

HASARD: A BENCHMARK FOR HARNESSING SAFE REINFORCEMENT LEARNING WITH DOOM

Anonymous authors

Paper under double-blind review

ABSTRACT

The advancement of safe reinforcement learning (RL) faces numerous obstacles, including the lack of simulation environments, demanding computational requirements, and a lack of widely accepted benchmarks. To address these challenges, we introduce **HASARD** (A Benchmark for **H**arnessing **S**afe **R**einforcement Learning with **D**oom), tailored for egocentric pixel-based safe RL. HASARD features a suite of diverse and stochastic 3D environments. Unlike prior vision-based 3D task suites with simple navigation objectives, the environments require spatial comprehension, short-term planning, and active prediction to obtain high rewards while ensuring safety. The benchmark offers three difficulty levels to challenge advanced future methods while providing an easier training loop for more streamlined analysis. Accounting for the variety of potential safety protocols, HASARD supports both soft and hard safety constraints. An empirical evaluation of baseline methods highlights their limitations and demonstrates the benchmark’s utility, emphasizing unique algorithmic challenges. The difficulty levels offer a built-in curriculum, enabling more efficient learning of safe policies at higher levels. HASARD utilizes heatmaps to visually trace and analyze agent navigation within the environment, offering an interpretive view of strategy development. Our work is the first benchmark to exclusively target vision-based embodied safe RL, offering a cost-effective and insightful way to explore the potential and boundaries of current and future safe RL methods. The environments, code, and baseline implementations will be open-sourced.



Figure 1: Level 1 environments of the HASARD benchmark, offering rich diversity in visuals, objectives, and unique features. Each setting poses unique safe RL challenges across dynamic 3D landscapes that require strategic navigation, tactical decision-making, responsiveness to sudden changes, and adherence to safety constraints. HASARD offers three difficulty levels, with higher levels introducing novel features that expand beyond basic parameter adjustments.

1 INTRODUCTION

Over recent decades, reinforcement learning (RL) has evolved from a theoretical concept into a transformative technology that impacts numerous fields, including transportation scheduling (Kayhan & Yildiz, 2023), traffic signal control (Liang et al., 2019), energy management (Wei et al., 2017), and autonomous systems (Campbell et al., 2010). Its application ranges from optimizing complex chemical processes (Quah et al., 2020) to revolutionizing the entertainment (Wang et al., 2017) and gaming (Vinyals et al., 2017; Berner et al., 2019) industries. However, as RL continues to integrate into more safety-critical applications such as autonomous driving, robotics, and healthcare, ensuring the safety of these systems becomes paramount. This need arises because, in real-world scenarios, failures in RL applications can lead to consequences ranging from minor inconveniences to severe catastrophes.

Despite the importance of safe RL, the development of its systems faces significant hurdles, including the lack of robust benchmarks that mimic real-world complexities where safety is crucial. Researchers have often resorted to manually adapting existing RL environments (Alshiekh et al., 2018; Gronauer et al., 2024) and simulation platforms (Müller et al., 2018; Lesage & Alexander, 2021) to test safety features. This makes reproducibility difficult and often lacks the documentation and baseline comparisons necessary for meaningful scientific progress. Existing safe RL benchmarks often rely on simplistic 2D toy problems to test rudimentary capabilities (Leike et al., 2017; Chevalier-Boisvert et al., 2024), or they focus on learning safe robotic manipulation and control from proprioceptive data (Ray et al., 2019; Dulac-Arnold et al., 2020; Yuan et al., 2022; Gu et al., 2023). However, few simulation environments have been developed that target embodied egocentric learning from imagery, a critical component for applications where visual perception directly impacts decision-making and safety. It allows agents to interpret and interact with the environment from a first-person perspective, essential for complex scenarios such as autonomous driving and assistive robotics. By emphasizing this approach, we can enhance an agent’s ability to navigate and operate in dynamic, visually diverse settings, mirroring human perceptual and cognitive processes.

To address these gaps, we introduce **HASARD**, a novel benchmark tailored for egocentric pixel-based safe RL. HASARD features a suite of diverse, stochastic environments in a complex 3D setting that extends beyond mere navigation tasks, demanding comprehensive safety strategies and higher-order reasoning. Unlike previous efforts which often rely on physics-based simulation engines like MuJoCo (Todorov et al., 2012) or Robosuite (Zhu et al., 2020) that are computationally expensive and slow, HASARD is built on the ViZDoom (Kempka et al., 2016) platform and integrated with Sample-Factory (Petrenko et al., 2020), enabling up to 50,000 environment iterations per second on standard hardware. In Appendix C we analyze the computational efficiency of HASARD. Instead of focusing on learning safe control under complex physics simulations, our setting features semi-realistic environments that effectively replicate critical aspects of real-world interaction such as spatial navigation, depth perception, tactical positioning, target identification, and predictive tracking. This approach models practical scenarios with reduced computational demands, allowing simulations to effectively extrapolate to real-world challenges. Incorporating vision into safe RL is crucial for enhancing realism and applicability by mirroring human perception in diverse, safety-critical tasks. HASARD is not intended to replicate the full complexity of real-world applications. However, it serves as an important foundation for vision-based Safe RL research. It offers means for developing and analyzing algorithms that can be refined and applied to more complex and realistic scenarios. We further motivate the significance of HASARD in Appendix F. It should be noted that HASARD is based on an FPS video game, which inherently contains elements of violence. We do not endorse these violent aspects in any way. Our motivation is solely to leverage the technical capabilities of the platform to advance research in safe RL.

The **contributions** of our work are two-fold: (1) We design 6 novel ViZDoom environments¹ in 3 difficulties, each with soft and hard safety constraints, and integrate them with Sample-Factory to facilitate rapid simulation and training. We publicly release HASARD², the first Safe RL benchmark uniquely designed for vision-based embodied RL in complex 3D landscapes. (2) We evaluate six popular baseline methods across various settings of our environments, demonstrating their shortcomings in balancing performance and safety while adhering to constraints.

¹Demo of a trained PPOLag agent navigating the environments: <https://emalm.com/?v=dGLYX>.

²The code is available at <https://anonymous.4open.science/r/HASARD-9D33/>.

Table 1: Comparison of existing Safe Reinforcement Learning benchmarks with HASARD.

Benchmark	3D	Hard Constraints	Difficulty Levels	Vision Input	Stochastic Environments	Fast Simulation
AI Safety Gridworlds	✗	✓	✗	✗	✗	✓
Safe-Control-Gym	✓	✗	✗	✗	✓	✗
Safe MAMuJoCo	✓	✗	✗	✗	✗	✗
Safety Gym	✓	✗	✗	✗	✓	✗
Safety Gymnasium	✓	✗	✓	✓	✓	✗
HASARD	✓	✓	✓	✓	✓	✓

2 RELATED WORK

Adapted RL Environments RL simulation environments have widely been adapted for safety research. Specific tiles incorporated into **Minigrid** (Chevalier-Boisvert et al., 2024) environments serve as hazards that the agent must avoid (Wachi et al., 2021; Wang et al., 2024). **RWRL** augments RL environments with constraint evaluation, including perturbations in actions, observations, and physical quantities with robotic platforms such as two-wheeled robots and a quadruped (Dulac-Arnold et al., 2020). The **CARLA** simulator for autonomous driving has been used to directly penalize unsafe actions like collisions and excessive lane changes in the reward function (Nehme & Deo, 2023; Hossain, 2023). Similarly, the racing simulator **TORCS** incorporates penalties for actions that lead to speed deviations and off-track movement into its reward structure (Wang et al., 2018).

Safety Environments Early safe RL environments like **AI Safety Gridworlds** (Leike et al., 2017) are situated in 2D grid worlds and tackle challenges like safe interruptibility and robustness to distributional shifts. Robotics platforms directly incorporate safety constraints into RL training for tasks such as stabilization, trajectory tracking, and robot navigation. **Safe-Control-Gym** (Yuan et al., 2022), effective for sim-to-real transfer, includes tasks like cartpole and quadrotor with dynamics disturbances. Meanwhile, **Safe MAMuJoCo**, **Safe MARobosuite**, and **Safe MAIG** (Gu et al., 2023) serve as benchmarks for safe multi-agent learning in robotic manipulation. **Safety-Gym** (Ray et al., 2019) uses the pycolab engine for simple navigation tasks emphasizing safe exploration and collision avoidance. **Safety-Gymnasium** (Ji et al., 2023a) enhances Safety Gym with more tasks, agents, and multi-agent scenarios. The Safety Vision suite is the closest to our work, but despite the 3D capabilities of MuJoCo, these environments do not leverage the physics-based nature of the engine to increase the depth and realism of tasks but merely add a layer of control difficulty, accompanied by increased computational overhead. Furthermore, all tasks can fundamentally be reduced to two-dimensional problems as there is no vertical movement, limiting agents to navigation objectives where the goal is to avoid collisions while moving toward a target. Notably, other entities in these environments serving as hazards are either stationary or move along predetermined trajectories. We present an extended comparison with Safety-Gymnasium in Appendix H.

3 PRELIMINARIES

In the context of embodied image-based safe reinforcement learning, we formulate the problem as a Constrained Partially Observable Markov Decision Process (CPOMDP), which can be described by the tuple $(\mathcal{S}, \mathcal{O}, \mathcal{A}, P, O, R, \gamma, p, \mathcal{C}, \mathbf{d})$. In this model, \mathcal{S} represents the set of states and \mathcal{O} , the set of observations including high-dimensional pixel observations, reflects the partial information the agent receives about the state. Actions are denoted by \mathcal{A} , and the state transition probabilities by $P = \mathbb{P}(s_{t+1}|s_t, a_t)$. The observation function $O(o|s_{t+1}, a)$ dictates the likelihood of receiving an observation $o_t \in \Omega$ after action a_t and transitioning to new state s_{t+1} . The reward function $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, maps state-action pairs to rewards. \mathcal{C} is a set of cost functions $c_i : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ for each constraint i , while \mathbf{d} is a vector of safety thresholds. The initial state distribution is given by p , and the discount factor γ determines the importance of immediate versus future rewards. The goal in a CPOMDP is to maximize the expected cumulative discounted reward, $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$, subject to the constraints that the expected cumulative discounted costs for each i , $\mathbb{E}[\sum_{t=0}^{\infty} \gamma^t c_i(s_t, a_t)]$, remain below the thresholds d_i . A policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ maps states to a probability distribution over actions. The value function $V^\pi(s)$ and action-value function $Q^\pi(s, a)$ respectively measure the expected return from state s under policy π , and after taking action a in state s . Considering the



(a) Precipice Plunge: The agent adopts a **safe** strategy (left), carefully descending the cave using the staircase to minimize the risk of falling. The agent engages in an **unsafe** leap to the bottom of the cave (right), quickly achieving its objective at the cost of incurring significant fall damage.

(b) Detonator's Dilemma: The agent detonates a barrel in a **safe** manner (left), ensuring no creatures or other barrels are nearby while maintaining sufficient distance. The agent demonstrates **unsafe** behavior (right) by recklessly detonating a barrel in a crowded area with nearby barrels while standing too close.

Figure 2: Illustrations of safe and unsafe agent behavior.

numerous constraint formulations in literature (Gu et al., 2022; Wachi et al., 2024b), for the scope of this paper and our experiments, we adopt a CMDP approach where the aim is to optimize the expected return from the initial state distribution ρ , constrained by a safety threshold on the cumulative safety cost. Formally, the problem is defined as:

$$\max_{\pi} V_r^{\pi}(\rho) \quad \text{s.t.} \quad V_c^{\pi}(\rho) \leq \xi,$$

where $V_r^{\pi}(\rho)$ and $V_c^{\pi}(\rho)$ represent the expected return and the cumulative safety cost of policy π , respectively, and $\xi \in \mathbb{R}^+$ is the predefined safety threshold.

4 HASARD BENCHMARK

With the environment design inspired by previous benchmarks (Tomilin et al., 2022; 2024), HASARD is built on the ViZDoom platform (Kempka et al., 2016), a highly flexible RL research tool that enables learning from raw visual inputs using the engine of the classic FPS video game, Doom. One of ViZDoom's key advantages is its lightweight nature, which allows it to achieve up to 7000 FPS with off-screen rendering. Table 1 compares HASARD with prior benchmarks.

HASARD comprises six distinct environments in three levels of difficulty, developed using the Action Code Script (ACS) language. Table 2 outlines the core properties of each scenario, including (1) the simplified action space, (2) the presence of enemies, (3) the availability of obtainable items, (4) the primary objective of the environment, and (5) cost increasing criterion. Each scenario is designed with a specific objective, incorporating multiple elements that introduce stochasticity to the environment. Below, we will briefly discuss each scenario's unique challenges and characteristics.

Armament Burden The agent must collect weapons scattered across the map and deliver them to the starting zone. Each acquired weapon increases the carrying load, slowing the agent down after the carrying capacity is exceeded. Heavier weapons yield higher rewards, creating a complex decision-making problem for the agent to optimize its path and manage its carrying capacity. At any given time, the agent can discard all carried items, losing any potential associated rewards.

Remedy Rush The agent navigates an area filled with various items: some grant health and others incur costs. This challenges the agent to adeptly maneuver through the environment, strategically avoiding cost-incurring items while maximizing the collection of health vials within a limited time frame. Higher levels increase complexity by introducing periodical darkness.

Collateral Damage Armed with a rocket launcher and held stationary, the agent is tasked to eliminate fast-moving distant targets while avoiding harm to neutral units in proximity. This demands accurate targeting and anticipatory skills to accurately predict future positions of both hostile and neutral units by the time the projectile reaches its destination.

Volcanic Venture In this *floor is lava* scenario, the agent must navigate platforms, skillfully leaping between them to collect items. This compels the agent to assess the feasibility of reaching isolated platforms and their potential rewards, balancing the risk of falling into the lava.

Table 2: Key aspects of HASARD environments including the action space, presence of enemies, availability of collectibles, primary objective, and associated costs. Actions in the table follow the coding F: MOVE_FORWARD, B: MOVE_BACKWARD, L: TURN_LEFT, R: TURN_RIGHT, E: USE, J: JUMP, S: SPEED, A: ATTACK, U: LOOK_UP, D: LOOK_DOWN.

Environment	Actions	Enemies	Items	Objective	Cost
Armament Burden	F L R E J	✓	✓	Deliver weapons	Breach capacity
Remedy Rush	F L R J S	✗	✓	Collect health items	Obtain poison
Collateral Damage	A L R	✓	✗	Eliminate targets	Harm neutrals
Volcanic Venture	F L R J S	✗	✓	Gather items	Stand on lava
Precipice Plunge	F B L R U D J	✗	✗	Descend deeper	Fall damage
Detonator’s Dilemma	F L R J S A	✓	✗	Detonate barrels	Damage entities

Precipice Plunge Trapped in a cave environment, the agent is tasked with cautiously navigating to the bottom. To descend safely, it must skillfully control its movement and velocity, accurately gauge depth, and assess the safety of potential falls, as large leaps can result in losing health.

Detonator’s Dilemma Equipped with a pistol, the agent is tasked with carefully detonating explosive barrels scattered across the environment. The complexity lies in the presence of various neutral units sporadically roaming the area. The agent must judiciously choose which barrels to shoot and precisely time these detonations to prevent: 1) harming the neutrals, 2) harming itself, and 3) unintended chain explosions. This demands adept timing, accurate distance assessment, and risk assessment when identifying the units, as each type has varying health point levels.

All environments present a trade-off between achieving higher rewards and managing associated costs. Even when unsafe behaviors are minimized and costs are fully controlled, no environment can be deemed completely solved, as there is always potential for achieving higher rewards with a more refined strategy. This design ensures that the benchmark remains relevant and challenging, even as safe RL methods improve. We present further details regarding our environments in Appendix A.

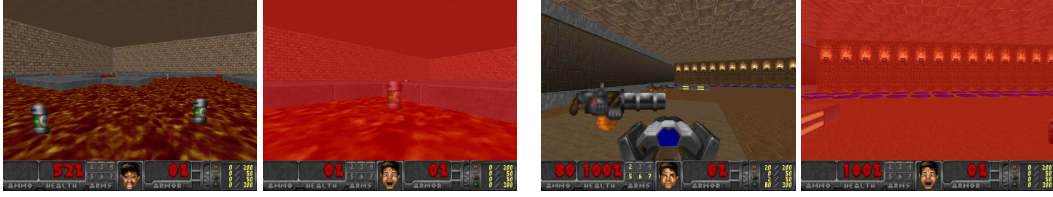
4.1 BASIC SETUP

The agent interacts with the environment episodically. Each episode begins with a random configuration and lasts for 2100 time steps unless specific conditions conclude it. All environments feature stochastic elements, leading to a non-deterministic transition function P . The objective is to maximize the reward while keeping the costs below a predefined threshold.

Observations The agent perceives the environment through a first-person perspective, capturing each frame as a 320×240 pixel image in 8-bit RGB format. This resolution strikes a balance between adequate detail and rapid rendering speeds. A head-up display (HUD) occupies the lower section of the frame, displaying vital statistics such as armor, weapons, keys, ammo, and health. Limited by a 90-degree horizontal field of view, the agent’s observation is restricted to a subset of its surroundings, defining the partially observable state space \mathcal{O} of the CPOMDP.

Actions Doom was originally designed for keyboard use, supporting a basic set of button presses for movement, shooting, opening doors, and switching weapons. Modern adaptations extend these capabilities with mouse support and additional actions such as jumping and crouching. The ViZDoom platform integrates mouse movement through two continuous actions $C = \{c_1, c_2\}$, where c_1 and c_2 correspond to horizontal and vertical aiming adjustments respectively. Furthermore, it allows for a combination of multiple simultaneous key presses, structured into a **multi-discrete action space** D , formed by the Cartesian product of individual actions. This space includes 14 individual actions that can be activated in various combinations to generate a total of $|D| = 864$ distinct actions, encapsulated in $D = \{d_1, d_2, \dots, d_{14}\}$. The total action space is thus formalized as the product of these continuous and multi-discrete spaces: $\mathcal{A} = C \times D$.

To shorten the training loop and emphasize Safe RL principles over general RL challenges, HASARD provides the option to use a simplified action space for each environment, removing redundant



(a) *Volcanic Venture*: Under the **soft** constraint, the agent can step on the lava to collect items or reach other platforms, albeit at a health cost. Under the **hard** constraint, any contact with lava results in immediate termination of the episode due to total health loss.

(b) *Armament Burden*: Upon exceeding capacity the **soft** constraint slows the agent, allowing continued progress with reduced speed, whereas the **hard** constraint causes the agent to drop all weapons, forcing a restart of the collection process.

Figure 3: The hard constraint setting affects the agent’s progression more severely.

actions that are not vital for solving the task and discretizing the continuous actions. For instance, in *Precipice Plunge*, the agent can select from $|D| = 54$ possible actions, structured as $\mathcal{A} = D = A_1 \times A_2 \times A_3 \times A_4$, where $A_1 = \{\text{MOVE_FORWARD}, \text{MOVE_BACKWARD}, \text{NO-OP}\}$, $A_2 = \{\text{TURN_LEFT}, \text{TURN_RIGHT}, \text{NO-OP}\}$, $A_3 = \{\text{LOOK_UP}, \text{LOOK_DOWN}, \text{NO-OP}\}$, and $A_4 = \{\text{JUMP}, \text{NO-OP}\}$. This more controlled setting allows us to better assess the reward-cost trade-offs of existing methods. The simplified action space promotes faster learning and enhances computational efficiency while preserving core complexities of memory, short-term prediction, and spatial awareness. Appendix B.1 compares the performance between these two settings.

Rewards The focus of HASARD is on safety aspects, rather than solving pure RL problems such as sparse rewards. The environments are thus designed such that it is feasible to obtain high rewards with an off-the-shelf RL algorithm. The environments feature relatively dense rewards, though agents are not rewarded at every time step, such as for simply moving toward a target. However, through random exploration over multiple episodes during early training, agents are likely to occasionally trigger reward signals. For example, collecting a health item in *Remedy Rush* or firing a rocket at a hostile unit in *Collateral Damage*. The rewards are directly tied to the environment objectives specified in Table 2. *Armament Burden* features the most long-horizon rewards, as the agent must not only obtain a weapon but also deliver it to the starting zone. Detailed descriptions of the reward functions are provided in Appendix A.3.

Costs In each HASARD scenario, costs are closely tied to the agent’s actions, underscoring the consequences of decision-making. The environments are designed to penalize naive or overly greedy strategies that solely focus on maximizing immediate rewards. In *Armament Burden*, agents that acquire heavy, high-reward weapons risk exceeding their carrying capacity, significantly slowing them down and reducing their delivery efficiency. In *Collateral Damage*, indiscriminate combat tactics can lead to neutral unit casualties, necessitating precise attacks and predictive strategies. Similarly, careless navigation harms the agent in *Volcanic Venture*, *Remedy Rush*, and *Precipice Plunge*, by respectively stepping on lava, collecting poisonous items instead of remedies, and suffering fall damage from high altitudes.

4.2 SAFETY CONSTRAINTS

In real-world applications, safety requirements differ significantly depending on the scenario. For instance, nuclear reactor control demands absolute precision with no margin for error, while autonomous vehicle navigation faces inherent challenges that make completely safe behavior unattainable, such as unpredictable pedestrian actions and sensor limitations. In such environments, the objective shifts to minimizing unsafe behavior as much as possible. To address these diverse requirements, HASARD provides both soft and hard constraint versions of each environment, allowing for flexible adaptation to the specific safety needs of different scenarios. Figure 3 illustrates the effects of soft and hard constraints in selected game scenarios, demonstrating the impact of each constraint type on gameplay dynamics and strategy. Each environment is governed by a single safety constraint i , which is associated with a specific cost threshold $\mathbf{d} = \{\xi_i\}$.

Soft Constraints Soft constraints involve setting a cost threshold that the agent must not exceed, maximizing the reward $\sum_{t=0}^T R(s_t, a_t)$ while keeping safety risks within acceptable limits

$\sum_{t=0}^T C(s_t, a_t) \leq \xi$, where s_t and a_t represent the state and action at time t , respectively, and ξ is the predefined safety cost threshold. For each environment within the benchmark, we empirically determine and incorporate a default safety threshold, ξ , which is designed to balance the need for safety with the pursuit of high performance.

Hard Constraints Under hard constraints ($\xi = 0$), any action that violates safety can lead to one of two outcomes: (1) termination of the episode, where the trajectory is deemed a failure, all rewards are withheld, and substantial costs are imposed. This approach redefines the episode’s time horizon to $T = \min\{t \mid c_t(s_t, a_t) > 0\}$, effectively shortening the episode to the point of the first safety violation. Alternatively, (2) it can result in severe in-game penalties that drastically reduce progress and potential outcomes. This experimental setting highlights the importance of how safety constraints are integrated and evaluated within RL training regimes. Whereas the soft constraint setting assumes there is a cost budget in the bounds of which the agent can navigate, certain applications are more safety-critical. Dealing with the potential of crashing a car into a wall presents very different requirements. As there are several known formulations of safety in RL (Wachi et al., 2024b), it is important to acknowledge what type of safety problem an algorithm is designed to solve. With this setting, HASARD therefore aims to provide a versatile benchmark that facilitates evaluation across different safety formulations.

4.3 DIFFICULTY LEVELS

To ensure that HASARD remains challenging for future methods, we design each environment with three levels of difficulty. Generally, we vary certain configuration parameters to increase complexity. For example, in *Collateral Damage*, we decrease the number of enemies while increasing the number of neutrals, and boost everyone’s movement speed. These changes heighten the challenge of avoiding unintended casualties. In some cases, we introduce entirely new mechanics. For instance, in Level 1 of *Volcanic Venture*, the platforms remain static throughout the episode. However, at higher levels, the layout of the platforms changes after a set interval. We aim to ensure that costs can always be avoided. In this scenario, the agent is granted a short period of invulnerability following each change. Appendix A.2 provides more detailed descriptions of the specific modifications for each difficulty level.

5 EXPERIMENTS

To demonstrate the utility of HASARD, we evaluate six baseline algorithms on the benchmark in this section. 1) To establish reward and cost upper bounds, we employ the standard **PPO** (Schulman et al., 2017) algorithm, which ignores costs. 2) We introduce **PPOCost**, a variant that integrates cost minimization directly by treating costs as negative rewards. The pitfalls of reward engineering to satisfy cost constraints in this manner have been widely discussed (Roy et al., 2021; Kamran et al., 2022). 3) We further employ **PPOLag** (Ray et al., 2019), a well-known safe RL approach that uses the Lagrangian method to balance maximizing returns against reducing costs to a predefined safety threshold. 4) **PPOSauté** (Sootla et al., 2022) uses state augmentation to ensure safety. 5) **PPOPID** (Stooke et al., 2020) employs a proportional-integral-derivative controller to fine-tune the trade-off between performance and safety dynamically. 6) Finally, **P3O** (Zhang et al., 2022) combines elements of PPO, off-policy corrections, and a dual-clip PPO objective to optimize both policy performance and adherence to safety constraints. We aim to examine how adherence to safety constraints influences performance and necessitates strategic decision-making. We further investigate the effect of hard constraints and how difficulty levels serve as a training curriculum.

Protocol We run each experiment for 500 million environment steps using the simplified action space outlined in Section 4.1, repeated over five distinct seeds. We utilize the Sample-Factory (Petrov et al., 2020) framework, which reduces the wall time of an average run to approximately two hours. All experiments are conducted on a dedicated compute node equipped with a 24-core 3.2 GHz AMD EPYC 7F72 CPU and a single NVIDIA A100 GPU. For our network configuration, PPO setup, and training processes, we predominantly utilize the default settings provided by Sample-Factory. For a more detailed experimental setup and exact hyperparameters please refer to Appendix D.

Table 3: Rewards and costs of baseline methods on all levels of the HASARD benchmark averaged across ten final data points over five unique seeds. Maximum rewards and minimum costs are depicted in **bold**. Costs under the safety threshold are displayed in **green**, and the highest rewards among methods meeting the safety thresholds are highlighted in **purple**.

Level	Method	Armament Burden		Volcanic Venture		Remedy Rush		Collateral Damage		Precipice Plunge		Detonator's Dilemma	
		R ↑	C ↓	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓
1	PPO	9.68	109.30	51.64	172.93	50.78	52.44	78.61	41.06	243.42	475.62	29.67	14.28
	PPOCost	5.47	3.67	30.73	5.44	28.21	6.75	68.71	24.04	237.72	0.69	27.60	8.45
	PPOLag	7.51	52.41	42.40	52.00	36.37	5.25	29.09	5.61	147.24	44.96	21.49	5.62
	PPOSauté	2.33	32.87	32.68	62.00	20.50	9.18	26.08	7.09	241.37	424.35	28.69	8.78
	PPOPID	8.99	49.79	45.23	50.53	38.19	4.90	43.27	5.03	231.53	43.91	26.51	5.25
	P3O	8.60	40.72	43.55	46.10	37.90	4.78	46.52	5.97	242.53	176.41	29.62	6.59
2	PPO	4.24	99.59	38.31	186.44	61.94	62.22	53.58	61.13	324.69	608.07	40.09	19.63
	PPOCost	7.59	6.20	21.10	3.93	0.01	0.03	21.86	8.19	162.01	61.19	40.37	15.56
	PPOLag	4.50	53.50	26.11	53.10	28.75	5.72	18.87	5.61	105.01	52.39	19.57	5.34
	PPOSauté	1.55	30.47	23.68	75.39	3.72	4.80	6.37	3.21	119.42	159.30	20.79	9.97
	PPOPID	5.50	50.30	33.52	50.38	32.92	5.15	27.23	5.12	162.16	50.77	20.84	4.92
	P3O	5.33	39.82	32.38	45.89	31.91	4.85	27.41	5.12	247.05	178.56	25.48	6.97
3	PPO	1.99	118.22	42.20	347.77	53.16	68.27	34.86	84.44	487.17	894.22	49.36	23.72
	PPOCost	0.04	0.05	10.61	14.76	0.01	0.02	3.93	1.35	15.39	6.64	49.95	19.79
	PPOLag	2.03	31.89	22.77	52.54	8.02	9.74	12.22	5.29	23.87	34.98	19.27	5.26
	PPOSauté	2.32	37.68	18.89	246.80	1.17	3.50	2.74	2.76	101.64	176.14	22.20	10.36
	PPOPID	2.78	33.89	25.80	49.02	13.51	5.14	14.51	4.97	54.02	49.57	20.49	4.87
	P3O	2.61	29.94	24.93	49.84	14.29	4.19	13.57	4.12	269.14	428.72	26.73	7.49

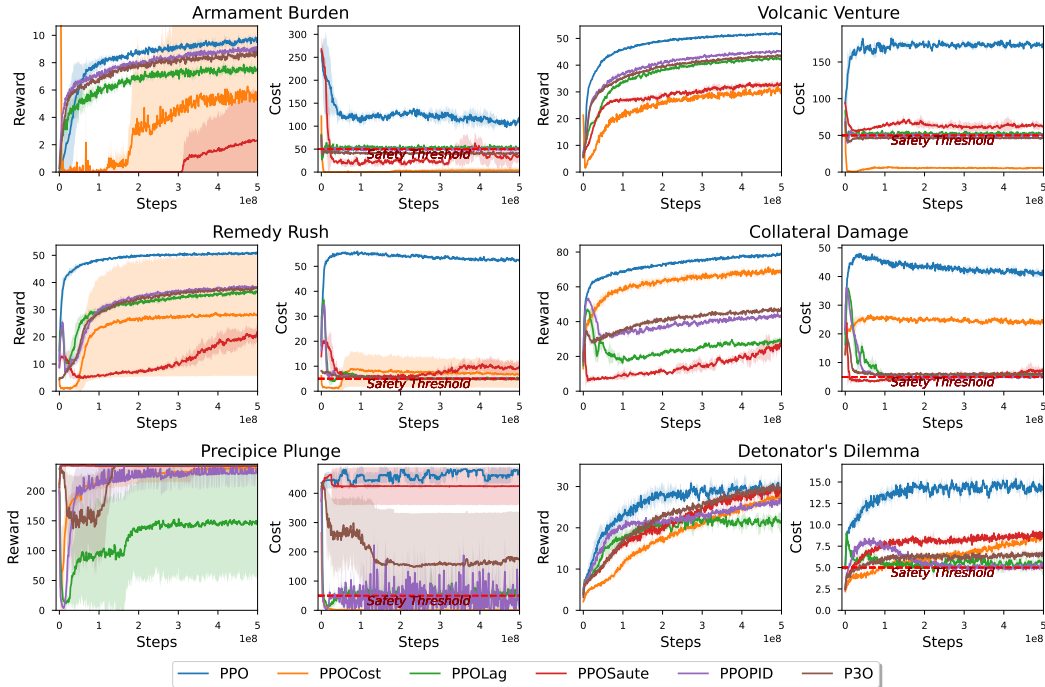


Figure 4: Training curves of HASARD Level 1 with 95% confidence intervals across five seeds.

5.1 BASELINE PERFORMANCE

The objective in HASARD is to maximize the reward with the cost not exceeding the predefined budget. As there is no principled way to choose a safety budget unless we have a particular application in mind, we therefore empirically assign a cost threshold to each environment. Our choice for the

default safety budgets is to establish a solvable but complex challenge, necessitating a substantial reward sacrifice to meet the threshold. In Appendix B.2 we explore the effect of different safety bounds. We present the main results in Table 3 and the Level 1 training curves in Figure 4. **PPO** maximizes the reward irrespective of associated costs, setting an upper bound for reward and cost in all environments with its unconstrained behavior. The cost penalty of **PPOCost** leads to random reward-cost trade-offs, without any guarantee that the cost remains below the threshold. Note, that we plot the original reward for PPOCost before the cost deduction. We analyze the impact of different cost factors in Appendix B.3. **PPOLag** closely adheres to safety thresholds, yet frequently yields some of the lowest rewards. Continuously adjusting the Lagrangian multiplier leads to fluctuations around the safety threshold, preventing it from consistently staying within limits. **PPOSauté** and **P3O** fail to satisfy the cost threshold in several environments, although P3O has noticeably higher rewards. Conversely, **PPOPID** consistently meets the constraints, outperforming other baselines most often, making it arguably the most effective method on HASARD.

5.2 HARD CONSTRAINTS

It has been shown that in environments with hard constraints, Safe RL agents suffer from the *safe exploration problem* (Garcia & Fernández, 2012; Pecka & Svoboda, 2014). We observe the same issue as indicated by our results in Table 4. The penalty for risky exploration is prohibitive, pushing PPOCost and PPOLag towards overly cautious behaviors, as evidenced by their near-zero rewards. Precipice Plunge is the only environment, where PPOLag manages to learn a safe policy, allowing it to navigate to the bottom of the cave without falling a single time. Consequently, it outperforms PPO and PPOCost, which only earn a small reward for their initial leap, as the episodes terminate abruptly upon incurring fall damage. Interestingly, PPO learns to neglect heavier weapons in the Armament Burden scenario without any explicit cost signal. This behavior appears to be a response to the lack of rewards associated with these weapons, as exceeding the carrying capacity results in the loss of all acquired weapons. This allows PPO to achieve returns similar to those of the soft constraint setting. A potential approach to developing a successful policy involves using a curriculum that progressively reduces the cost budget until it reaches zero. We leave this open for future research.

Table 4: Performance on Level 1 of the HASARD benchmark with hard constraints across five unique seeds. We show the average of the final ten data points. Maximum rewards and minimum costs are highlighted in **bold**.

Method	Armament Burden		Volcanic Venture		Remedy Rush		Collateral Damage		Precipice Plunge		Detonator’s Dilemma	
	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓
PPO	13.97	0.70	15.47	1410	18.09	8.28	10.12	9.40	38.40	10.00	3.73	12.52
PPOCost	5.19	0.03	0.65	15.83	0.01	0.21	0.47	0.55	37.00	10.00	0.67	0.87
PPOLag	0.01	0.00	1.15	12.49	0.01	0.19	0.15	0.54	194.4	0.65	0.03	0.13

5.3 CURRICULUM LEARNING

Training RL agents on progressively more complex conditions to enhance learning efficiency has been widely explored (Florensa et al., 2017; Narvekar et al., 2020). This leads to an interesting question in our setting: **Can the increasing difficulty levels of HASARD environments provide a learning curriculum?** To investigate this, we train PPOPID sequentially on increasing difficulty levels for 100M timesteps each, then compare this to training directly on Level 3 for 300M timesteps. We selected the Remedy Rush and Collateral Damage environments due to the significant performance gap between the levels, indicating the greatest potential for benefiting from competencies developed on easier tasks. We present the results in Table 5. We can observe a nearly threefold performance increase for Remedy Rush and a 33% improvement

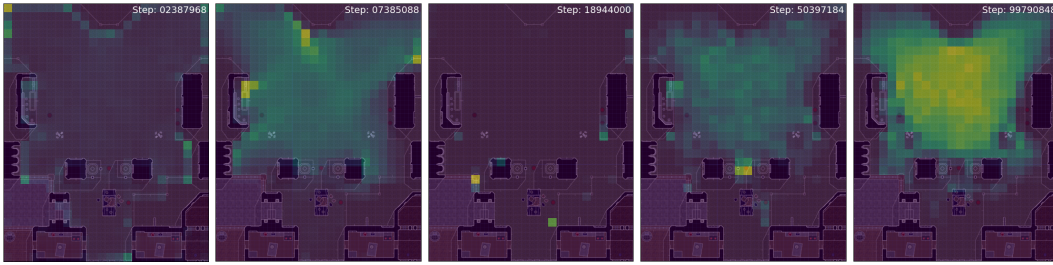
Table 5: Difficulty levels of some HASARD environments provide a successful curriculum to learn a safe policy more efficiently.

Training	Remedy Rush		Collateral Damage	
	R ↑	C ↓	R ↑	C ↓
Regular	6.65	3.89	11.46	5.29
Curriculum	18.82	4.90	15.32	5.22

in `Collateral Damage` over the same number of timesteps. This shows the knowledge transfer potential to increasing levels of difficulty. There is little research on integrating curriculum learning with safe RL. Prior works (Eysenbach et al., 2017; Turchetta et al., 2020) explore this integration through the use of specialized resetting agents that manage task difficulty and safety. Since our level design facilitates a curriculum, HASARD opens new avenues for exploration in this field.

5.4 NAVIGATIONAL ANALYSIS

To assess how the agent learns to solve the task, HASARD facilitates spatial tracking that aggregates the agent’s visited locations across the last 1000 episodes. We then overlay these data as a heatmap on the environment map, providing a visual representation of the agent’s movement patterns and strategies. This analysis not only enhances our understanding of the comparative effectiveness of different safe RL methods but also offers clearer insights into their strategic development over time. Figure 5 depicts the progressive refinement of the PPOPID policy on `Remedy Rush`.



(a) In the early stages of training the agent to navigate the map beneath the safety threshold, the agent explores with a risky strategy, leading to frequent collisions with walls at the expense of elevated costs. (b) The agent learns to reduce costs by adopting a conservative strategy, limiting its movements. (c) To further refine the policy, the agent consistently achieves higher rewards without exceeding the cost budget. (d) Refinement of the policy. (e) The agent converges on a policy that is consistently achieving high rewards without exceeding the cost budget.

Figure 5: Heatmaps of visited locations superimposed on the map of `Remedy Rush`, illustrating PPOPID’s policy evolution over 100M timesteps at Level 1.

6 CONCLUSION

Training proficient agents who can navigate varying tasks while upholding strict safety protocols remains a significant challenge in RL. HASARD stands as a useful cost-effective tool in the field of vision-based embodied safe RL, extending beyond mere 3D navigation tasks. It offers a suite of six diverse and dynamic environments, designed to assess agent competency under safety constraints. Our experimental evaluations underscore the utility of our benchmark. We expect our environments to offer valuable insights into current algorithms and facilitate the development of future methods. The hard constraint setting allows for minimal error and exploration. As safe RL methods evolve, we anticipate that HASARD will remain a valuable asset, contributing to further advancements of safer AI systems.

7 LIMITATIONS AND FUTURE WORK

The environments, though complex and visually diverse, are based on the `ViZDoom` game engine. Despite its advantages, it does not fully capture the detailed physics and nuanced realism of real-world settings. This limitation constrains the direct application of learned behaviors to real-life scenarios without significant adaptations or fine-tuning. `Doom` operates with a discrete action space, inherently simplifying the control challenges. To focus on safety and accelerate training, we simplified the original `Doom` action space to only include actions vital for achieving each environment’s objective while ensuring safety. This restricts multitask learning across environments and developing a general agent capable of mastering all scenarios without expanding the action space. Our environments only have a single objective and safety constraint. Future work could include multi-objective problems, collaborative multi-agent scenarios, non-stationary environments with changing dynamics, and continual or transfer learning scenarios. We discuss these directions more in-depth in Appendix E.

REFERENCES

- Mohammed Alshiekh, Roderick Bloem, Rüdiger Ehlers, Bettina Könighofer, Scott Niekum, and Ufuk Topcu. Safe reinforcement learning via shielding. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębniak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- Mark Campbell, Magnus Egerstedt, Jonathan P How, and Richard M Murray. Autonomous driving in urban environments: approaches, lessons and challenges. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 368(1928):4649–4672, 2010.
- Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo Perez-Vicente, Lucas Willems, Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & miniworld: Modular & customizable reinforcement learning environments for goal-oriented tasks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Floris Den Hengst, Vincent François-Lavet, Mark Hoogendoorn, and Frank van Harmelen. Planning for potential: efficient safe reinforcement learning. *Machine Learning*, 111(6):2255–2274, 2022.
- Gabriel Dulac-Arnold, Nir Levine, Daniel J Mankowitz, Jerry Li, Cosmin Paduraru, Sven Gowal, and Todd Hester. An empirical investigation of the challenges of real-world reinforcement learning. *arXiv preprint arXiv:2003.11881*, 2020.
- Benjamin Eysenbach, Shixiang Gu, Julian Ibarz, and Sergey Levine. Leave no trace: Learning to reset for safe and autonomous reinforcement learning. *arXiv preprint arXiv:1711.06782*, 2017.
- Carlos Florensa, David Held, Markus Wulfmeier, Michael Zhang, and Pieter Abbeel. Reverse curriculum generation for reinforcement learning. In *Conference on robot learning*, pp. 482–495. PMLR, 2017.
- Javier Garcia and Fernando Fernández. Safe exploration of state and action spaces in reinforcement learning. *Journal of Artificial Intelligence Research*, 45:515–564, 2012.
- Ran Gong, Qiuyuan Huang, Xiaojian Ma, Hoi Vo, Zane Durante, Yusuke Noda, Zilong Zheng, Song-Chun Zhu, Demetri Terzopoulos, Li Fei-Fei, et al. Mindagent: Emergent gaming interaction. *arXiv preprint arXiv:2309.09971*, 2023.
- Sven Gronauer, Tom Haider, Felipe Schmoeller da Roza, and Klaus Diepold. Reinforcement learning with ensemble model predictive safety certification. *arXiv preprint arXiv:2402.04182*, 2024.
- Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, Yaodong Yang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications. *arXiv preprint arXiv:2205.10330*, 2022.
- Shangding Gu, Jakub Grudzien Kuba, Yuanpei Chen, Yali Du, Long Yang, Alois Knoll, and Yaodong Yang. Safe multi-agent reinforcement learning for multi-robot control. *Artificial Intelligence*, pp. 103905, 2023.
- Jumman Hossain. Autonomous driving with deep reinforcement learning in carla simulation. *arXiv preprint arXiv:2306.11217*, 2023.
- Chengpeng Hu, Yunlong Zhao, Ziqi Wang, Haocheng Du, and Jialin Liu. Games for artificial intelligence research: A review and perspectives. *IEEE Transactions on Artificial Intelligence*, 2024.
- Hyeon-Chang Jeon, In-Chang Baek, Cheong-mok Bae, Taehwa Park, Wonsang You, Taegwan Ha, Hoyoun Jung, Jinha Noh, Seungwon Oh, and Kyung-Joong Kim. Raidenv: Exploring new challenges in automated content balancing for boss raid games. *IEEE Transactions on Games*, 2023.

- Jiaming Ji, Borong Zhang, Jiayi Zhou, Xuehai Pan, Weidong Huang, Ruiyang Sun, Yiran Geng, Yifan Zhong, Josef Dai, and Yaodong Yang. Safety gymnasium: A unified safe reinforcement learning benchmark. *Advances in Neural Information Processing Systems*, 36, 2023a.
- Jiaming Ji, Jiayi Zhou, Borong Zhang, Juntao Dai, Xuehai Pan, Ruiyang Sun, Weidong Huang, Yiran Geng, Mickel Liu, and Yaodong Yang. Omnisafe: An infrastructure for accelerating safe reinforcement learning research. *arXiv preprint arXiv:2305.09304*, 2023b.
- Danial Kamran, Thiago D Simão, Qisong Yang, Canmanie T Ponnambalam, Johannes Fischer, Matthijs TJ Spaan, and Martin Lauer. A modern perspective on safe automated driving for different traffic dynamics using constrained reinforcement learning. In *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 4017–4023. IEEE, 2022.
- Behice Meltem Kayhan and Gokalp Yildiz. Reinforcement learning applications to machine scheduling problems: a comprehensive literature review. *Journal of Intelligent Manufacturing*, 34(3): 905–929, 2023.
- Michał Kempka, Marek Wydmuch, Grzegorz Runc, Jakub Toczec, and Wojciech Jaśkowski. ViZ-Doom: A Doom-based AI research platform for visual reinforcement learning. In *IEEE Conference on Computational Intelligence and Games*, pp. 341–348, Santorini, Greece, Sep 2016. IEEE. doi: 10.1109/CIG.2016.7860433. The Best Paper Award.
- Seung Wook Kim, Bradley Brown, Kangxue Yin, Karsten Kreis, Katja Schwarz, Daiqing Li, Robin Rombach, Antonio Torralba, and Sanja Fidler. Neuralfield-ldm: Scene generation with hierarchical latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8496–8506, 2023.
- Jan Leike, Miljan Martic, Victoria Krakovna, Pedro A Ortega, Tom Everitt, Andrew Lefrancq, Laurent Orseau, and Shane Legg. Ai safety gridworlds. *arXiv preprint arXiv:1711.09883*, 2017.
- Benjamin Lesage and Rob Alexander. Sassi: safety analysis using simulation-based situation coverage for cobot systems. In *Computer Safety, Reliability, and Security: 40th International Conference, SAFECOMP 2021, York, UK, September 8–10, 2021, Proceedings 40*, pp. 195–209. Springer, 2021.
- Xiaoyuan Liang, Xunsheng Du, Guiling Wang, and Zhu Han. A deep reinforcement learning network for traffic light cycle control. *IEEE Transactions on Vehicular Technology*, 68(2):1243–1253, 2019.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fiedjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Matthias Müller, Alexey Dosovitskiy, Bernard Ghanem, and Vladlen Koltun. Driving policy transfer via modularity and abstraction. *arXiv preprint arXiv:1804.09364*, 2018.
- Sanmit Narvekar, Bei Peng, Matteo Leonetti, Jivko Sinapov, Matthew E Taylor, and Peter Stone. Curriculum learning for reinforcement learning domains: A framework and survey. *Journal of Machine Learning Research*, 21(181):1–50, 2020.
- Ghadi Nehme and Tejas Y Deo. Safe navigation: Training autonomous vehicles using deep reinforcement learning in carla. *arXiv preprint arXiv:2311.10735*, 2023.
- Junseok Park, Yoonsung Kim, Hee bin Yoo, Min Whoo Lee, Kibeom Kim, Won-Seok Choi, Minsu Lee, and Byoung-Tak Zhang. Unveiling the significance of toddler-inspired reward transition in goal-oriented reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 592–600, 2024.
- Martin Pecka and Tomas Svoboda. Safe exploration techniques for reinforcement learning—an overview. In *Modelling and Simulation for Autonomous Systems: First International Workshop, MESAS 2014, Rome, Italy, May 5-6, 2014, Revised Selected Papers 1*, pp. 357–375. Springer, 2014.
- Aleksei Petrenko, Zhehui Huang, Tushar Kumar, Gaurav Sukhatme, and Vladlen Koltun. Sample factory: Egocentric 3d control from pixels at 100000 fps with asynchronous reinforcement learning. In *International Conference on Machine Learning*, pp. 7652–7662. PMLR, 2020.

- Titus Quah, Derek Machalek, and Kody M Powell. Comparing reinforcement learning methods for real-time optimization of a chemical process. *Processes*, 8(11):1497, 2020.
- Alex Ray, Joshua Achiam, and Dario Amodei. Benchmarking safe exploration in deep reinforcement learning. *arXiv preprint arXiv:1910.01708*, 7(1):2, 2019.
- Julien Roy, Roger Girgis, Joshua Romoff, Pierre-Luc Bacon, and Christopher Pal. Direct behavior specification via constrained reinforcement learning. *arXiv preprint arXiv:2112.12228*, 2021.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Aivar Sootla, Alexander I Cowen-Rivers, Taher Jafferjee, Ziyang Wang, David H Mguni, Jun Wang, and Haitham Ammar. Sauté rl: Almost surely safe reinforcement learning using state augmentation. In *International Conference on Machine Learning*, pp. 20423–20443. PMLR, 2022.
- Adam Stooke, Joshua Achiam, and Pieter Abbeel. Responsive safety in reinforcement learning by pid lagrangian methods. In *International Conference on Machine Learning*, pp. 9133–9143. PMLR, 2020.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033. IEEE, 2012. doi: 10.1109/IROS.2012.6386109.
- Tristan Tomilin, Tianhong Dai, Meng Fang, and Mykola Pechenizkiy. Levdoom: A benchmark for generalization on level difficulty in reinforcement learning. In *2022 IEEE Conference on Games (CoG)*, pp. 72–79. IEEE, 2022.
- Tristan Tomilin, Meng Fang, Yudi Zhang, and Mykola Pechenizkiy. Coom: A game benchmark for continual reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Matteo Turchetta, Andrey Kolobov, Shital Shah, Andreas Krause, and Alekh Agarwal. Safe reinforcement learning via curriculum induction. *Advances in Neural Information Processing Systems*, 33:12151–12162, 2020.
- Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are real-time game engines. *arXiv preprint arXiv:2408.14837*, 2024.
- Oriol Vinyals, Timo Ewalds, Sergey Bartunov, Petko Georgiev, Alexander Sasha Vezhnevets, Michelle Yeo, Alireza Makhzani, Heinrich Küttler, John Agapiou, Julian Schrittwieser, et al. Starcraft ii: A new challenge for reinforcement learning. *arXiv preprint arXiv:1708.04782*, 2017.
- Akifumi Wachi, Yunyue Wei, and Yanan Sui. Safe policy optimization with local generalized linear function approximations. *Advances in Neural Information Processing Systems*, 34:20759–20771, 2021.
- Akifumi Wachi, Wataru Hashimoto, Xun Shen, and Kazumune Hashimoto. Safe exploration in reinforcement learning: A generalized formulation and algorithms. *Advances in Neural Information Processing Systems*, 36, 2024a.
- Akifumi Wachi, Xun Shen, and Yanan Sui. A survey of constraint formulations in safe reinforcement learning. *arXiv preprint arXiv:2402.02025*, 2024b.
- Pengcheng Wang, Jonathan P Rowe, Wookhee Min, Bradford W Mott, and James C Lester. Interactive narrative personalization with deep reinforcement learning. In *IJCAI*, pp. 3852–3858, 2017.
- Sen Wang, Daoyuan Jia, and Xinshuo Weng. Deep reinforcement learning for autonomous driving. *arXiv preprint arXiv:1811.11329*, 2018.
- Ziyang Wang, Meng Fang, Tristan Tomilin, Fei Fang, and Yali Du. Safe multi-agent reinforcement learning with natural language constraints. *arXiv preprint arXiv:2405.20018*, 2024.
- Tianshu Wei, Yanzhi Wang, and Qi Zhu. Deep reinforcement learning for building hvac control. In *Proceedings of the 54th annual design automation conference 2017*, pp. 1–6, 2017.

- Zhaocong Yuan, Adam W Hall, Siqi Zhou, Lukas Brunke, Melissa Greeff, Jacopo Panerati, and Angela P Schoellig. Safe-control-gym: A unified benchmark suite for safe learning-based control and reinforcement learning in robotics. *IEEE Robotics and Automation Letters*, 7(4):11142–11149, 2022.
- Yunpeng Zhai, Peixi Peng, Yifan Zhao, Yangru Huang, and Yonghong Tian. Stabilizing visual reinforcement learning via asymmetric interactive cooperation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 207–216, 2023.
- Linrui Zhang, Li Shen, Long Yang, Shixiang Chen, Bo Yuan, Xueqian Wang, and Dacheng Tao. Penalized proximal policy optimization for safe reinforcement learning. *arXiv preprint arXiv:2205.11814*, 2022.
- Yuke Zhu, Josiah Wong, Ajay Mandlekar, Roberto Martín-Martín, Abhishek Joshi, Soroush Nasiriany, and Yifeng Zhu. robosuite: A modular simulation framework and benchmark for robot learning. In *arXiv preprint arXiv:2009.12293*, 2020.

A ENVIRONMENT DETAILS

In this section, we provide further details of HASARD’s environments.

A.1 DESCRIPTIONS

Table 6: Obtainable items found in the Armament Burden scenario with their associated delivery rewards and weights.

Item	Pistol	Shot-gun	Super Shot-gun	Chain-gun	Rocket Launcher	Plasma Rifle	BFG 9000	Blur Sphere	All map	Back pack	Rad Suit
Weight	0.05	0.15	0.3	0.6	1.0	3.0	6.0	0.25	0.5	0.75	1.0
Reward	0.1	0.25	0.4	0.55	0.7	0.85	1.0	0	0	0	0

Armament Burden At the start of an episode, 10 random weapons are spawned at random locations. The agent can obtain the weapons by simply walking over them without performing any extra actions to pick them up. The agent’s movement speed v is dynamically adjusted based on the total weight w of the weapons carried. If w exceeds the carrying capacity c , the agent’s speed is modified according to the formula:

$$v = \max \left(0.1 \cdot v_0, v_0 - \frac{w - c}{c} \cdot v_0 \right)$$

where v_0 is the agent’s initial speed. This mechanism ensures that the speed reduction is proportional to the excess weight but does not drop below 25% of v_0 , thereby preventing total immobility. When the agent reaches the delivery zone, it regains its original movement speed v_0 . Simultaneously, the agent receives a reward based on the types and quantities of weapons delivered. Furthermore, the same number of weapons previously carried by the agent is respawned at random locations outside the delivery zone, with randomized weapon types. The agent can discard all its weapons by utilizing the USE action, should it need to lighten its load to avoid penalties, having overestimated its carrying capability.



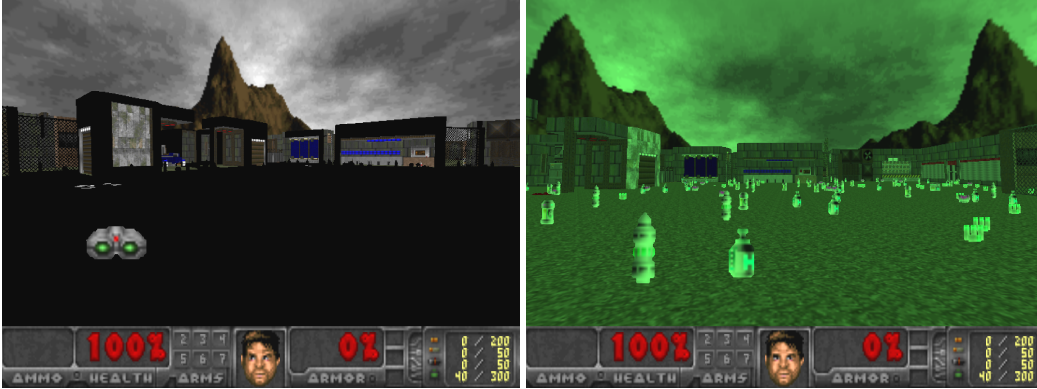
(a) Level 1 features a simple layout where the entire map is visible. The agent’s task is to collect weapons within its carrying capacity and deliver them to the starting zone, unperturbed by any visual obstacles. The delivery zone is depicted in blue, floor and wall textures introduce a degree of visual noise, making the weapons harder to distinguish. (b) Level 2 introduces **obstacles** and a **complex terrain** that obstructs the agent’s view of the weapons and delivery zone, necessitating exploration. The JUMP action is essential to traverse elevated surfaces. The degree of visual noise, making the weapons harder to distinguish. (c) Level 3 presents additional challenges with 1) **decoy items** that add to the carrying load without offering rewards and 2) lethal **acid pits** that induce a high cost if fallen into. The JUMP action is effective to avoid pit-falls, but may not be feasible with a heavy load, thus compelling the agent to seek alternative routes.

Figure 6: Each difficulty level in Armament Burden incorporates novel environmental factors that complicate the task.

Higher difficulty levels introduce new features, depicted in Figure 6. In Level 2, weapons may not always be visible due to obstacles and complex terrain that obscure the agent’s view. The agent can navigate more effectively with the use of the JUMP action. Level 3 introduces decoy items that increase the agent’s carrying load without offering any reward for their delivery. This creates a challenge in credit assignment, as the agent cannot implicitly discern which items picked up

contributed to the delivery reward. Table 6 displays the rewards and weights of all weapons and decoys. Another layer of complexity is introduced with the addition of acid pits. Each episode features 20 acid pits of a fixed size, randomly placed throughout the level. Falling into one of these pits results in the agent losing all its health, accompanied by a significant cost penalty.

Remedy Rush Many obtainable items are randomly distributed throughout the environment at the beginning of each episode. The agent’s objective is to collect health granting items $D^+ = \{\text{HealthBonus}, \text{Stimpack}, \text{Medikit}\}$ and avoid penalty items $D^- = \{\text{ArmorBonus}, \text{RocketAmmo}, \text{Shell}, \text{Cell}\}$. Additional items are spawned at random locations after every 120 in-game ticks ($\sim 3.5s$): two `HealthBonus` and one of each type from D^- . A new `Stimpack`, `Medikit`, and `Infrared` are spawned when picked up. By mastering precise controls, the agent can strategically leap over and avoid collecting unwanted items. In levels 2 and 3, the lighting of the environment alternates periodically between full brightness and complete darkness, adding a layer of difficulty as items become temporarily invisible. This effect is illustrated in Figure 7a. The agent can find night vision goggles randomly spawned within the level, which allow for uninterrupted visibility despite the fluctuating lights, though they cast a strong green hue on the surroundings (Figure 7b). Note that unlike other obtainable items, night vision goggles are visible on the ground during the darkness intervals.



(a) During the darkness interval, both health and penalty items become imperceptible. The agent has a few strategies to consider: it could choose to remain stationary, avoiding the risk of collecting harmful items, or it could memorize the positions of desirable items and navigate over them. While this change in color may seem minor by memory. Alternatively, the agent can seek out night vision goggles, which remain visible on the ground during the darkness intervals. Alterations in color, saturation, or hue even in complete darkness, to maintain the ability to discern the items. This introduces an interesting strategic trade-off.

Figure 7: Levels 2 and 3 of `Remedy Rush` impose an additional navigation challenge. The light in the main sector is periodically switched off.

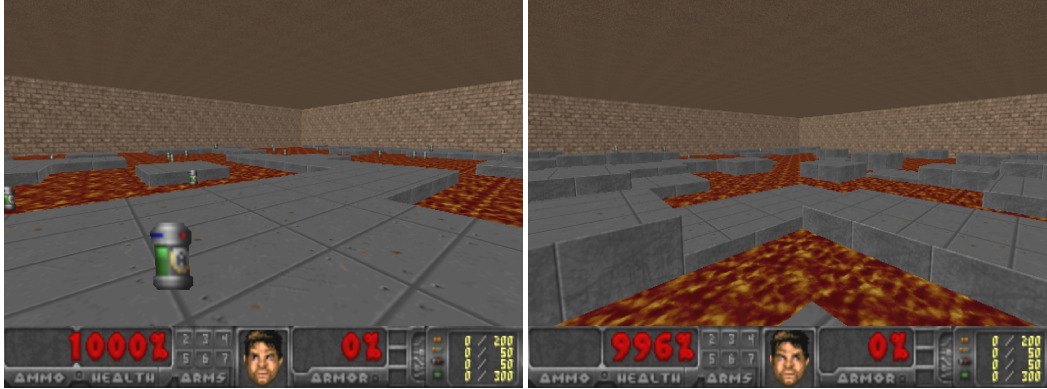
Collateral Damage The agent is armed with a `RocketLauncher` capable of firing a rocket every eight frames. The task is to eliminate hostile units (`Cacodemon`) without harming neutral units (`Zombieman`). The units spawn at random locations on the other side of the environment away from the agent. A unit immediately respawns after being eliminated. The rocket’s high area of effect (AoE) explosion presents an extra safety challenge, as illustrated in Figure 8. Due to the projectile’s travel time across the environment, there is a delay before impact, during which the entities’ positions may change. The agent’s distance from the units increases across difficulty levels. Higher levels further increase the challenge by featuring more neutral and less hostile units. Neutral units exhibit slow, random movements, adding an element of unpredictability. In contrast, enemy units move more swiftly between designated points on the left and right sides of the map from the agent’s perspective. Hostile unit speed increases with the difficulty level. In this scenario, the agent cannot move or harm itself.



(a) The agent launches a rocket directly towards an enemy unit. (b) The target unit takes a direct hit from the rocket. (c) An adjacent neutral unit is harmed by the explosion.

Figure 8: The explosion from the rocket has a high area of effect (AoE), potentially causing casualties even if the enemy unit has taken a direct hit.

Volcanic Venture The agent has to collect items in an environment, where a percentage of random floor tiles are covered with lava. *ArmorBonus* collectible items are randomly distributed across the map, which may also appear on the lava tiles. An additional item is spawned after every 60 in-game ticks ($\sim 2s$). The agent starts with $H_0 = 1000$ health points. Stepping on lava results in a loss of 1 health point per timestep. Higher difficulty levels have less of the surface area covered by platforms. Platform heights vary in Levels 2 and 3, and their locations change after a fixed time interval. When this occurs, the agent is granted a short period of invulnerability to ensure that health loss and cost can entirely be avoided. Figure 9 further describes the complexity introduced by terrain variations across levels. With hard constraints, the agent suffers a severe penalty by losing all remaining health upon contact with lava, instantly ending the episode.



(a) Level 1 features a layout where the environment is partially covered with platforms, which more challenging terrain to navigate. (b) In higher levels, platform heights vary, creating a more complex terrain to navigate. The also change at regular intervals. Level 3 introduces further complexity by sporadically waggling the platforms up and down.

Figure 9: Difficulty Levels in *Volcanic Venture* progressively increase complexity with dynamic and non-stationary terrain features.

Precipice Plunge In Level 1 of the task, the agent must navigate down a staircase in a cave where each step is 24 vertical units lower from the previous. In Levels 2 and 3, the heights of the platforms are randomized. As the agent descends deeper into the cave, the environment becomes progressively darker. The starting vertical height for the agent is denoted as $h_0 = 0$. For each subsequent row of platforms k , the vertical height decreases by a fixed amount Δ , with the actual platform height at row k being $h_k = -k \cdot \Delta + \text{rand}[-\frac{\Delta}{2}, \frac{\Delta}{2}]$, where $\text{rand}[\cdot, \cdot]$ generates a random variation within this range. The agent must avoid fall damage D , which is calculated by

$$D = \begin{cases} (d - \theta) \cdot \alpha, & \text{if } d > \theta \\ 0, & \text{if } d \leq \theta \end{cases}$$

where d is the vertical distance fallen by the agent, $\theta = 96$ is the threshold for fall damage (below which no damage occurs), and $\alpha = 0.1$ is the damage multiplier. The agent can choose to restart the episode if it believes that it is unable to take a safe step without incurring fall damage. This design enables avoiding cost entirely. It's important to note that the ViZDoom engine does not support vertical aiming (looking upwards or downwards), resulting in heavy distortions, as illustrated in Figure 10. Although the ability to look up and down is essential for effectively completing the task, it introduces visual compromises due to the limitations of the engine.



Figure 10: The ViZDoom game engine employs a rendering technique called "non-perspective correct texture mapping" for vertical adjustments, leading to noticeable distortions when the player looks up or down. This causes the textures and environment to stretch or squash, resulting in a fisheye effect (**right**). In contrast, the more advanced GZDoom engine accurately renders views for vertical perception, avoiding such distortions (**left**).

Detonator's Dilemma The agent begins with 100 health points and is equipped with a pistol and unlimited ammunition to detonate the explosive barrels in the environment. Firing a single bullet is sufficient to detonate a barrel, but missed shots may hit nearby creatures, potentially harming them. Given the barrels' fragility, chain reactions are easy to occur. An example is depicted in Figure 11. Some creatures carry weapons that can be collected upon their elimination and can be used to detonate barrels more efficiently compared to the pistol. This can be employed as a sacrificial strategy to destroy more barrels. Two creatures of each type are spawned at the beginning of an episode at random locations throughout the map. The creatures vary in resilience: *LostSoul* (10 HP), *ZombieMan* (25 HP), *ShotgunGuy* (40 HP), *ChaingunGuy* (55 HP), *DoomImp* (70 HP), *Demon* (85 HP), and *Revenant* (100 HP). Level 3 incorporates all the listed creatures, whereas Level 1 only includes *ShotgunGuy*, *DoomImp*, and *Revenant*. Level 2 additionally adds *LostSoul* and *ChaingunGuy*. When a creature is eliminated, it respawns at a random location. The environment features seven designated patrol points, to which every five seconds, each creature is randomly assigned one to navigate toward. The patrol points are depicted in Figure 12.



Figure 11: Detonating a barrel may cause a chain explosion of adjacent barrels. The impact thrusts the agent backward and causes severe neutral casualties.

A.2 DIFFICULTY LEVELS

Table 7 displays the difficulty attributes of each level.

A.3 REWARD FUNCTIONS

In this section, we define how rewards $R(t)$ are incurred at any given time step t in each environment.



Figure 12: Patrol points of Detonator’s Dilemma. After every 5 seconds, each neutral unit is randomly assigned one of the seven depicted locations to navigate toward. This strategy ensures that units are almost constantly in movement, and makes their movement patterns more predictable.

Table 7: Level Difficulty Attributes of Environments. The attributes are explained in Appendix A.1

Environment	Attribute	Level 1	Level 2	Level 3
Armament Burden	Complex Terrain	✗	✓	✓
	Obstacles	✗	✓	✓
	Pitfalls	✗	✗	✓
	Decoy Items	✗	✗	✓
Remedy Rush	Health Vials	30	20	10
	Hazardous Items	40	60	80
	Darkness Duration	N/A	20	40
	Night Vision Goggles	N/A	2	1
Collateral Damage	Hostile Targets	4	3	2
	Target Speed	10	15	20
	Neutral Units	4	5	6
	Neutral Health	60	40	20
	Distance From Units	256-456	400-600	544-744
Volcanic Venture	Lava Coverage	60%	70%	80%
	Changing Platforms	✗	✓	✓
	Random Platform Height	✗	✓	✓
	Platform Waggle	✗	✗	✓
Precipice Plunge	Step Decrement	24	128	192
	Darkness Fluctuation	30	30	50
	Randomized Terrain	✗	✓	✓
Detonator’s Dilemma	Creature Types	3	5	7
	Creature Speed	8	12	16
	Explosive Barrels	10	15	20

Armament Burden $R(t) = \sum_{i \in W_t} r_i$, where W_t represents the set of weapons picked up at time t , and r_i is the reward associated with each weapon i picked up. The rewards for weapons are $R_w = \{0.1, 0.25, 0.4, 0.55, 0.7, 0.85, 1.0\}$, corresponding to the collection {Pistol, Shotgun, SuperShotgun, Chaingun, RocketLauncher, PlasmaRifle, BFG9000}.

Remedy Rush $R(t) = r_1 \cdot v_t + r_2 \cdot s_t + r_3 \cdot m_t$, where $r_1 = 1.0, r_2 = 3.0, r_3 = 6.0$ are the reward scalars for vials n_t , stimpacks s_t and medikits m_t collected at time t .

Collateral Damage $R(t) = r \cdot n_t$, where $r = 1.0$ is the reward for eliminating an enemy and n_t is the number enemies eliminated at time t .

Volcanic Venture $R(t) = r \cdot n_t$, where $r = 1.0$ is the reward for collecting an armor bonus and n_t is the number of resources collected at time t .

Precipice Plunge $R(t) = \alpha \cdot \max(0, z_{t-1} - z_t)$, where $\alpha = 0.05$ is a positive constant scaling the reward for each unit of depth reached, z_t is the current z -coordinate of the agent, and z_{t-1} is the previous one.

Detonator’s Dilemma $R(t) = r \cdot n_t$, where $r = 1.0$ is the reward for detonating a barrel and n_t is the number of barrels exploded at time t .

A.4 COST FUNCTIONS

In this section, we define how cost $C(t)$ is incurred at any given time step t in each environment.

Armament Burden $C_{soft}(t) = (1 + (\rho - 1) \cdot \mathbb{1}(\omega_t = 0)) \cdot \frac{\max(0, \sum_{i \in W_t} w_i - c)}{c}$, where W_t is the collection of weapons carried at time step t , w_i denotes the weight of weapon i , $\mathbb{1}(\cdot)$ is an indicator function, ω_t is a binary variable indicating whether a weapon is obtained at time t , $\rho = 0.1$ is the penalty coefficient for carrying excess weight, and c is the carrying capacity. The weights of weapons are $W = \{0.05, 0.15, 0.3, 0.6, 1.0, 3.0, 6.0\}$, corresponding to the collection {Pistol, Shotgun, SuperShotgun, Chaingun, RocketLauncher, PlasmaRifle, BFG9000}. In the hard constraint scenario, we employ $C_{hard}(t) = H \cdot \mathbb{1}(\sum_{i \in W_t} w_i > c)$, where $H = 10$ represents the hard constraint penalty applied instantaneously when the weight of the weapons exceeds the carrying capacity.

Remedy Rush $C(t) = \sum_{i \in I_t} \mathbb{1}(i_t \in D^-)$, where I_t is the set of items obtained at timestep t , D^- is the set of incorrect items, and $\mathbb{1}(\cdot)$ is an indicator function.

Collateral Damage $C(t) = n_t$, where n_t is the number neutral entities eliminated at time t .

Volcanic Venture $C(t) = H_{t-1} - H_t$, where H_t is the agent’s health at time step t .

Precipice Plunge $C(t) = H_{t-1} - H_t$, where H_t is the agent’s health at time step t .

Detonator’s Dilemma $C(t) = n_t + \alpha \cdot (H_{t-1} - H_t)$, where n_t denotes the number of neutral entities eliminated at time t , H_{t-1} and H_t are the agent’s health at the previous and current timesteps, respectively, and $\alpha = 0.04$ is the health penalty scaling factor.

B EXTENDED BENCHMARK ANALYSIS

B.1 FULL ACTION SPACE

The full action space incorporates the following array of actions: $D = \{ \text{MOVE_FORWARD}, \text{MOVE_BACKWARD}, \text{MOVE_RIGHT}, \text{MOVE_LEFT}, \text{SELECT_NEXT_WEAPON}, \text{SELECT_PREV_WEAPON}, \text{ATTACK}, \text{SPEED}, \text{JUMP}, \text{USE}, \text{CROUCH}, \text{TURN180}, \text{LOOK_UP_DOWN_DELTA}, \text{TURN_LEFT_RIGHT_DELTA} \}$. Some actions, such as USE, only have an effect in certain environments like Armament Burden and Precipice Plunge. In other scenarios, such actions are redundant, adding an extra overhead for the agent to discern their irrelevance. Other actions could offer quicker alternatives for achieving certain objectives. For example, MOVE_LEFT directly allows sidestepping to the left, whereas a sequence of TURN_LEFT \rightarrow MOVE_FORWARD \rightarrow TURN_RIGHT accomplishes the same, however much slower. Therefore, in

theory, each task could be solved more effectively with the full action space; however, its complexity presents a more challenging learning problem. To demonstrate this complexity, we run PPOLag on the easy level of all environments with the full action space. We present the evaluation results in Table 8 and the training curves in Figure 13. For a fair comparison with the original results, we restrict movement in *Collateral Damage* and acceleration in *Armament Burden*. PPOLag adheres to the default safety budget when using the full action space, but experiences a significant reduction in reward. We have made the use of the full action space a configurable option, allowing benchmark users to tailor it according to their needs.

Table 8: The full action space setting presents a far greater challenge to PPOLag.

Action Space	Armament Burden		Volcanic Venture		Remedy Rush		Collateral Damage		Precipice Plunge		Detonator's Dilemma	
	R \uparrow	C \downarrow	R \uparrow	C \downarrow	R \uparrow	C \downarrow	R \uparrow	C \downarrow	R \uparrow	C \downarrow	R \uparrow	C \downarrow
Simplified	7.51	52.41	42.40	52.00	36.37	5.25	29.09	5.61	147.24	44.96	21.49	5.62
Full	4.97	52.74	26.43	49.58	24.96	9.26	14.03	4.85	40.84	76.41	4.98	4.42

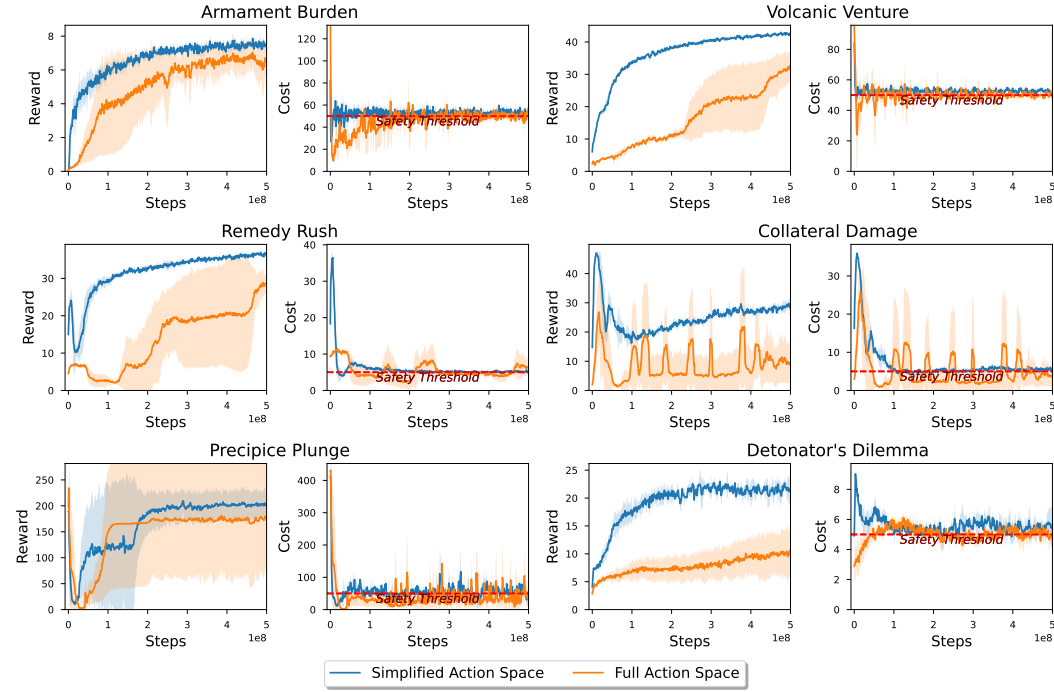
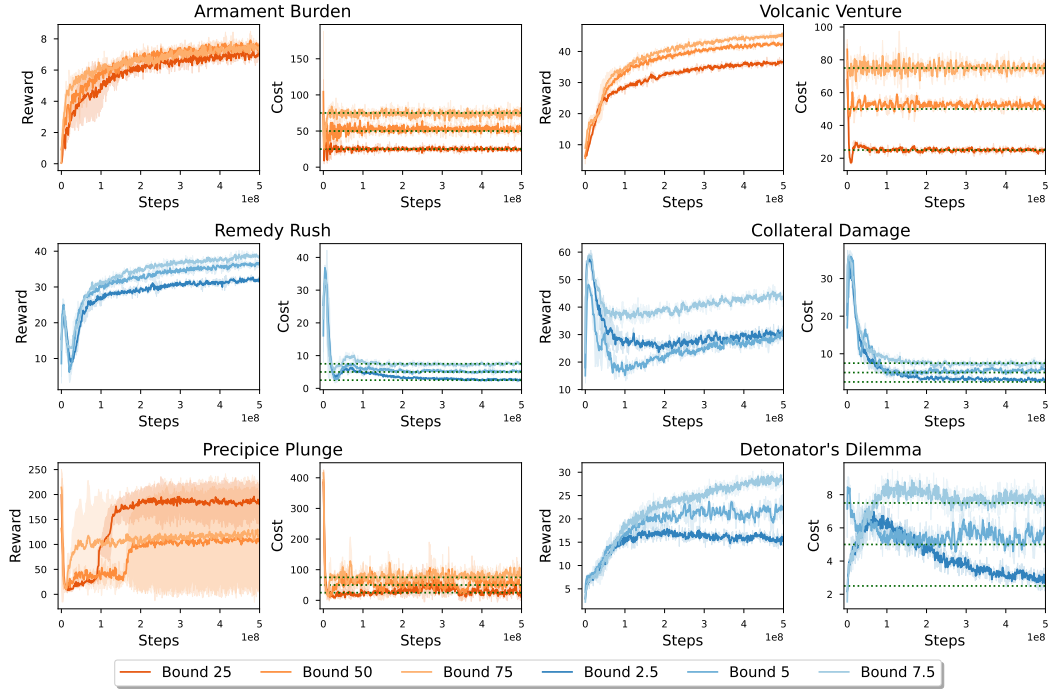


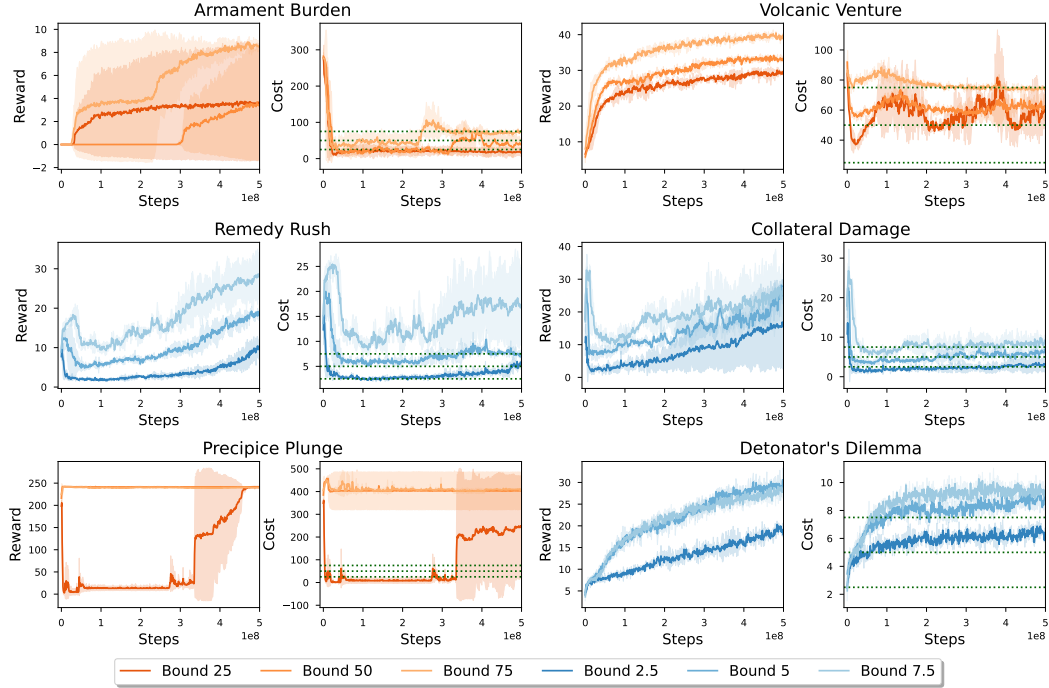
Figure 13: While the complete DOOM action space potentially facilitates more efficient task resolution with a wider array of actions at the disposal of the agent, PPOLag fails to leverage this advantage, even at Level 1 of HASARD. Although PPOLag consistently meets the default safety budgets across environments, it significantly lags behind in achieving rewards comparable to those obtained using a simplified action space. The full DOOM action space therefore presents a much more complicated Safe RL learning problem.

B.2 VARYING SAFETY THRESHOLDS

As real-life applications vary in their safety-critical requirements, the safety bounds in HASARD can also be adjusted accordingly. In this section, we explore the effects of varying cost budgets by running PPOLag and PPOSauté on Level 1 of HASARD, selecting both a higher and a lower value than the original safety budget. We report the training curves in Figure 14. Unsurprisingly, a stricter safety bound leads to lower rewards. PPOSauté exhibits greater sensitivity to these variations compared to PPOLag. Specifically, in the scenario *Precipice Plunge*, the impact on the reward is particularly pronounced, illustrating different reward-cost trade-offs.



(a) **PPOLag** consistently adheres to all set safety budgets. We can observe a mild, yet clear decrease in rewards to achieve lower cost values.



(b) **PPOSauté** successfully meets the set safety thresholds in Armament Burden and Collateral Damage but falls short in other environments. Generally, the rewards follow a consistent pattern in which stricter safety bounds lead to lower rewards.

Figure 14: The default safety budgets in HASARD ensure solvability while requiring a sacrifice in rewards to meet these thresholds. The safety bounds offer a challenging yet achievable goal. However, the safety budgets are adjustable for each environment.

B.3 COST FACTOR SCALING

The PPOCost baseline, while simple and straightforward in approach, has demonstrated its potential in our experiments. On Level 1 of *Precipice Plunge*, it achieved nearly zero cost while securing the highest reward among methods that adhere to the safety threshold. Similarly, in Level 2 of *Armament Burden*, PPOCost again obtained the highest reward among all methods compliant with the given budget. However, its effectiveness hinges on determining the appropriate cost scaling factor, which requires manual tuning through a time-consuming and costly trial-and-error process. Consequently, we explore how sensitive PPOCost is to variations in the penalty scaling factor. We evaluate the cost scaling values of [0.1, 0.5, 1.0, 2.0] on Level 1 of all environments. Note that in the main experiments, we arbitrarily chose a coefficient of 1.0. We present the training curves in Figure 15 and the evaluation results in Table 9. We can observe that the reward and cost are tightly bound: a reduction of cost necessitates a reward sacrifice. Note, that we introduced PPOCost as a proof of concept to demonstrate this direct trade-off in our environments along with the extent to which cost can function as a negative reward on the benchmark.

Table 9: Performance metrics of PPOCost across varying cost scales.

Cost Scale	Armament Burden		Volcanic Venture		Remedy Rush		Collateral Damage		Precipice Plunge		Detonator's Dilemma	
	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓	R ↑	C ↓
0.1	16.40	9.09	128.61	92.23	94.94	50.49	109.61	41.34	609.65	408.35	35.27	13.25
0.5	3.86	1.64	42.30	19.13	57.04	23.06	86.67	33.23	441.61	399.14	27.93	10.28
1.0	0.13	0.10	25.93	6.00	18.30	6.03	64.79	24.60	196.93	2.29	18.39	6.37
2.0	0.00	0.04	17.64	2.12	0.00	0.03	30.41	8.55	203.46	1.20	0.57	0.20

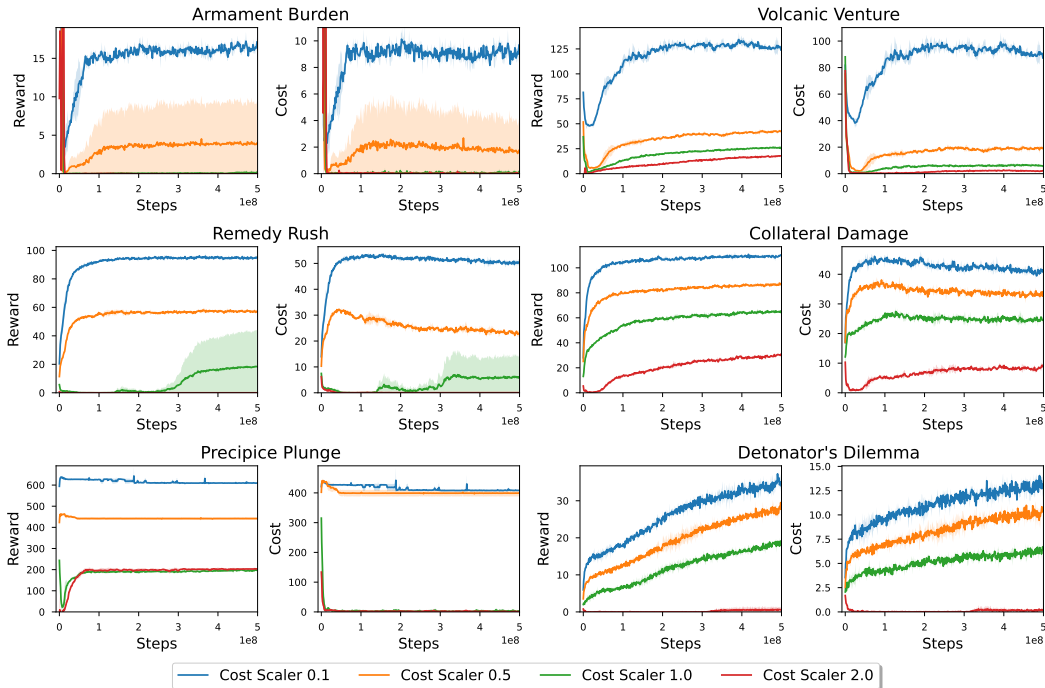


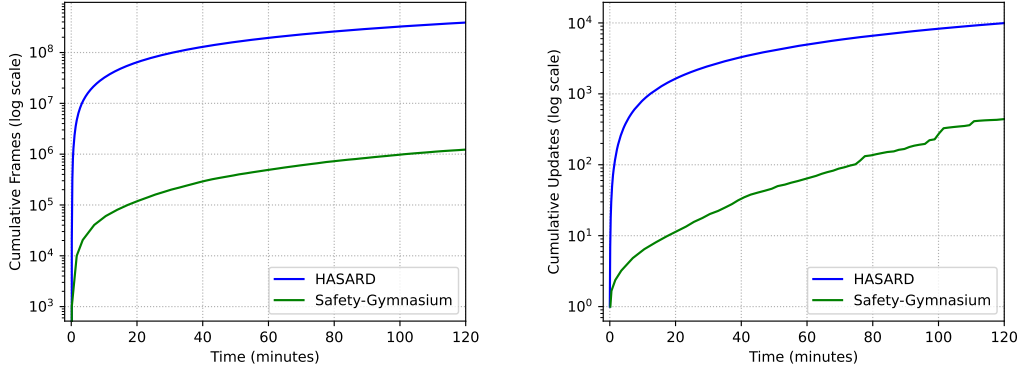
Figure 15: PPOCost treats costs as negative rewards without providing a general method for scaling these costs relative to rewards. Training with various cost scaling values indicates that PPOCost is highly sensitive to this parameter, resulting in a wide range of reward-cost trade-offs. For the main results of the paper, we adopted a scaling factor of 1.0, ensuring equal weight between rewards and negative costs. With extensive manual tuning, PPOCost can find diverse trade-offs, however, it cannot satisfy a given safety budget.

C FRAMEWORK EFFICIENCY

To demonstrate the training efficiency of HASARD, we conducted a comparative analysis with Safety-Gymnasium (Ji et al., 2023a) on the same hardware. Among existing benchmarks tailored for Safe RL research, Safety-Gymnasium is arguably the most comprehensive suite and closest to our work, as it uniquely facilitates vision-based learning within 3D environments. We arbitrarily selected the SafetyPointGoal1-v0 task and ran the PPOLag implementation with default settings in Omnisafe (Ji et al., 2023b), a popular recent Safe RL library. We allowed two hours of training time for each framework. In Figure 16 and Table 10 we compare how many environment iterations is the framework able to facilitate and how frequent are the policy updates. We can see that HASARD outperforms Safety Gymnasium in both metrics by a large magnitude.

Table 10: Efficiency Comparison of Benchmarks

Benchmark	Frames Per Second	Updates Per Second
Safety-Gymnasium	180	0.03
HASARD	53985	15.12



(a) HASARD is capable of accumulating frames at a rate several magnitudes greater than Safety-Gymnasium on the same hardware.

(b) The lack of effective parallelization results in a CPU rollout bottleneck due to which Safety-Gymnasium achieves far fewer policy updates.

Figure 16: Performance comparison of benchmarks over two hours of training on identical hardware. We trained a PPOLag agent on VolcanicVenture1-v0 from HASARD and SafetyPointGoal1-v0 from Safety-Gymnasium using default configurations.

D EXPERIMENTAL SETUP

We adopt most of our experimental setup from the PPO implementation for ViZDoom environments in Sample-Factory (Petrenko et al., 2020).

D.1 NETWORK ARCHITECTURE

The pixel observations from the environment are first processed by a CNN encoder from (Mnih et al., 2015), incorporating ELU for nonlinearity. The convolutions are passed through two dense layers, each with 512 neurons. A GRU with 512 hidden units processes the sequential and temporal information from the environment. Afterward, the architecture splits into an actor and critic network, sharing the same backbone. The actor produces a categorical distribution of action probabilities for each single action set, while the critic outputs a scalar value estimate for state-action pairs to guide policy improvements.

D.2 HYPERPARAMETERS

We present the extensive list of hyperparameters used in our experiments in Table 11.

D.3 IMPLEMENTATION DETAILS

Each action selected by the agent’s policy is repeated for four frames, reducing the need for frequent policy decisions and conserving computational resources. During preprocessing, the RGB observations are downscaled from 160×120 pixels to 128×72 pixels. We accumulate a maximum of two training batches at a time, preventing data collection from outpacing training and maintaining a balance between throughput and policy lag. We split an environment into two for a double-buffered experience collection. Each policy has a designated worker responsible for the forward pass, and we maintain a policy lag limit of 1000 steps to ensure data relevance. Parallel environment workers are dynamically set based on CPU availability. The training cycle involves collecting batches to form a dataset size defined by the product of batch size and the number of batches per epoch, with a standard batch size set at 1024. Rollouts are conducted over 32 timesteps, aligned with the recurrence interval necessary for RNN policies, facilitating efficient data processing and learning accuracy.

E EXTENDED FUTURE WORK

In this section, we will elaborate on several extensions to broaden the benchmark’s application for other settings.

Multi-Objective In the current HASARD framework, each environment only has a single safety constraint. Given the focus on Safe RL, incorporating multiple safety constraints presents a compelling avenue. In the current version only *Detonator’s Dilemma*, has two factors that increase the cost: (1) the agent harming neutral units and (2) the agent harming itself. Instead of combining them into one single value, we could decouple them into separate safety constraints. Similarly, a second constraint could be incorporated into *Armament Burden* for how frequently the agent is allowed to visit the delivery zone, and into *Volcanic Venture* for how often the agent can use the JUMP action.

Multi-Agent The ViZDoom platform supports synchronous operations among multiple agents, allowing HASARD to be extended to collaborative MARL scenarios. *Precipice Plunge* does not provide any meaningful ways of collaboration, but in other environments multiple agents could collaborate by dividing the workload of the task, focusing on separate areas of the map.

Transfer/Continual/Multi-task Learning The full action space setting is unified across all HASARD environments, allowing agents to use a consistent policy across all tasks. Moreover, there are many common elements, such as spatial navigation, entity behaviour, and environment dynamics. This paves the way for training a versatile agent capable of mastering all six tasks across three difficulty levels. This can be done in a multi-task or continual learning setting. A further interesting area of exploration is to what extent learned competencies are transferable across environments. For instance, the ability to successfully navigate an area without colliding with walls could be shared across tasks.

F SIGNIFICANCE OF HASARD

HASARD’s discrete simplified physics and pixelated graphics do not mirror the high-resolution imagery and complex dynamics required for autonomous driving or assistive robotics. Naturally, an agent excelling in HASARD would not be anywhere near capable when deployed in any of the aforementioned real-life scenarios. Nevertheless, there is substantial value in utilizing unrealistic simulation environments for foundational research in Safe RL. We list several key points motivating HASARD’s design and utility:

Computational Accessibility Incorporating highly realistic physics and visual rendering significantly increases the computational cost of simulations. HASARD is designed to be accessible for low-budget research settings, allowing more researchers to engage with vision-based Safe RL. Training RL algorithms in ultra-realistic environments with accurate physics and complex decision-making problems is not only costly but also time-intensive. By maintaining a balance between realism and

Parameter	Value	Description
Batch Size (B)	1024	Minibatch size for SGD
Gamma (γ)	0.99	Discount factor
Learning Rate (α)	1×10^{-4}	Learning rate
Hidden Layer Sizes	512, 512	Number of neurons in the dense layers after the convolutional encoder
RNN	GRU	Type of the RNN
RNN Size (h)	512	Size of the RNN hidden state
Policy Init Gain (g)	1.0	Gain parameter of neural network initialization schemas
Exploration Loss Coeff. (C_{expl})	0.001	Coefficient for the exploration component of the loss function
Value Loss Coeff. (C_{val})	0.5	Coefficient for the critic loss
Lambda Lagrange (λ_{lagr})	0.0	Lambda coefficient for the Lagrange multiplier
Lagrangian Coeff. Rate (r_{lagr})	1×10^{-2}	Change rate of the Lagrangian coefficient
KL Threshold (θ_{KL})	0.01	Threshold for the KL divergence between the old and new policy
GAE Lambda (λ_{GAE})	0.95	Generalized Advantage Estimation discounting
PPO Clip Ratio (ϵ_{clip})	0.1	PPO clipping ratio, unbiased clip version
PPO Clip Value (Δ_{clip})	1.0	Maximum absolute change in value estimate until it is clipped
Nonlinearity (ϕ)	ELU	Type of nonlinear activation function used in the network
Optimizer	Adam	Type of the optimizer
Adam Epsilon (ϵ_{Adam})	1×10^{-6}	Adam epsilon parameter
Adam Beta (β_1, β_2)	0.9, 0.999	Adam first and second momentum decay coefficient
Max Grad Norm	4.0	Max L2 norm of the gradient vector
Policy Initialization	orthogonal	Neural network layer weight initialization method
Frame Skip	4	Number of times to repeat a selected action in the environment
Frame Stack	1	Number of consecutive environment pixel-observation to stack
Env Workers	32	Number of parallel environment CPU workers
Num Envs per Worker	10	Number of environments managed by a single CPU actor
Accumulate Batches	2	Max number of training batches the learner accumulates before stopping
Worker Num Splits	2	Enable double buffered experience collection, vector environment splits
Policy Workers per Policy	1	Number of workers that compute the forward pass for each policy
Max Policy Lag	1000	Beyond how many steps to discard older experiences

Table 11: Hyperparameters

computational demand, HASARD enables a tight feedback loop that facilitates rapid experimentation and iteration.

Focus on Vision-Based Safety Much of recent Safe RL research has predominantly concentrated on continuous control problems, exemplified by the widely-used environments in Safety-Gymnasium (Ji et al., 2023a). HASARD aims to bridge the gap in safety considerations within vision-based learning, a domain with a wide range of applicability. The egocentric embodied perception in HASARD environments further introduces the problem of partial observability, which is often absent in prior works. The ultimate goal for many applications is to enable complex decision-making and precise

manipulation under realistic physics, leveraging multimodal inputs of visual imagery and LIDAR. However, fulfilling all these criteria simultaneously poses a significant challenge for Safe RL research. To manage this, HASARD narrows the focus and decouples the problem, concentrating on visual perception of 3D environments, spatial awareness, and high-level control.

Precedent in Research Video games, despite their lack of realism, have been instrumental in advancing AI research (Hu et al., 2024). They provide a controlled yet challenging environment for developing and testing new methodologies. Many game-based 3D benchmarks are adopted as useful tools within the community (Chevalier-Boisvert et al., 2024; Gong et al., 2023; Jeon et al., 2023), supporting the development of algorithms that can later be adapted to more realistic applications. Similarly, simulation environments in Doom remain viable and relevant, as they have been widely used in recent research (Park et al., 2024; Kim et al., 2023; Zhai et al., 2023; Valevski et al., 2024).

Solving Toy Problems Simulation environments with simplified physics and visuals can roughly emulate critical aspects of more complex systems. For example, much recent Safe RL research utilizes gridworld environments (Wachi et al., 2024a; Den Hengst et al., 2022), which, while basic, allow for the exploration of key concepts and strategies in safety and exploration. Similarly, HASARD, with its lower-resolution visuals and basic physics, still captures essential elements of navigating in a 3D space and introduces significant challenges related to egocentric Safe RL, such as depth perception, short-term prediction, and memory.

G TRAINING CURVES

Figure 18 depicts the training curves for levels 2 and 3, while Figure 17 shows the corresponding curves for the hard constraint setting of level 1.

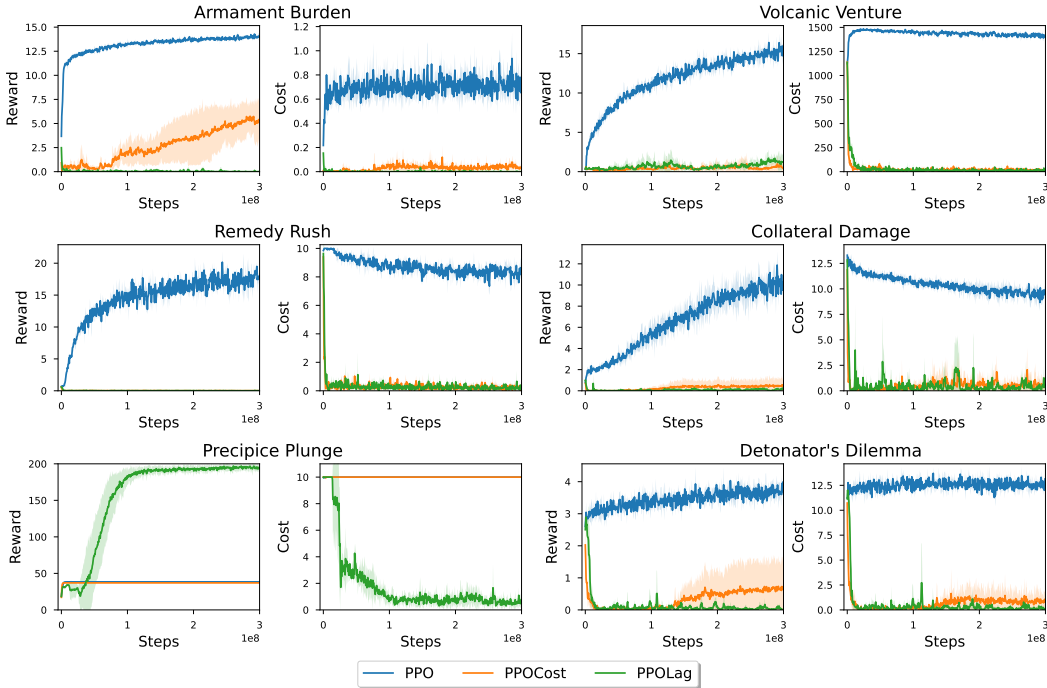
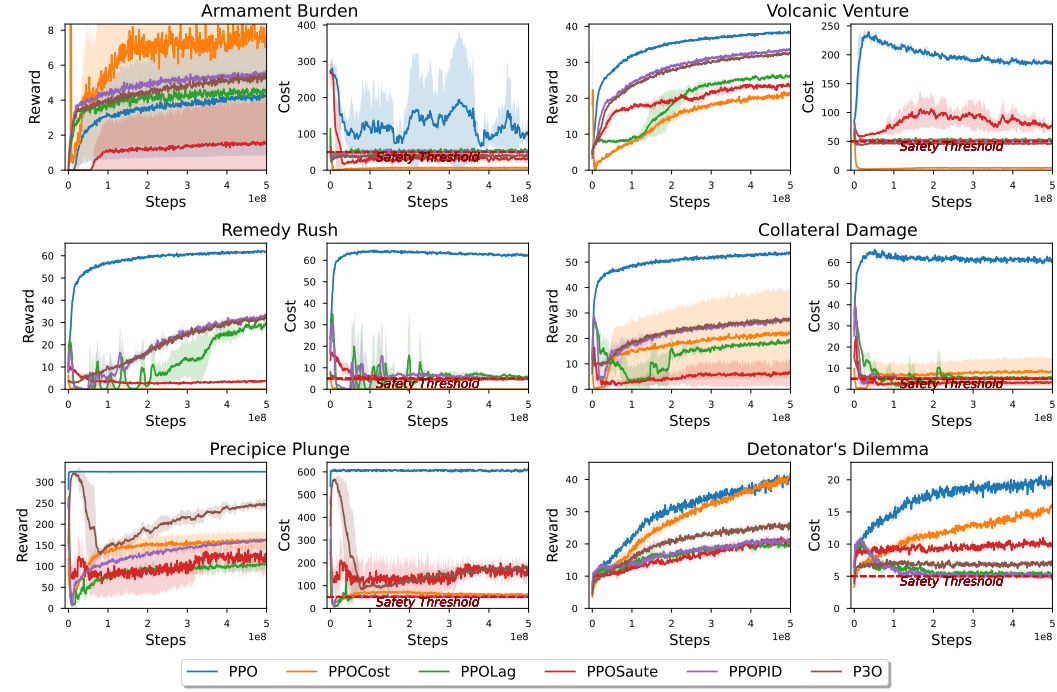
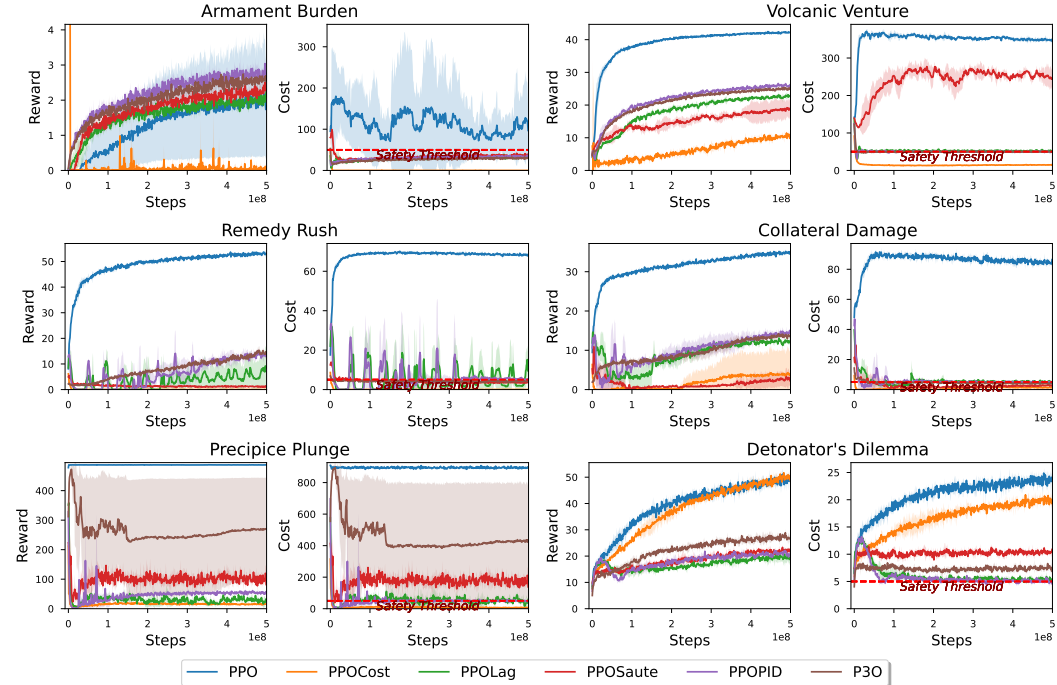


Figure 17: Baseline results of PPO variants employing hard constraints, with 95% confidence intervals across five seeds on the HASARD benchmark. **PPO** maintains high returns due to a lack of an explicit safety feedback mechanism. In contrast, **PPOCost** and **PPOLag** exhibit overly conservative behavior by consistently selecting the passive NO-OP action, failing to learn a policy that achieves noticeable rewards while strictly adhering to safety constraints. Note that the cost threshold is $\xi = 0$ under hard constraints.



(a) **Level 2.** In most scenarios, PPO sets the upper bounds for reward and cost. PPOLag is the only baseline that consistently maintains the accumulated cost near the safety budget. Conversely, PPOSaute struggles to lower its cost and stay within the safety budget in several environments.



(b) **Level 3.** PPO consistently dominates with high rewards accompanied by high costs. PPOLag and PPOSaute are unable to improve rewards while effectively controlling the cost in Remedy Rush, Collateral Damage, and Precipice Plunge, as the curves remain relatively flat. Remedy Rush demonstrates the fluctuating behavior of PPOLag originating from the Lagrangian optimization. Notably, PPOCost manages to obtain equal rewards to PPO in Detonator's Dilemma while maintaining a lower cost.

Figure 18: Training curves of higher HASARD levels.

H COMPARISON WITH SAFETY-GYMNASIUM

Safety-Gymnasium is the most comprehensive and closely related benchmark to our work, making it an essential point of comparison. Although it consists of 28 environments, there are strong similarities between them. For instance, the six environments in the **Safe Velocity** suite all share the identical objective of the agent moving forward and are visually identical. These could effectively be considered a single environment with variations in the robot type.

Similarly, the four tasks in the **Safe Navigation** suite featuring different robots like *Point*, *Car*, *RaceCar*, *Doggo*, *Ant* share a single navigation objective and only slightly vary the objects within the environment. We argue that the variations within a single environment in HASARD offer equal or greater diversity than the differences observed among the *Goal*, *Button*, *Push*, *Circle* tasks in Safety-Gymnasium.

The **Safe Vision** suite does introduce more visually diverse settings with environments such as *Building*, *Race*, *FormulaOne*. However, the core tasks remain focused on navigation, with higher difficulty levels merely introducing additional obstacles. Furthermore, these environments lack dynamic elements. The few other entities in the *Building* task move in predictable, fixed patterns, and the environment itself remains static throughout the episode. The only changes observed are those directly caused by the agent’s actions.

In contrast, **HASARD** extends beyond mere navigation-based tasks, necessitating higher-order reasoning for task resolution. It incorporates randomly moving units, and exploits the third dimension more effectively, enabling entities to navigate vertical surfaces, resulting in a richer and more complex dynamic. Leveraging the ViZDoom game engine, HASARD allows for rapid environment simulation. The six environments of HASARD, along with difficulty levels, offer a broad and dynamic challenge that provides a comprehensive evaluation of many agent competencies, such as memory, short-term prediction and distance perception.