) \cdot \mathbf{P} \cdot \mathbf{e}(k))}{\ln \alpha} \right\rceil. \tag{55}$$

Based on Equation (55), we can conclude that if the condition in Equation (52) is satisfied, then the inequality in Equation (54) also holds. Consequently, we have $\alpha^{t-k} \cdot \mathbf{e}^\top(k) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{e}(k) \le \delta$. This, together with Equation (53), implies $\mathbf{e}^\top(t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t) \le \delta$. Furthermore, if we consider $\sigma(k)$ as a piece-wise signal (i.e., $\sigma(m) = \sigma(k)$ for $m \in \mathbb{T}_{\sigma(k)} = \{k, k+1, \ldots, k+\tau\}$) and $\tau = t-k$, then $\mathbf{e}^\top(t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t) \le \delta$ can be rewritten as $\mathbf{e}^\top(k+\tau) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{e}(k+\tau) \le \delta$. This completes the proof. □

The real-time tracking error $\mathbf{e}(t)$ represents the distance to the patch center $\chi \cdot \widehat{\mathbf{s}}_{\sigma(k)}$. Therefore, $\mathbf{e}^\top(t) \cdot \widehat{\mathbf{P}}_{\sigma(t)} \cdot \mathbf{e}(t)$ serves as a measurement metric for proximity to the patch center. Additionally, $\mathbf{e}^\top(k+\tau) \cdot \widehat{\mathbf{P}}_{\sigma(k)} \cdot \mathbf{e}(k+\tau) \le \delta$ can be interpreted as a safety criterion for returning to HP-Student. Consequently, Corollary C.1 implies that $\tau$ is computed to ensure that the real plant, under the control of HA-Teacher, satisfies the preset safety criteria rather than infinitely approaching the patch center.

## D Unknown Unknown: Randomized Beta Distribution

In the real world, plants can encounter a multitude of unknown variables, each with unique characteristics. To tackle this challenge, we propose utilizing a variant of the Beta distribution [27] to effectively model one type of these unknowns. This approach holds promise in mathematically defining and addressing these uncertainties.

**Definition D.1** (Randomized Beta Distribution). The disturbance, noise, or fault, denoted by $\mathbf{d}(k)$, is considered to be a bounded unknown if (i) $\mathbf{d}(k) \sim Beta(\alpha(k), \beta(k), c, a)$, and (ii) $\alpha(k)$ and $\beta(k)$ are random parameters. In other words, the disturbance $\mathbf{d}(k)$ is within the range of [a, c], and its probability density function (pdf) is given by

$$f(\mathbf{d}(k); \alpha(k), \beta(k), a, c) = \frac{(\mathbf{d}(k)-a)^{\alpha-1}(c-\mathbf{d}(k))^{\alpha(k)-1}\Gamma(\alpha(k)+\beta(k))}{(c-a)^{\alpha(k)+\beta(k)-1}\Gamma(\alpha(k))\Gamma(\beta(k))}, \tag{56}$$

where $\Gamma(\alpha(k)) = \int_0^\infty t^{\alpha(k)-1}e^{-t}dt$, $\mathrm{Re}(\alpha(k)) > 0$, $\alpha(k)$ and $\beta(k)$ are randomly given at every $k$.

The randomized Beta distribution defined in Definition D.1 is crucial for describing a certain type of unknown unknown. This is due to two critical reasons. First, the characteristics of unknown unknowns involve minimal historical data and unpredictable time and distributions. This leads to unavailable models for scientific discoveries and understanding. In the example shown in Figure 9, the parameters $\alpha$ and $\beta$ directly influence the probability density function (pdf) of the distribution, and consequently, the mean and variance. Suppose $\alpha$ and $\beta$ are randomized (expressed as $\alpha(k)$ and $\beta(k)$). In that case, the distribution of $\mathbf{d}(k)$ can take the form of a uniform distribution, exponential distribution, truncated Gaussian distribution, or a combination of these. However, the specific distribution is unknown. Therefore, the randomized $\alpha(k)$ and $\beta(k)$, which result in a randomized Beta distribution, can effectively capture the characteristics of "unavailable model" and "unforeseen" traits associated with unknown unknowns in both time and distribution. Furthermore, the randomized Beta distributions are bounded, with the bounds denoted as $a$ and $c$. This is motivated by the fact that, in general, there are no probabilistic solutions for handling unbounded unknowns, such as earthquakes and volcanic eruptions.
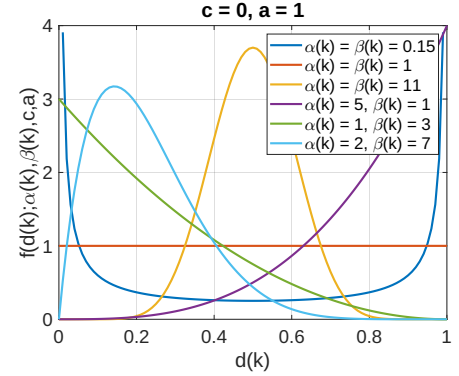


Fig. 9. $\alpha(k)$ and $\beta(k)$ control the robability density the function of the distribution.

## E Experiment: Cart-Pole System

### E.1 Configurations of Pre-training and Runtime Learning

The pre-training configurations for HP-Student, other DRL agents, and the runtime learning mimicking real plants are all the same to ensure fair comparisons. Specifically, we utilize the DDPG algorithm [33] to pre-train DRL and Phy-DRL models and to support runtime learning. The actor and critic networks are implemented as Multi-Layer Perceptrons

(MLPs) with four fully connected layers. The output dimensions of the critic and actor networks are 256, 128, 64, and 1, respectively. The activation functions of the first three neural layers are ReLU, while the output of the last layer is the Tanh function for the actor network and Linear for the critic network. The input of the critic network is $[\mathbf{s}; \mathbf{a}]$, while the input of the actor network is $\mathbf{s}$. In more detail, we set the discount factor $\gamma = 0.9$ and the learning rates of the critic and actor networks to be the same at 0.0003. We set the batch size to 200. The episode consists of 1000 steps, and the sampling frequency is 30 Hz.

### E.2 Design of HP-Student and HA-Teacher

HP-Student is built upon the Phy-DRL agent, allowing us to directly apply its design in the cart-pole system experiment detailed in Appendix K of [12].

Compared to HP-Student, HA-Teacher possesses a more comprehensive understanding of system dynamics in physics, directly and equivalently derived from the dynamics model in [20]:

$$
\frac{d}{dt}\underbrace{\begin{bmatrix} x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{bmatrix}}_{\mathbf{s}} = \underbrace{\begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & \frac{-m_p g \sin\theta \cos\theta}{\theta[\frac{4}{3}(m_c+m_p)-m_p\cos^2\theta]} & \frac{\frac{4}{3}m_p l \sin\theta\dot{\theta}}{\frac{4}{3}(m_c+m_p)-m_p\cos^2\theta} \\ 0 & 0 & 0 & 1 \\ 0 & 0 & \frac{g\sin\theta(m_c+m_p)}{l\theta[\frac{4}{3}(m_c+m_p)-m_p\cos^2\theta]} & \frac{-m_p \sin\theta \cos\theta\dot{\theta}}{\frac{4}{3}(m_c+m_p)-m_p\cos^2\theta} \end{bmatrix}}_{\widehat{\mathbf{A}}(\mathbf{s})} \cdot \begin{bmatrix} x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{bmatrix} + \underbrace{\begin{bmatrix} 0 \\ \frac{\frac{4}{3}}{\frac{4}{3}(m_c+m_p)-m_p\cos^2\theta} \\ 0 \\ \frac{-\cos\theta}{l[\frac{4}{3}(m_c+m_p)-m_p\cos^2\theta]} \end{bmatrix}}_{\widehat{\mathbf{B}}(\mathbf{s})} \cdot \underbrace{F}_{\mathbf{a}},
$$

(57)

where $\widehat{\mathbf{A}}(\mathbf{s})$ and $\widehat{\mathbf{B}}(\mathbf{s})$ are known to the HA-Teacher. The sampling technique transforms the continuous-time dynamics model (59) to the discrete-time one:

$$
\mathbf{s}(k+1) = (\mathbf{I}_4 + T \cdot \widehat{\mathbf{A}}(\mathbf{s})) \cdot \mathbf{s}(k) + T \cdot \widehat{\mathbf{B}}(\mathbf{s}) \cdot \mathbf{a}(k),
$$

from which we obtain the model knowledge $\mathbf{A}(\widehat{\mathbf{s}}_{\sigma(k)})$ and $\mathbf{B}(\widehat{\mathbf{s}}_{\sigma(k)})$ in Equation (15) as

$$
\mathbf{A}(\widehat{\mathbf{s}}_{\sigma(k)}) = \mathbf{I}_4 + T \cdot \widehat{\mathbf{A}}(\widehat{\mathbf{s}}_{\sigma(k)}), \ \mathbf{B}(\bar{\mathbf{s}}^*) = T \cdot \widehat{\mathbf{B}}(\widehat{\mathbf{s}}_{\sigma(k)}),
$$

(58)

where $T = \frac{1}{30}$ second, i.e., the sampling frequency is 30 Hz.

We currently have $\mathbf{A}(\widehat{\mathbf{s}}_{\sigma(k)})$ and $\mathbf{B}(\widehat{\mathbf{s}}_{\sigma(k)})$ in Equation (60). To satisfy Assumption 6.2, we set $\kappa$ to be 0.01. For the inequalities in Equations (17) to (20), we assign the values $\alpha = 0.99$, $\chi = 0.3$, $\gamma_1 = 1$, and $\gamma_2 = 0.1$. Finally, based on the given safety set $\mathbb{X} = \{\mathbf{s} \in \mathbb{R}^n \mid |x| \leq 1, |\theta| < 0.8\}$ and the action space of HA-Teacher $\mathbb{A} = \{\mathbf{a}_{HA} \in \mathbb{R} \mid |\mathbf{a}_{HA}| \leq 40\}$, we obtain the following information for the inequalities in Equations (23) and (26): $\mathbf{D} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1/0.8 & 0 \end{bmatrix}$ and $\mathbf{C} = 1/40$. Finally, according to Theorem 6.3, the real-time patch of HA-Teacher can be directly obtained.

## F Experiment: Real Quadruped Robot

In the real quadruped robot experiment, we utilized a Python-based framework designed for the Unitree A1 robot, which was released on GitHub by [50]. This framework consists of a Pybullet-based simulation, an interface for direct simulation-to-real transfer, and an implementation of the Convex Model Predictive Controller for fundamental motion control.

## F.1 Policy Learning

The runtime learning machine and Phy-DRL are designed to achieve the safe mission described in Section 7.2. The policy observation consists of a 10-dimensional tracking error vector between the robot's state vector and the mission vector. Both systems are based on the DDPG algorithm [33]. The actor and critic networks are implemented as Multi-Layer Perceptrons (MLPs) with four fully connected layers. The output dimensions of the critic network are 256, 128, 64, and 1, while the output dimensions of the actor network are 256, 128, 64, and 6. The input for the critic-network consists of the tracking error vector and the action vector, while the input for the actor network is the tracking error vector. The activation functions for the first three neural layers are ReLU, and the output of the last layer is the Tanh function for the actor network and Linear for the critic network. Additionally, the discount factor $\gamma$ is set to 0.9, and the learning rates for the critic and actor networks are both 0.0003. Finally, the batch size is set to 512.

## F.2 Design of HP-Student and HA-Teacher

We directly apply the design in the Appendix L of [12] to HP-Student here.

Compared to HP-Student, HA-Teacher possesses a deeper understanding of system dynamics, which is directly and equivalently derived from the dynamics model in [18] as

$$
\frac{d}{dt}\underbrace{\begin{bmatrix} h \\ \widetilde{e} \\ v \\ w \end{bmatrix}}_{s} = \underbrace{\begin{bmatrix} \mathbf{O}_{1\times1} & \mathbf{O}_{1\times3} & 1 & \mathbf{O}_{1\times5} \\ \mathbf{O}_{3\times3} & \mathbf{O}_{3\times3} & \mathbf{O}_{3\times3} & \mathbf{R}(\phi,\theta,\psi) \\ \mathbf{O}_{3\times3} & \mathbf{O}_{3\times3} & \mathbf{O}_{3\times3} & \mathbf{O}_{3\times3} \\ \mathbf{O}_{3\times3} & \mathbf{O}_{3\times3} & \mathbf{O}_{3\times3} & \mathbf{O}_{3\times3} \end{bmatrix}}_{\widehat{A}(s)} \cdot \begin{bmatrix} h \\ \widetilde{e} \\ v \\ w \end{bmatrix} + \underbrace{\begin{bmatrix} \mathbf{O}_3 & \mathbf{O}_3 & \mathbf{O}_3 & \mathbf{O}_3 \\ \mathbf{O}_3 & \mathbf{O}_3 & \mathbf{O}_3 & \mathbf{O}_3 \\ \mathbf{O}_3 & \mathbf{O}_3 & \mathbf{I}_3 & \mathbf{O}_3 \\ \mathbf{O}_3 & \mathbf{O}_3 & \mathbf{O}_3 & \mathbf{I}_3 \end{bmatrix}}_{\widehat{B}(s)} \cdot a + \begin{bmatrix} 0 \\ \mathbf{O}_{3\times1} \\ \mathbf{O}_{3\times1} \\ \widetilde{g} \end{bmatrix}, \tag{59}
$$

where system states include the position of the body's CoM height (h), the CoM velocity (v) represented as a 3D vector [CoM x-velocity, CoM y-velocity, CoM z-velocity], the Euler angles ($\widetilde{e}$) described as a 3D vector [roll, pitch, yaw], and the angular velocity in world coordinates (w). and the $\widetilde{g} = [0; 0; -g] \in \mathbb{R}^3$, with $g$ being the gravitational acceleration. We note that the $\mathbf{R}(\phi,\theta,\psi) = \mathbf{R}_z(\psi) \cdot \mathbf{R}_y(\theta) \cdot \mathbf{R}_x(\phi) \in \mathbb{R}^{3\times3}$ in Equation (59) is the rotation matrix, with

$$
\mathbf{R}_x(\phi) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos\phi & -\sin\phi \\ 0 & \sin\phi & \cos\phi \end{bmatrix}, \; \mathbf{R}_y(\theta) = \begin{bmatrix} \cos\theta & 0 & \sin\theta \\ 0 & 1 & 0 \\ -\sin\theta & 0 & \cos\theta \end{bmatrix}, \; \mathbf{R}_z(\psi) = \begin{bmatrix} \cos\psi & -\sin\psi & 0 \\ \sin\psi & \cos\psi & 0 \\ 0 & 0 & 1 \end{bmatrix}.
$$

The $\widehat{A}(s)$ and $\widehat{B}(s)$ are known to the HA-Teacher for his real-time patch design. The sampling technique transforms the continuous-time dynamics model in Equation (59) to the discrete-time one:

$$
s(k+1) = (\mathbf{I}_4 + T \cdot \widehat{A}(s)) \cdot s(k) + T \cdot \widehat{B}(s) \cdot a(k) + T \cdot g(s),
$$

from which we obtain the knowledge of $A(\widehat{s}_{\sigma(k)})$ and $B(\widehat{s}_{\sigma(k)})$ in Equation (15) as

$$
A(\widehat{s}_{\sigma(k)}) = \mathbf{I}_4 + T \cdot \widehat{A}(\widehat{s}_{\sigma(k)}) \text{ and } B(\widehat{s}_{\sigma(k)}) = T \cdot \widehat{B}(\widehat{s}_{\sigma(k)}). \tag{60}
$$

Meanwhile, for the patch in Equation (11), the model mismatch in Assumption 6.2, and the dwell time in Equation (8), we let $\chi = 0.25$, $\kappa = 0.01$, and $\tau = 100$. For LMIs in Equations (17) to (20), we let $\alpha = 0.9$, $\gamma_1 = 1$, and $\gamma_2 = 0.45$. Finally, according to the given safety set $\mathbb{X} = \{ s \mid |\text{CoM x-velocity} - 0.3 \text{ m/s}| \leq 0.3 \text{ m/s}, |\text{CoM z-height} - 0.24 \text{ m}| \leq 0.15 \text{ m} \}$ and the action space of HA-Teacher $\mathbb{A} = \{ a_{HA} \mid |a_{HA}| \leq [30, 30, 30, 60, 60, 60]^\top \}$, we obtain following knowledge for

the LMIs in Equations (23) and (26):

$$
\mathbf{D} = \begin{bmatrix} 1/0.15 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1/0.3 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}, \mathbf{C} = \begin{bmatrix} 1/30 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/30 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1/30 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/60 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1/60 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1/60 \end{bmatrix}.
$$

Finally, according to Theorem 6.3, the real-time patch of HA-Teacher can be directly obtained.

## G Computation Resources

In all case studies, we trained and tested the deep reinforcement learning (DRL) algorithm on a desktop computer running Ubuntu 22.04. The desktop was equipped with a 12th Gen Intel(R) Core(TM) i9-12900K 16-core processor, 64 GB of RAM, and an NVIDIA GeForce GTX 3090 GPU. The DRL algorithm was implemented in Python using the TensorFlow framework. We utilized the open-source Python CVX solver to solve LMI (Linear Matrix Inequalities) problems.

In our system architecture, the computation of $\widehat{\mathbf{F}}_{\sigma(k)}$ and $\widehat{\mathbf{P}}_{\sigma(k)}$ for the HA-Teacher at each patch needs to be performed when the Safety Coordinator is triggered. To ensure real-time computation of CVX and interaction with the environment, we have implemented a multi-processing pipeline to control the robot and solve LMIs in parallel in real-time. For solving LMIs, we always allow the solver to use the most recent state so that when the safety coordinator is triggered, the latest $\widehat{\mathbf{F}}$ and $\widehat{\mathbf{P}}$ are readily available. We have taken into account the delay issue and formulated it in the LMI problems.

We observed that the MATLAB-based CVX solver consistently solved the LMIs problem better than the Python-based solver, providing more reliable solutions. However, transferring data between MATLAB and Python could cause additional delays when updating $\widehat{\mathbf{F}}_{\sigma(k)}$ and $\widehat{\mathbf{P}}_{\sigma(k)}$ for HA-Teacher. Additionally, implementing multiprocessing in both MATLAB and Python posed technical challenges due to software compatibility issues. As a result, we opted for the Python-based CVX solver for real-time real-world experiments, while recommending the MATLAB-based solver for less time-sensitive applications.