

TRAINING VERIFIABLY ROBUST AGENTS USING SET-BASED REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep reinforcement learning uses neural networks to solve complex control tasks. However, neural networks are sensitive to input perturbations, which makes their deployment in safety-critical environments challenging and thus their formal verification necessary. This work lifts recent results from formal verification of neural networks to reinforcement learning in continuous state and action spaces. While previous work mainly focuses on adversarial attacks for robust reinforcement learning, we augment reinforcement learning with set-based computing: We enclose all possible outputs for a set of perturbed inputs and compute a gradient set for training, i.e., each possible output has a different gradient. Thereby, we control the size of the propagated sets, yielding favorable worst-case bounds for actions and value functions that enable formal verification across different verification frameworks for up to 9 times larger input perturbations. Our work addresses the gap between state-of-the-art adversarial training methods and formal verification to train verifiably robust agents, making them applicable in safety-critical environments.

1 INTRODUCTION

In recent years, deep reinforcement learning using neural networks has significantly improved solving complex control tasks (Mnih et al., 2015; OpenAI et al., 2020; Lillicrap et al., 2016). In many control tasks, state-action spaces are continuous, high-dimensional, and influenced by inherent system uncertainties, modeling errors, and sensor noise (Kober et al., 2013). However, such uncertainties are challenging for reinforcement learning when parameterizing policies as neural networks, which are sensitive to small input perturbations (Szegedy et al., 2014). This may lead to instabilities and safety violations of the controlled system: Fig. 1(a) shows an example of a navigation task where small input perturbations lead to trajectories that enter an unsafe set. Robustness to such perturbations is therefore essential. For deployment in safety-critical applications, it is additionally necessary to demonstrate safety and guaranteed worst-case performance. This necessitates formal verifiability, which is not inherently guaranteed by current state-of-the-art robust reinforcement learning methods.

Related works on robust reinforcement learning (Zhang et al., 2021a; Huang et al., 2017; Deshpande et al., 2021; Mandlkar et al., 2017; Lütjens et al., 2020; Zhang et al., 2021b) propose a competitive framework with an adversary (Moos et al., 2022; Pinto et al., 2017): In observation-robust algorithms, the adversary exploits the sensitivity of neural networks to choose a worst-case observation for the policy (Moos et al., 2022). Computing the worst-case observation is often intractable (Madry et al., 2018). Consequently, a variety of naive, gradient-based (Pattanaik et al., 2018; Mandlkar et al., 2017;

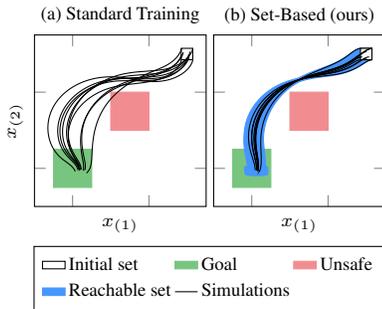


Figure 1: Comparison of standard and our set-based reinforcement learning on a navigation task. (a) Trajectories of the standard agent intersect with the obstacle. (b) We can formally verify our robust agent.¹

¹Code: <https://anonymous.4open.science/r/RobustSetBasedRL>
Video: <https://t1p.de/zjayi>

Huang et al., 2017), learning-based (Zhang et al., 2021a), and convex relaxation methods (Zhang et al., 2020), have been proposed to approximate adversarial observations. For instance, gradient-based methods typically employ the Fast Gradient Sign Method (FGSM) to approximate the most adversarial input (Goodfellow et al., 2015). In contrast, our method does not rely on a single adversarial instance. Instead, it updates the agent parameters using the entire set of admissible perturbations, thereby improving robustness and enabling stronger formal guarantees of verifiability.

Safe deployment in safety-critical environments requires more than empirical robustness, it demands formal guarantees. While prior work often relies on adversarial attacks or randomized smoothing (Wu et al., 2022) to determine a probabilistic upper bound on the worst-case performance, these approaches do not provide a provable lower bound. In contrast, our work focuses on formal verification of neural network-based control policies. Recent advances by Manzanas Lopez et al. (2023); Brix et al. (2023) make it possible to verify entire neural network control systems: This is often achieved by (i) modeling the disturbed state of the system as a continuous set, (ii) computing the corresponding output set of the neural networks, and (iii) enclosing the dynamics of the environment over time using reachability analysis. If the obtained reachable set does not violate specifications, the neural network control system is verified as shown in Fig. 1(b).

For the training of robust neural networks, a recent work proposes set-based training: For a set of possible inputs, the set of possible outputs is enclosed and the neural network is trained with a gradient set (Koller et al., 2025), i.e., each possible output has a different gradient. By picking gradients that point toward the center of the output set, the size of the output sets can be controlled. Thereby, the trained neural network is more robust and easier to formally verify with set-based verification algorithms, because smaller propagated sets reduce the conservatism of the verification algorithm. In this work, we lift set-based training to reinforcement learning. Our main contributions are:

- (i) A novel set-based reinforcement learning algorithm that, for the first time, computes a gradient set, which contains a different gradient for each possible output given input perturbations. Our algorithm trains agents that are provably more robust and formally verifiable across different available reachability-analysis frameworks.
- (ii) A rigorous analysis of the underlying set propagation to derive a novel set-based loss function for regression tasks, which is used to compute gradient sets that optimize over entire output sets and thus achieve formal verifiability of the trained agents given input perturbations.
- (iii) An extensive evaluation including a comparison with state-of-the-art adversarial training algorithms and an ablation study to justify our design choices.

2 PRELIMINARIES

2.1 NOTATION

We write vectors as lowercase letters, matrices as uppercase letters, sets as calligraphic letters, and probability distributions as script font letters. The i -th entry of $v \in \mathbb{R}^n$ is written as $v_{(i)}$. The entry in the i -th row and j -th column of a matrix is $M_{(i,j)}$; $M_{(i,\cdot)}$ is the i -th row, and $M_{(\cdot,j)}$ the j -th column. The horizontal concatenation of matrices $A \in \mathbb{R}^{n \times m}$ and $B \in \mathbb{R}^{n \times p}$ is denoted by $[A \ B]$. The identity matrix is denoted by $I_n \in \mathbb{R}^{n \times n}$, and the vector containing only ones or zeros is denoted by $\mathbf{1}$ or $\mathbf{0}$. The set of natural numbers up to $n \in \mathbb{N}$ is written as $[n] = \{1, 2, \dots, n\} \subset \mathbb{N}$. We denote a multidimensional interval by $\mathcal{I} = [l, u] = \{x \in \mathbb{R}^n \mid \forall i \in [n]: l_{(i)} \leq x_{(i)} \leq u_{(i)}\}$. The gradient of a function f w.r.t. a variable x is denoted by $\nabla_x f(x, \cdot)$. The operator $\text{diag} : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$ returns a diagonal matrix with the vector elements on its diagonal. The expected value of a random variable x under condition $y \sim \mathcal{Y}$ is $\mathbb{E}_{y \sim \mathcal{Y}}[x(y)]$.

2.2 NEURAL NETWORKS

A feed-forward neural network $N_\theta : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_\kappa}$ with learnable parameters θ consists of $\kappa \in \mathbb{N}$ alternating linear and activation layers, where the k -th layer has $n_k \in \mathbb{N}$ output neurons. The output $\hat{y} = N_\theta(x) \in \mathbb{R}^{n_\kappa}$ is computed by propagating an input $x \in \mathbb{R}^{n_0}$ through all layers.

Definition 2.1 (Neural Network, (Bishop & Nasrabadi, 2006, Sec. 5.1)). *Given a neural network N_θ and an input $x \in \mathbb{R}^{n_0}$, the output $\hat{y} = N_\theta(x) \in \mathbb{R}^{n_\kappa}$ is given by*

$$h_0 = x, \quad h_k = L_k(h_{k-1}) = \begin{cases} W_k h_{k-1} + b_k & \text{if } k\text{-th layer is linear,} \\ \sigma_k(h_{k-1}) & \text{otherwise,} \end{cases} \quad \text{for } k \in [\kappa], \quad \hat{y} = h_\kappa,$$

with weights $W_k \in \mathbb{R}^{n_k \times n_{k-1}}$, biases $b_k \in \mathbb{R}^{n_k}$, and elementwise activation functions $\sigma_k(\cdot)$.

2.3 DEEP DETERMINISTIC POLICY GRADIENT

We focus on continuous control tasks with a multidimensional state and action space \mathcal{S}, \mathcal{A} (Januszewski et al., 2021; Recht, 2019). Our set-based reinforcement learning approach is based on the deep deterministic policy gradient algorithm (DDPG) (Lillicrap et al., 2016) which consists of an actor $\mu_\phi: \mathcal{S} \rightarrow \mathcal{A}$ with parameters ϕ and a critic $Q_\theta: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ with parameters θ . Starting from an initial state s_0 , the actor observes the state $s_t \in \mathcal{S}$ at time step t and returns an action $a_t = \mu_\phi(s_t)$, which controls the system until the next time step $t + 1$. Using a reward function $r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and the state transition probabilities $p(s_{t+1}|s_t, a_t)$, the environment returns a reward r_t and its next state s_{t+1} (Fig. 2). These transitions (s_t, a_t, r_t, s_{t+1}) are stored in a buffer \mathcal{B} (Fig. 2). The training objective is to find the policy that maximizes the discounted cumulative reward (Silver et al., 2014, Eq. 8):

$$\max_{\phi} J(\phi, \rho^\mu) = \max_{\phi} \mathbb{E}_{s_t \sim \rho^\mu} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, \mu_\phi(s_t)) \right], \quad (1)$$

with discount factor $\gamma \in [0, 1]$ and discounted state visitation distribution ρ^μ for policy μ (Lillicrap et al., 2016, Sec. 2). The critic Q_θ approximates the expected total discounted reward for action a_t in state s_t (Lillicrap et al., 2016, Eq. 3):

$$Q_\theta(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho^\mu} [Q_\theta(s_{t+1}, \mu_\phi(s_{t+1}))]. \quad (2)$$

The critic is trained off-policy with a stochastic policy β (Lillicrap et al., 2016, Eq. 4) and targets y_t (Lillicrap et al., 2016, Eq. 5) using the following bootstrapped Q-iterations:

$$L(\theta) = \mathbb{E}_{s_t \sim \rho^\beta, a_t \sim \beta} [(Q_\theta(s_t, a_t) - y_t)^2], \quad y_t = r(s_t, a_t) + \gamma Q_\theta(s_{t+1}, \mu_\phi(s_{t+1})). \quad (3)$$

The actor is trained using the policy gradient (Lillicrap et al., 2016, Eq. 6):

$$\nabla_{\phi} J(\phi, \rho^\beta) \approx \mathbb{E}_{s_t \sim \rho^\beta} \left[\nabla_{a_t} Q_\theta(s_t, a_t) \Big|_{a_t = \mu_\phi(s_t)} \nabla_{\phi} \mu_\phi(s_t) \right]. \quad (4)$$

2.4 SET-BASED COMPUTATIONS

We model uncertainties using zonotopes, a continuous set representation, due to their favorable computational complexity, propagating the sets of inputs and gradients through the neural networks:

Definition 2.2 (Zonotope (Girard, 2005)). *Given center $c \in \mathbb{R}^n$ and generators $G \in \mathbb{R}^{n \times q}$, we define*

$$\mathcal{Z} = \langle c, G \rangle_{\mathcal{Z}} = \{c + G\beta \mid \beta \in [-1, 1]^q\} \subset \mathbb{R}^n.$$

The affine map $x \mapsto Ax + b$ of a zonotope $\mathcal{Z} = \langle c, G \rangle_{\mathcal{Z}}$ is computed by (Althoff, 2010, Sec. 2.4)

$$A\mathcal{Z} + b = \{Ax + b \mid x \in \mathcal{Z}\} = \langle Ac + b, AG \rangle_{\mathcal{Z}}. \quad (5)$$

The enclosing interval of zonotope $\mathcal{Z} = \langle c, G \rangle_{\mathcal{Z}}$, i.e., $\mathcal{Z} \subseteq [l, u]$ and its diameter $\text{dia}(\mathcal{Z})$, are computed by (Althoff, 2010, Prop. 2.2)

$$l = c - |G| \mathbf{1}, \quad u = c + |G| \mathbf{1}, \quad \text{dia}(\mathcal{Z}) := u - l = 2|G| \mathbf{1}. \quad (6)$$

For some derivations, we write $\ln \text{Dia}(G) := \ln(2|G| \mathbf{1})$ and $\ln \text{Dia}'(G) := \nabla_G \ln \text{Dia}(G) = \text{diag}(|G| \mathbf{1})^{-1} \text{sign } G$ to avoid clutter, $\text{sign}(G)$ returns the sign of each entry of matrix G .

The Minkowski sum of a zonotope $\mathcal{Z} = \langle c, G \rangle_{\mathcal{Z}}$ and an interval $\mathcal{I} = [l, u]$ is computed by (Althoff, 2010, Prop. 2.1 and Sec. 2.4):

$$\mathcal{Z} \oplus \mathcal{I} = \{x_1 + x_2 \mid x_1 \in \mathcal{Z}, x_2 \in \mathcal{I}\} = \langle c + 1/2(u + l), [G \text{ dia}(1/2(u - l))] \rangle_{\mathcal{Z}}. \quad (7)$$

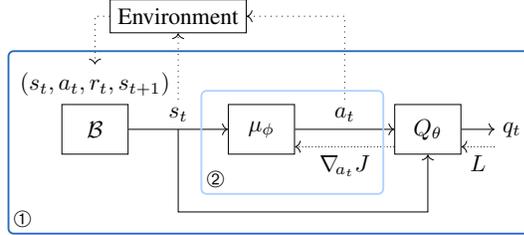


Figure 2: Illustration of the structure of the deep deterministic policy gradient algorithm; ① and ② show the components that are augmented by our set-based training (introduced in Sec. 3).

2.5 SET PROPAGATION THROUGH NEURAL NETWORKS

Computing the exact output set $\mathcal{Y}^* = N_\theta(\mathcal{X})$ of a neural network for a given input set $\mathcal{X} \subset \mathbb{R}^{n_0}$ is computationally hard, i.e., for neural networks with ReLU-activation it is NP-hard (Katz et al., 2017). Thus, an enclosure of the output set $\mathcal{Y} \supseteq \mathcal{Y}^*$ is computed by conservatively propagating the input set through the layers of the neural network:

Proposition 2.3 (Neural Network Set Propagation (Singh et al., 2018)). *Given an input set \mathcal{X} , the output set of a neural network can be enclosed as:*

$$\mathcal{H}_0 = \mathcal{X}, \quad \mathcal{H}_k = \text{enclose}(L_k, \mathcal{H}_{k-1}) \quad \text{for } k \in [\kappa], \quad \mathcal{Y}^* \subseteq \mathcal{Y} = \mathcal{H}_\kappa.$$

The operation $\text{enclose}(L_k, \mathcal{H}_{k-1})$ encloses the output set of the k -th layer given the input set \mathcal{H}_{k-1} . If the k -th layer is linear (Def. 2.1), the affine map (5) is applied:

$$\text{enclose}(L_k, \mathcal{H}_{k-1}) = W_k \mathcal{H}_{k-1} + b_k. \quad (8)$$

The output set of an activation layer is enclosed by approximating its activation function element-wise with a linear function with slope $m_k \in \mathbb{R}^{n_k}$ and approximation error $[d_k, \bar{d}_k]$ Koller et al. (2025):

$$\text{enclose}(L_k, \mathcal{H}_{k-1}) = \text{diag}(m_k) \mathcal{H}_{k-1} \oplus [d_k, \bar{d}_k], \quad m_k = \frac{\sigma_k(u_{k-1}) - \sigma_k(l_{k-1})}{u_{k-1} - l_{k-1}}. \quad (9)$$

Intuitively, we can thereby compute the output of the neural network for an infinite number of inputs in a single forward pass, thereby enabling us to adapt learning to account for worst-case outcomes.

2.6 PROBLEM STATEMENT

We are given a reinforcement learning task with uncertain initial states $s_0 \in S_0 \subset \mathbb{R}^n$. Moreover, we model state uncertainties with an ℓ_∞ -ball of radius $\epsilon \in \mathbb{R}_+$. The set of all perturbations is $\mathcal{V}_\epsilon^\infty := \{\nu: \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}^n \mid \forall s_t \in \mathbb{R}^n, \forall t \in \mathbb{R}_+: \|\nu(s_t, t) - s_t\|_\infty \leq \epsilon\}$. Our goal is to use set-based training (Koller et al., 2025) to robustly maximize the reward under adversarial perturbations

$$\max_{\phi} \min_{\nu \in \mathcal{V}_\epsilon^\infty} J(\phi, \rho^{\mu^{\circ\nu}}), \quad (10)$$

where $\rho^{\mu^{\circ\nu}}$ denotes the state visitation distribution perturbed with adversary $\nu \in \mathcal{V}_\epsilon^\infty$. Further, let \mathcal{R} enclose all reachable states for the time interval $[0, t_{\text{end}}]$. We derive a verified performance as a lower bound on the reward by evaluating the worst trajectory from \mathcal{R} .

3 SET-BASED REINFORCEMENT LEARNING

We augment actor-critic reinforcement learning with set-based training using a gradient set, i.e., each possible output has a different gradient. Therefore, we propagate sets through the actor and the critic and train both, actor and critic, with gradient sets (SA-SC), which corresponds to a set-based evaluation of ① in Fig. 2. To this end, we extend the critic loss (3), policy gradient (4), and the replay buffer to sets. The main steps are given in Alg. 1.

Given s_t , we enclose all possible adversaries:

$$\mathcal{S}_t = \{\nu(s_t, t) \mid \nu \in \mathcal{V}_\epsilon^\infty\} = \langle s_t, \epsilon I \rangle_Z. \quad (11)$$

We enclose the set of possible actions by propagating \mathcal{S}_t through the actor (Prop. 2.3):

$$\mathcal{A}_t = \langle c_{\mathcal{A}_t}, G_{\mathcal{A}_t} \rangle_Z = \text{enclose}(\mu_\phi, \mathcal{S}_t). \quad (12)$$

For the off-policy training of the critic, we perturb the set of actions with random exploration noise e_t (Lillicrap et al., 2016, Eq. 7):

$\tilde{\mathcal{A}}_t := \mathcal{A}_t + e_t$, and compute the corresponding set of critic outputs:

$$\mathcal{Q}_t = \langle c_{\mathcal{Q}_t}, G_{\mathcal{Q}_t} \rangle_Z = \text{enclose}(Q_\theta, \mathcal{S}_t \times \tilde{\mathcal{A}}_t), \quad (13)$$

Algorithm 1 Set-based reinforcement learning.

- 1: Randomly initialize Q_θ, μ_ϕ with θ, ϕ
 - 2: Initialize replay buffer \mathcal{B}
 - 3: **for** episode = 1, . . . , maxEpisodes **do**
 - 4: Get initial observation s_0
 - 5: **for** $t = 0, 1, \dots, \text{maxSteps}$ **do**
 - 6: $\mathcal{S}_t \leftarrow \langle s_t, \epsilon I \rangle_Z$ {perturb state Eq. (11)}
 - 7: $\mathcal{A}_t \leftarrow \mu_\phi(\mathcal{S}_t)$ {evaluate actor Eq. (12)}
 - 8: $\tilde{\mathcal{A}}_t \leftarrow \mathcal{A}_t + e_t, e_t \sim \mathbb{P}(E)$ {Eq. (13)}
 - 9: $r_t \leftarrow r(s_t, c_{\tilde{\mathcal{A}}_t})$
 - 10: $s_{t+1} \leftarrow \text{execute action } c_{\tilde{\mathcal{A}}_t}$
 - 11: Store transition $(s_t, \tilde{\mathcal{A}}_t, r_t, s_{t+1})$ in \mathcal{B}
 - 12: Sample batch of n transitions from \mathcal{B}
 - 13: Compute targets y_i using Eq. (3)
 - 14: Update critic using Prop. 3.1
 - 15: Update actor using Def. 3.3
 - 16: **end for**
 - 17: **end for**
-

where the Cartesian product (\times) respects the dependencies between two zonotopes (Lützw & Althoff, 2023, Sec. II.C). The environment receives the perturbed center of the action set $c_{\tilde{\mathcal{A}}_t} := c_{\mathcal{A}_t} + e_t$ and returns the reward $r(s_t, c_{\tilde{\mathcal{A}}_t})$ and next state s_{t+1} , which are stored in the replay buffer \mathcal{B} as transition $(s_t, \tilde{\mathcal{A}}_t, r_t, s_{t+1})$. To obtain temporally uncorrelated samples, we randomly sample n transitions from the buffer for the training of the critic neural network. For each transition $i \in [n]$, we compute the target y_i using (3) and compute gradient sets using a set-based loss function:

Proposition 3.1 (Set-Based Regression Loss). *Given output set $\mathcal{Q}_i = \langle c_{\mathcal{Q}_i}, G_{\mathcal{Q}_i} \rangle_Z \subset \mathbb{R}$ and target $y_i \in \mathbb{R}$, the set-based regression loss is defined as*

$$L_{Reg}(y_i, \mathcal{Q}_i) = \underbrace{1/2(c_{\mathcal{Q}_i} - y_i)^2}_{\text{standard training loss}} + \underbrace{\eta_Q/\epsilon \ln\text{Dia}(G_{\mathcal{Q}_i})}_{\text{robustness}},$$

with weighting factor $\eta_Q \in \mathbb{R}_+$ and perturbation radius $\epsilon \in \mathbb{R}_+$. The gradient of L_{Reg} w.r.t. \mathcal{Q}_i is:

$$\nabla_{\mathcal{Q}_i} L_{Reg}(y_i, \mathcal{Q}_i) = \langle c - y_i, \eta_Q/\epsilon \ln\text{Dia}'(G_{\mathcal{Q}_i}) \rangle_Z.$$

Proof. See Appendix B. □

The set-based regression loss combines the half-squared error (Bishop & Nasrabadi, 2006, Eq. 5.14) of the center with a robustness loss to reduce the size of the output set. The actor neural network is trained using a set-based policy gradient which contains a different gradient for each possible output.

Definition 3.2 (Set-Based Policy Gradient SA-SC). *Given states \mathcal{S}_i with the corresponding actions $\mathcal{A}_i = \langle c_{\mathcal{A}_i}, G_{\mathcal{A}_i} \rangle_Z$ and critic outputs $\mathcal{Q}_i = \langle c_{\mathcal{Q}_i}, G_{\mathcal{Q}_i} \rangle_Z$, a set-based policy gradient is defined as*

$$\nabla_{\mathcal{A}_i} J_{Set}(\mu_\phi) := \left\langle \underbrace{\nabla_{c_{\mathcal{A}_i}} J_{Set}(\mu_\phi)}_{\text{policy gradient}}, \underbrace{\nabla_{G_{\mathcal{A}_i}} J_{Set}(\mu_\phi)}_{\text{robustness}} \right\rangle_Z,$$

where $\nabla_{c_{\mathcal{A}_i}} J_{Set}(\mu_\phi) = \mathbb{E}_{s_i \sim \rho^\beta} [\nabla_{c_{\mathcal{A}_i}} c_{\mathcal{Q}_i}]$,

and $\nabla_{G_{\mathcal{A}_i}} J_{Set}(\mu_\phi) = -\eta_\mu/\epsilon \mathbb{E}_{s_i \sim \rho^\beta} [\omega \ln\text{Dia}'(G_{\mathcal{A}_i}) + (1 - \omega) \nabla_{G_{\mathcal{A}_i}} \ln\text{Dia}'(G_{\mathcal{Q}_i})]$,

with weights $\eta_\mu \in \mathbb{R}_+$, $\omega \in [0, 1]$ and perturbation $\epsilon \in \mathbb{R}_+$.

The set-based policy gradient is a zonotope consisting of a center, which corresponds to the standard policy gradient (4), and a generator matrix. The generator matrix $\nabla_{G_{\mathcal{A}_i}} J_{Set}(\mu_\phi)$ uses the factor ω to connect a robustness loss for the action set \mathcal{A}_i and the gradients of the robustness loss of the critic outputs \mathcal{Q}_i in the action space (Fig. 2). Similar to (Zhang et al., 2021a, Sec. 3.3), the set-based policy gradient reduces the size of the action set to mitigate performance loss caused by an adversary. Finally, given the gradients w.r.t. the output of the actor and the critic, we can update the respective parameters using set-based backpropagation (Koller et al., 2025, Sec. IV-B.). We can omit the set propagation through the critic with $\omega = 1$. In this case, only the robustness loss of the action set is used (SA-PC; set-evaluation of ② in Fig. 2). Fig. 3 visualizes the gradient set for samples of a zonotope: The gradients of the robustness loss point toward the center and thereby enforce smaller output sets.

Definition 3.3 (Set-Based Policy Gradient SA-PC). *Given states \mathcal{S}_i with the corresponding actions $\mathcal{A}_i = \langle c_{\mathcal{A}_i}, G_{\mathcal{A}_i} \rangle_Z$, a set-based policy gradient is defined as*

$$\nabla_{\mathcal{A}_i} J_{Set}(\mu_\phi) := \left\langle \underbrace{\nabla_{c_{\mathcal{A}_i}} J_{Set}(\mu_\phi)}_{\text{policy gradient}}, \underbrace{\nabla_{G_{\mathcal{A}_i}} J_{Set}(\mu_\phi)}_{\text{robustness}} \right\rangle_Z,$$

where $\nabla_{c_{\mathcal{A}_i}} J_{Set}(\mu_\phi) = \mathbb{E}_{s_i \sim \rho^\beta} [\nabla_{c_{\mathcal{A}_i}} Q_\theta(s_i, c_{\mathcal{A}_i})]$, $\nabla_{G_{\mathcal{A}_i}} J_{Set}(\mu_\phi) = -\eta_\mu/\epsilon \mathbb{E}_{s_i \sim \rho^\beta} [\ln\text{Dia}'(G_{\mathcal{A}_i})]$,

with weight $\eta_\mu \in \mathbb{R}_+$ and perturbation $\epsilon \in \mathbb{R}_+$.

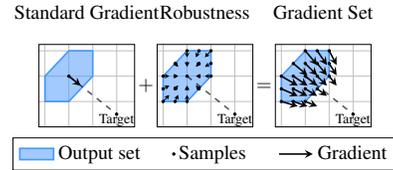


Figure 3: Visualization of the policy gradient set (Def. 3.3): each point in the output set has a different gradient.

4 DERIVATION OF SET-BASED LOSS FUNCTIONS

In this section, we motivate our choices for the set-based loss function (Prop. 3.1) and our set-based policy gradients (Def. 3.2 and 3.3) that are used to compute the gradient sets. While maximizing the likelihood of outputs is a standard procedure to derive loss functions (Bishop & Nasrabadi, 2006, Sec. 1.2.5), lifting this theory to set-based computing has the unique advantage of integrating formal methods into the training process. To this end, we connect set-based computing and probability theory by assuming a probability distribution over the considered closed sets. For our derivations, we use a conditional posterior distribution which can be rewritten to be proportional to a likelihood function and a prior distribution (Bishop & Nasrabadi, 2006, Eq. 1.44):

$$\text{cond. posterior} \propto \text{likelihood} \cdot \text{prior}. \quad (14)$$

This reformulation is used as the posterior and is not directly obtainable, but we can assume distributions for the likelihood and the prior to obtain an estimate. In our case, the likelihood corresponds to the standard (point-based) training goal, and the prior penalizes the volume of the computed sets. As zonotopes represent these sets and thus are point-symmetric, we additionally assume that the expected value over a zonotope is its center:

$$\mathbb{E}_{z \sim \mathcal{Z}}[z] = c. \quad (15)$$

4.1 SET-BASED REGRESSION LOSS

We sample random transitions $i \in [n]$ from the buffer \mathcal{B} to obtain a state s_i and the corresponding actions $\tilde{a}_i \sim \tilde{\mathcal{A}}_i$. For each transition, we use (13) to obtain the critic output \mathcal{Q}_i and use (3) to obtain the target y_i for each critic output $q_i \sim \mathcal{Q}_i$ using the rewards and next states stored in the buffer. To train Q_θ , we want to maximize the probability $p_\theta(q_i|y_i, s_i, \tilde{a}_i, \beta^{-1})$. Since this probability can not be computed directly, we model this probability as a conditional posterior using (14):

$$\underbrace{p_\theta(q_i|y_i, s_i, \tilde{a}_i, \beta^{-1})}_{\text{cond. posterior}} \propto \underbrace{p(y_i|q_i, \beta^{-1})}_{\text{likelihood}} \underbrace{p_\theta(q_i|s_i, \tilde{a}_i)}_{\text{prior}}. \quad (16)$$

For the prior, we assume that q_i is uniformly distributed over the interval $[l_{\mathcal{Q}_i}, u_{\mathcal{Q}_i}] \supseteq \mathcal{Q}_i \subset \mathbb{R}$, as Prop. 2.3 is also defined over the enclosing interval. Thus, the prior is given by

$$p_\theta(q_i|s_i, \tilde{a}_i) = \mathcal{U}(q_i|l_{\mathcal{Q}_i}, u_{\mathcal{Q}_i}) = \text{dia}(\mathcal{Q}_i)^{-1}. \quad (17)$$

As the critic learns the bootstrapped Q-iterations by regression, we assume that y_i is normally distributed with mean q_i and variance β^{-1} with the likelihood (Bishop & Nasrabadi, 2006, Eq. 1.60):

$$p(y_i|q_i, \beta^{-1}) = \mathcal{N}(y_i|q_i, \beta^{-1}) = \sqrt{\beta/2\pi} \exp(-\beta/2 (q_i - y_i)^2). \quad (18)$$

Since we observe not a single q_i but an entire set \mathcal{Q}_i , we use the expected value $\mathbb{E}_{q_i \sim \mathcal{Q}_i}[q_i] = c_{\mathcal{Q}_i}$ (15) in the likelihood function (Bishop & Nasrabadi, 2006, Sec. 10.3). Thus, we obtain the following term to be maximized:

$$p_\theta(q_i|y_i, s_i, \tilde{a}_i, \beta^{-1}) \propto p(y_i|c_{\mathcal{Q}_i}, \beta^{-1}) p_\theta(q_i|s_i, \tilde{a}_i) \quad (19)$$

Finally, we apply the negative logarithm and set $\beta = (\eta\alpha/\epsilon)^{-1}$ to obtain our set-based loss (Prop. 3.1):

$$\begin{aligned} & -\ln(p(y_i|c_{\mathcal{Q}_i}, \beta^{-1}) p_\theta(q_i|s_i, \tilde{a}_i)) \stackrel{(18), (17)}{=} -\ln(\mathcal{N}(y_i|q_i, \beta^{-1}) \text{dia}(\mathcal{Q}_i)^{-1}) \\ & \propto \beta/2 (c_{\mathcal{Q}_i} - y_i)^2 + \ln \text{Dia}(G_{\mathcal{Q}_i}) \stackrel{\text{Prop. 3.1}}{\propto} L_{\text{Reg}}(y_i, \mathcal{Q}_i). \end{aligned} \quad (20)$$

We choose β that way for easier fine-tuning (Koller et al., 2025, Def. 5).

4.2 SET-BASED POLICY GRADIENT

For a state s_i , we derive the set-based policy gradient analogous to (Xiao & Wang, 2022), to maximize the probability of an action $a_i \sim \mathcal{A}_i$ being optimal given s_i – which is again not directly obtainable. Thus, we introduce a binary variable o_i indicating whether a_i is optimal and abbreviate $o_i = 1$ by o_i (Xiao & Wang, 2022, Sec. 3.1). We again model this probability as a conditional posterior over the action output using (14):

$$\underbrace{p_{\phi, \theta}(a_i|o_i, s_i, q_i, \alpha)}_{\text{cond. posterior}} \propto \underbrace{p(o_i|q_i, \alpha)}_{\text{likelihood}} \underbrace{p_{\phi, \theta}(a_i|q_i, s_i)}_{\text{prior}}. \quad (21)$$

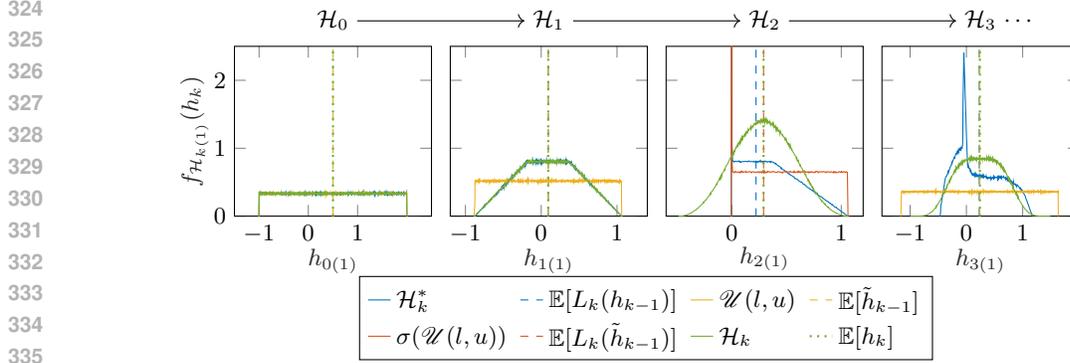


Figure 4: Propagated probability density func. $f_{\mathcal{H}_k^{(1)}}(h_k)$ of zonotopes for a ReLU-neural net: True sampled density func.(blue), interval enclosure (yellow), and density of sets from Prop. 2.3 with $\beta_j \sim \mathcal{U}(-1, 1)$ (Def. 2.2) (green).

We assume the likelihood to be exponentially distributed with $\alpha \in \mathbb{R}_+$ (Xiao & Wang, 2022, Sec. 3.2):

$$p(o_i|a_i, q_i, \alpha) = \exp(\alpha^{-1} q_i). \quad (22)$$

The prior in (21) is not directly computable. Thus, we model it again as a conditional posterior (14):

$$\underbrace{p_{\phi, \theta}(a_i|q_i, s_i)}_{\text{cond. posterior}} = \underbrace{p_{\theta}(q_i|a_i, s_i)}_{\text{likelihood}} \underbrace{p_{\phi}(a_i|s_i)}_{\text{prior}}, \quad (23)$$

where the likelihood function of q_i and the prior for a_i are uniform distributions over the enclosing intervals $[l_{\mathcal{Q}_i}, u_{\mathcal{Q}_i}] \supseteq \mathcal{Q}_i \subset \mathbb{R}$ and $[l_{\mathcal{A}_i}, u_{\mathcal{A}_i}] \supseteq \mathcal{A}_i \subset \mathbb{R}^{n_{\mathcal{A}_i}}$ analogous to (17):

$$\begin{aligned} p_{\theta}(q_i|a_i, s_i) &= \mathcal{W}(q_i|l_{\mathcal{Q}_i}, u_{\mathcal{Q}_i}) = \text{dia}(\mathcal{Q}_i)^{-1}, \\ p_{\phi}(a_i|s_i) &= \mathcal{W}(a_i|l_{\mathcal{A}_i}, u_{\mathcal{A}_i}) = \prod_{j=1}^{n_{\mathcal{A}_i}} \text{dia}(\mathcal{A}_i)_{(j)}^{-1}. \end{aligned} \quad (24)$$

Please note that these two probabilities correspond to the evaluation of the actor and the critic, respectively. For the likelihood function of (21), the expected value $\mathbb{E}_{q_i \sim \mathcal{Q}_i}[q_i] = c_{\mathcal{Q}_i}$ (15) is used, and taking the logarithm obtains us:

$$\ln(p(o_i|c_{\mathcal{Q}_i}, \alpha) p_{\phi}(a_i|s_i) p_{\theta}(q_i|a_i, s_i)) \stackrel{(22), (24)}{=} \alpha^{-1} c_{\mathcal{Q}_i} - \mathbf{1}^{\top} \ln \text{Dia}(G_{\mathcal{A}_i}) - \ln \text{Dia}(G_{\mathcal{Q}_i}). \quad (25)$$

The set-based policy gradient for SA-SC (Def. 3.2) is derived by differentiation, where we again set the weighting factor $\alpha = (\eta_{\mu}/\epsilon)^{-1}$ for easier fine-tuning (Koller et al., 2025, Def. 5). Moreover, we introduce a factor $\omega \in [0, 1]$ to weigh the gradients of the prior terms of \mathcal{A}_i and \mathcal{Q}_i .

Derivation of SA-PC. For SA-PC, only the actor is trained using set-based training while the critic uses standard (point-based) training (① in Fig. 2). Thus, we drop the prior for the critic output q_i as this is no longer evaluated set-based and use the expected value $\mathbb{E}_{a_i \sim \mathcal{A}_i}[a_i] = c_{\mathcal{A}_i}$ for the likelihood:

$$\underbrace{p(a_i|o_i, s_i, \alpha, \phi)}_{\text{cond. posterior}} \propto \underbrace{p(o_i|s_i, c_{\mathcal{A}_i}, \alpha)}_{\text{likelihood}} \underbrace{p_{\phi}(a_i|s_i)}_{\text{prior}}. \quad (26)$$

Taking the logarithm while keeping our assumption on the likelihood function and the prior leads to

$$\ln(p(o_i|s_i, c_{\mathcal{A}_i}, \alpha) p_{\phi}(a_i|s_i)) \stackrel{(22), (24)}{=} \alpha^{-1} Q_{\theta}(s_i, c_{\mathcal{A}_i}) - \mathbf{1}^{\top} \ln \text{Dia}(G_{\mathcal{A}_i}). \quad (27)$$

The set-based policy gradient for SA-PC (Def. 3.3) is derived by differentiation and choosing α as above. This also corresponds to setting $\omega = 1$ in Def. 3.2.

4.3 EXPECTATION PRESERVING IMAGE ENCLOSURE

For (19), (25) and (27), we simplify the likelihood with the expected value of the neural network output, i.e., the center. This simplification is justified with Prop. 4.1. In Fig. 4, we plot the probability

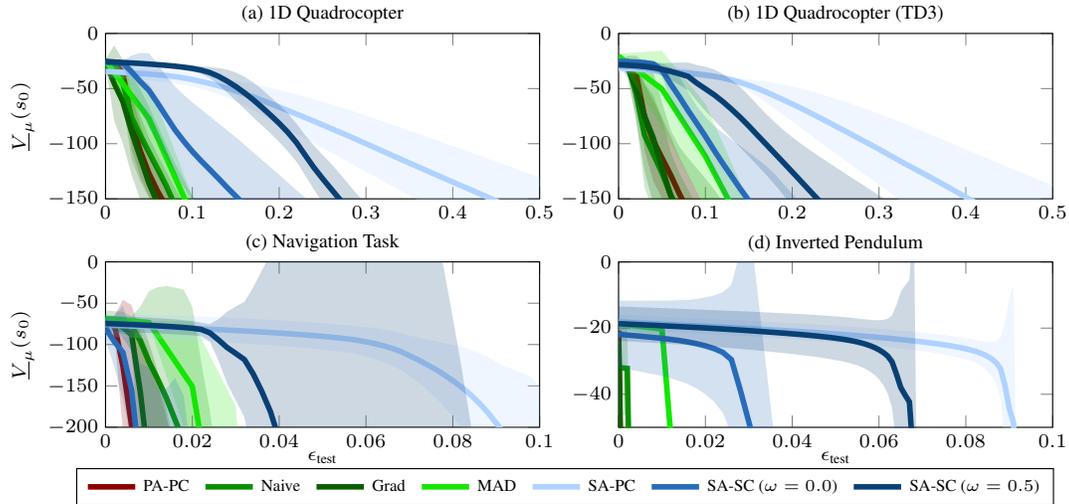


Figure 5: Comparison of verified performance $\underline{V}_\mu(s_0)$ for the (a) *1D Quadrotor*, (c) *Navigation Task*, and (d) *Inverted Pendulum*. The TD3 implementation is compared in (b) for the *1D Quadrotor*.

distributions of a set propagation and visualize the expected value for the first neuron of each layer. We compare the empirically evaluated density of a uniform input disturbance (blue) with the corresponding set-based zonotope propagation (green) across the network layers. The set propagation can be interpreted as a form of variational inference in which the output is constrained to be a zonotope. To illustrate the mean-preserving property of this single-layer variational approximation, we additionally plot the interval enclosures before activation (yellow) and the transformed intervals after activation (red). The expected value is trivially preserved for linear layers as the affine transformation for zonotopes is computed in closed form. For nonlinear layers, we observe that the expected value of the interval enclosure (red vertical line in third plot) is preserved through the enclosure (green vertical line; Prop. 2.3), with only small deviations to the expected value of the empirically evaluated density function (blue). Let us formally state this observation:

Proposition 4.1 (Tight Expectation-Preserving Set Propagation). *Given a neural network with ReLU-activations and an input set \mathcal{H}_{k-1} with the enclosing interval $[l_{k-1}, u_{k-1}] \supseteq \mathcal{H}_{k-1}$, the expected value of the enclosure of the k -th layer is*

$$\mathbb{E}_{h_k \sim \mathcal{H}_k}[h_k] = \mathbb{E}_{h_{k-1} \sim \mathcal{U}(l_{k-1}, u_{k-1})}[L_k(h_{k-1})],$$

with $\mathcal{H}_k = \text{enclose}(L_k, \mathcal{H}_{k-1})$ having minimal approximation errors.

Proof. See Appendix B. □

5 EVALUATION

We use the MATLAB toolbox CORA (Althoff, 2015) to implement our novel set-based reinforcement learning algorithm and compare the *SA-PC* and the *SA-SC* implementation against standard (point-based) training (*PA-PC*) and three state-of-the-art adversarial methods: *Naive*-, *Grad*- and *MAD* (Maximum Action Difference)-based implementations (Pattanaik et al., 2018, Alg. 2 and 4)(Zhang et al., 2020), which compute adversarial attacks to approximate the worst-case observation within a perturbation set. Similar to previous work (Yuan et al., 2022; Krasowski et al., 2023) on robust reinforcement learning, we consider four benchmarks in our evaluation: *Navigation Task*, *1D Quadrotor*, *Inverted Pendulum*, and *2D Quadrotor*. A detailed description of all benchmarks can be found in Appendix A. As in previous work (Yuan et al., 2022; Krasowski et al., 2023), we use neural networks with ReLU activations and two hidden layers of 64 and 32 neurons for the actor and critic networks. The output layer of the actor has a tanh activation. We provide the mean results and the 95% confidence interval across five different random seeds. Hyperparameters and evaluation details are given in Appendix A.

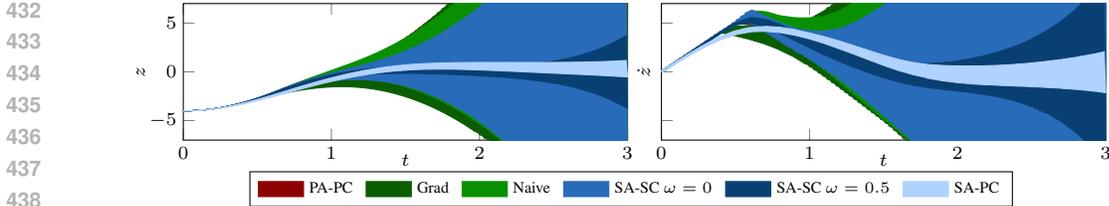


Figure 6: *Quadcopter 1D*: Comparison of reachable altitudes and vertical speeds for $\epsilon_{\text{test}} = 0.15$.

Verified performance. Following our problem statement (Sec. 2.6), we evaluate our agents based on their worst-case cumulative reward given uncertainties, which we compute using reachability analysis in CORA. Starting at an initial point s_0 , reachability analysis encloses all reachable states within a time interval $[0, t_{\text{end}}]$ with time horizon $t_{\text{end}} \in \mathbb{R}_+$ (Fig. 1b). Uncertainties (11) are added at each time step $t \in \{0, \dots, t_{\text{end}}\}$ and carried through until the time horizon is reached. Please note that, in general, this process is outer-approximative due to the continuous time and nonlinearities within the system and networks. Thus, reachability analysis obtains a formally verified lower bound of the worst-case cumulative reward, which we refer to as *verified performance* from now on. For reward functions of the form $r(s_t, a_t) = w^\top |s_t - s^*|$, we can use set-based computing to obtain the verified performance $\underline{V}_\mu(s_0)$ computed from the set of states $\mathcal{S}_t = \langle c_t, G_t \rangle_Z$ obtained by CORA:

$$\underline{V}_\mu(s_0) = \sum_{t=0}^{t_{\text{end}}} \gamma^t \max_{s_t \in \mathcal{S}_t} w^\top |s_t - s^*| \stackrel{(5),(6)}{=} \sum_{t=0}^{t_{\text{end}}} \gamma^t (w^\top |c_t - s^*| + \text{dia}(w^\top \mathcal{S}_t)/2). \quad (28)$$

Intuitively, $\underline{V}_\mu(s_0)$ is high if we can formally verify the robustness of our agent and drops otherwise. Previous works evaluated their agents solely through adversarial attacks, with learned probabilistic dynamic models (Yang et al., 2024) or randomized smoothing (Wu et al., 2022), which provide an empirical upper bound or a probabilistic bound of the worst-case reward and thus do not determine a formal lower bound.

Main results. We present the verified performance with increasing perturbation radius ϵ_{test} (11) of the considered training methods in Fig. 5. The set-based algorithms *SA-PC* and *SA-SC* train agents that can be verified for up to 9 times larger perturbation radii ϵ_{test} than the other methods. Thus, we can show robust performance of those agents, although much larger disturbances are considered. We showcase this in Fig. 6 using a large perturbation radius, where the reachable set of the *SA-PC* agent remains much smaller and thus the stability can be formally verified. As ϵ_{test} increases, the set-based training methods show a higher verified performance across all benchmarks, indicating reduced performance degradation under growing disturbance (Fig. 5).

Ablation study. Let us continue with our ablation study on various components of our algorithm:

1) *Influence of weighting factor ω (Def. 3.2)*: Recall that this parameter is used to determine where the volume of the set is penalized: With $\omega = 0$, only the output set of the critic is penalized (*SA-SC*) and with $\omega = 1$, only the output set of the actor is penalized, which corresponds to the *SA-PC* method. The *SA-PC* implementation trains more conservative actors, which perform best for large ϵ_{test} while having a worse verified performance for small ϵ_{test} . For example, the *SA-PC* actor for the *1D Quadrotor* is the most robust but reaches the goal later using lower vertical velocities (Fig. 6). By fine-tuning ω , a balance can be found where the verified performance for small ϵ_{test} can be regained while still being much more robust for larger ϵ_{test} (Fig. 5).

2) *Performance under different attacks*: It is well-known that networks trained using adversarial attacks are robust against their respective attack method but might not be robust against other methods. We observe in Fig. 7 that under *Naive* and *Grad* attacks, also their respectively trained agents outperform standard point-based training with DDPG Lillicrap et al. (2016). In particular, *SA-SC* with $\omega = 0$ and $\omega = 0.5$, perform on-par with these attacks, while *SA-PC* trained actors are more conservative and under-perform under these attacks. This additionally motivates the modularity of *SA-SC*, which is investigated in ablation 1). However, verified performance captures all attacks and, thus, *Grad*, *Naive*, and *MAD* methods perform similarly to standard point-based training (*PA-PC*) across all benchmarks (Fig. 5), while *SA-SC* and *SA-PC* demonstrate their robustness against the entire set of all possible perturbations.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

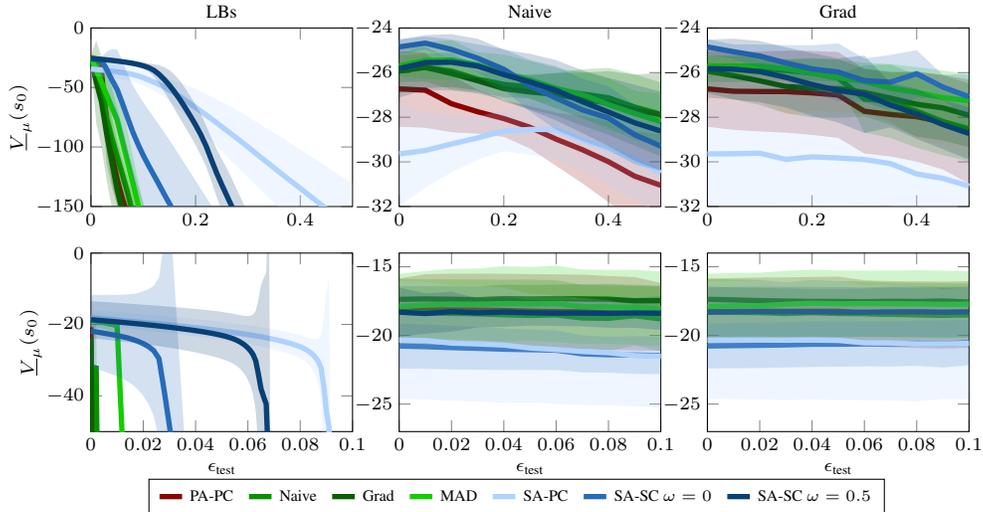


Figure 7: (Top) *Quadcopter 1D* and (bottom) *Inverted Pendulum*: Comparison of the return under *Naive* and *Grad* attacks with $\underline{V}_\mu(s_0)$.

Table 1: Time [s] for verification of Navigation task agents. Entries with – indicate that the verification toolbox is not able to verify the system.

Toolbox	CORA	CROWN-Reach	JuliaReach	NVV
Standard	423.45	130.12	–	–
Set-Based	1.99	20.65	4.49	1903.24

3) *Performance under different verification methods*: Since different formal verification techniques may already be in place for various safety-critical systems, we also show in Tab. 1 that our set-based trained agents can be verified using alternative methods that do not rely on zonotopes. All toolboxes (Althoff, 2015; Xiangru Zhong, 2024; Bogomolov et al., 2019; Tran et al., 2020) verified the set-based agent, while only two verified the standard trained agent, taking much longer to do so. Hence, our set-based agents are easier to verify.

4) *Extension to ensemble methods*: Our proposed set-based reinforcement learning can be directly extended for the Twin Delayed Deep Deterministic policy gradient algorithm (TD3) (Fujimoto et al., 2018) and other ensemble algorithms (Januszewski et al., 2021). Fig. 5b shows overall a similar performance as with DDPG and for the *SA-PC* version, a better performance for small perturbation radii ϵ_{test} of the set-based TD3 algorithm for the *1D Quadrotor*.

5) *Scalability*: In Appendix A.3, we show that our algorithm scales to more complex benchmarks with intricate dynamics and larger state and action spaces (Todorov et al., 2012). While the worst-case reward cannot be directly evaluated for every benchmark, we show that for systems lacking formal verification toolboxes, our algorithm scales effectively in practice and obtains robust neural networks.

6 CONCLUSION

We introduce the first set-based reinforcement learning algorithm. Unlike other algorithms that rely on adversarial inputs, our approach is unique in using set-based neural network training, working with entire sets of inputs and gradients. The gradient sets used for training are computed from set-based loss functions that are motivated by a rigorous analysis of the underlying set propagation. Our experimental results on different benchmarks demonstrate the efficacy of our approach. Particularly, our trained controllers can be formally verified for large perturbation sets, which is essential for their deployment in safety-critical environments. Consequently, set-based reinforcement learning is an effective, novel approach for training robust neural network controllers.

REFERENCES

- 540
541
542 Matthias Althoff. *Reachability analysis and its application to the safety assessment of autonomous*
543 *cars*. PhD thesis, Technical University of Munich, 2010.
- 544
545 Matthias Althoff. An introduction to CORA 2015. In *Proc. of Int. Workshop on Applied Verification*
546 *of Continuous and Hybrid Systems*, pp. 120–151, 2015.
- 547
548 Chi Au and Judy Tam. Transforming variables using the dirac generalized function. *The American*
549 *Statistician*, 53(3):270–272, 1999.
- 550
551 Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4.
Springer, 2006.
- 552
553 Sergiy Bogomolov, Marcelo Forets, Goran Frehse, Kostiantyn Potomkin, and Christian Schilling.
Juliareach: a toolbox for set-based reachability. In *Proc. of the Int. Conf. on Hybrid Systems:*
554 *Computation and Control (HSCC)*, pp. 39–44. Association for Computing Machinery, 2019.
- 555
556 Christopher Brix, Stanley Bak, Changliu Liu, and Taylor T. Johnson. The fourth int. verification of
557 neural networks competition (VNN-COMP 2023): Summary and results, 2023.
- 558
559 Aditya M Deshpande, Ali A Minai, and Manish Kumar. Robust deep reinforcement learning for
560 quadcopter control. *Int. Federation of Automatic Control (IFAC-PapersOnLine)*, 54(20):90–95,
2021.
- 561
562 Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in
563 actor-critic methods. In *Proc. of the Int. Conf. on Machine Learning (ICML)*, 2018.
- 564
565 Antoine Girard. Reachability of uncertain linear systems using zonotopes. In *Proc. of the Int. Conf.*
566 *on Hybrid Systems: Computation and Control (HSCC)*, pp. 291–305, 2005.
- 567
568 Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial
examples. In *Proc. of the Int. Conf. on Learning Representations (ICLR)*, 2015.
- 569
570 Geoffrey E Hinton and Zoubin Ghahramani. Generative models for discovering sparse distributed
571 representations. *Philosophical Transactions of the Royal Society of London. Series B: Biological*
572 *Sciences*, 352(1358):1177–1190, 1997.
- 573
574 Sandy H. Huang, Nicolas Papernot, Ian J. Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial
575 attacks on neural network policies. *Proc. Int. Conf. on Learning Representations (ICLR), Workshop*
Track, 2017.
- 576
577 Piotr Januszewski, Mateusz Olko, Michał Królikowski, Jakub Swiatkowski, Marcin Andrychowicz,
578 Łukasz Kuciński, and Piotr Miłoś. Continuous control with ensemble deep deterministic policy
579 gradients. In *Deep RL Workshop NeurIPS*, 2021.
- 580
581 Guy Katz, Clark Barrett, David L Dill, Kyle Julian, and Mykel J Kochenderfer. Reluplex: An efficient
582 smt solver for verifying deep neural networks. In *Proc. of the Int. Conf. on Computer Aided*
Verification (CAV), pp. 97–117, 2017.
- 583
584 Jens Kober, J Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The Int.*
585 *Journal of Robotics Research*, 32(11):1238–1274, 2013.
- 586
587 Lukas Koller, Tobias Ladner, and Matthias Althoff. Set-based training for neural network verification.
Transactions on Machine Learning Research, 2025.
- 588
589 Hanna Krasowski, Jakob Thumm, Marlon Müller, Lukas Schäfer, Xiao Wang, and Matthias Al-
590 thoff. Provably safe reinforcement learning: Conceptual analysis, survey, and benchmarking.
Transactions on Machine Learning Research, 2023.
- 591
592 Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,
593 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *Proc.*
of the Int. Conf. on Learning Representations (ICLR), 2016.

- 594 Diego Manzananas Lopez, Matthias Althoff, Luis Benet, Xin Chen, Jiameng Fan, Marcelo Forets,
595 Chao Huang, Taylor T Johnson, Tobias Ladner, Wenchao Li, Christian Schilling, and Qi Zhu.
596 ARCH-COMP22 category report: Continuous and hybrid systems with linear continuous dynamics.
597 In *Proc. of the Int. Workshop on Applied Verification of Continuous and Hybrid Systems (ARCH22)*,
598 volume 90, pp. 142–184, 2022.
- 599 Björn Lütjens, Michael Everett, and Jonathan P How. Certified adversarial robustness for deep
600 reinforcement learning. In *Proc. of the Conf. on Robot Learning (CoRL)*, pp. 1328–1337, 2020.
- 602 Laura Lützwow and Matthias Althoff. Reachability analysis of armax models. In *Proc. of the IEEE
603 Conf. on Decision and Control (CDC)*, pp. 7027–7034, 2023.
- 604 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
605 Towards deep learning models resistant to adversarial attacks. In *Proc. of the Int. Conf. on
606 Learning Representations (ICLR)*, 2018.
- 608 Ajay Mandlekar, Yuke Zhu, Animesh Garg, Li Fei-Fei, and Silvio Savarese. Adversarially robust
609 policy learning: Active construction of physically-plausible perturbations. In *Proc. of the IEEE/RSS
610 Int. Conf. on Intelligent Robots and Systems (IROS)*, pp. 3932–3939, 2017.
- 612 Diego Manzananas Lopez, Matthias Althoff, Marcelo Forets, Taylor T. Johnson, Tobias Ladner, and
613 Christian Schilling. ARCH-COMP23 category report: Artificial intelligence and neural network
614 control systems (AINNCS) for continuous and hybrid systems plants. In *Proc. of Int. Workshop on
615 Applied Verification of Continuous and Hybrid Systems (ARCH23)*, 2023.
- 616 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,
617 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control
618 through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 619 Janosch Moos, Kay Hansel, Hany Abdulsamad, Svenja Stark, Debora Clever, and Jan Peters. Robust
620 reinforcement learning: A review of foundations and recent advances. *Machine Learning and
621 Knowledge Extraction*, 4(1):276–315, 2022.
- 623 OpenAI, Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew,
624 Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning
625 dexterous in-hand manipulation. *The Int. Journal of Robotics Research*, 39:3–20, 2020.
- 626 Anay Pattanaik, Zhenyi Tang, Shuijing Liu, Gautham Bommanan, and Girish Chowdhary. Robust
627 deep reinforcement learning with adversarial attacks. In *Proc. of the Int. Conf. on Autonomous
628 Agents and Multiagent Systems (AAMAS)*, pp. 2040–2042, 2018.
- 630 Lerrel Pinto, James Davidson, Rahul Sukthankar, and Abhinav Gupta. Robust adversarial rein-
631 forcement learning. In *Proc. of the Int. Conf. on Machine Learning (ICML)*, pp. 2817–2826,
632 2017.
- 633 Benjamin Recht. A tour of reinforcement learning: The view from continuous control. *Annual
634 Review of Control, Robotics, and Autonomous Systems*, 2:253–279, 2019.
- 636 David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller.
637 Deterministic policy gradient algorithms. In *Proc. of the Int. Conf. on Machine Learning (ICML)*,
638 volume 32, pp. 387–395, 2014.
- 639 Gagandeep Singh, Timon Gehr, Matthew Mirman, Markus Püschel, and Martin Vechev. Fast and
640 effective robustness certification. In *Advances in Neural Information Processing Systems (NeurIPS)*,
641 2018.
- 643 Nicholas Socci, Daniel Lee, and H Sebastian Seung. The rectified gaussian distribution. *Advances in
644 Neural Information Processing Systems (NeurIPS)*, 10, 1997.
- 646 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow,
647 and Rob Fergus. Intriguing properties of neural networks. In *Proc. of the Int. Conf. on Learning
Representations (ICLR)*, 2014.

- 648 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
649 In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033,
650 2012.
- 651 Hoang-Dung Tran, Xiaodong Yang, Diego Manzananas Lopez, Patrick Musau, Luan Viet Nguyen,
652 Weiming Xiang, Stanley Bak, and Taylor T. Johnson. NNV: The neural network verification tool
653 for deep neural networks and learning-enabled cyber-physical systems. In *Proc. of the Int. Conf.
654 on Computer Aided Verification (CAV)*, 2020.
- 655 Fan Wu, Linyi Li, Zijian Huang, Yevgeniy Vorobeychik, Ding Zhao, and Bo Li. CROP: Certifying
656 robust policies for reinforcement learning through functional smoothing. In *International Confer-
657 ence on Learning Representations*, 2022. URL [https://openreview.net/forum?id=
658 HOjLHr1Zhmx](https://openreview.net/forum?id=HOjLHr1Zhmx).
- 659 Huan Zhang Xiangru Zhong, Yuhao Jia. CROWN-Reach. [https://github.com/
660 Verified-Intelligence/CROWN-Reach](https://github.com/Verified-Intelligence/CROWN-Reach), 2024.
- 661 Teng Xiao and Suhang Wang. Towards off-policy learning for ranking policies with logged feedback.
662 *Proc. of the AAAI Conf. on Artificial Intelligence (AAAI)*, 36(8):8700–8707, 2022.
- 663 Chenxi Yang, Greg Anderson, and Swarat Chaudhuri. Certifiably Robust Reinforcement Learning
664 through Model-Based Abstract Interpretation. In *2024 IEEE Conf. on Secure and Trustworthy
665 Machine Learning (SaTML)*, pp. 233–251. IEEE Computer Society, 2024.
- 666 Zhaocong Yuan, Adam W. Hall, Siqi Zhou, Lukas Brunke, Melissa Greeff, Jacopo Panerati, and
667 Angela P. Schoellig. Safe-control-gym: A unified benchmark suite for safe learning-based control
668 and reinforcement learning in robotics. *IEEE Robotics and Automation Letters*, 7:11142–11149,
669 2022.
- 670 Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane Boning, and Cho-Jui Hsieh.
671 Robust deep reinforcement learning against adversarial perturbations on state observations. In
672 *Proc. of the 34th Int. Conf. on Neural Information Processing Systems (NeurIPS)*, 2020.
- 673 Huan Zhang, Hongge Chen, Duane Boning, and Cho Jui Hsieh. Robust deep reinforcement learning
674 against adversarial perturbations on state observations. In *Int. Conf. on Learning Representations
675 (ICLR)*, 2021a.
- 676 Huan Zhang, Hongge Chen, Duane Boning, and Cho-Jui Hsieh. Robust reinforcement learning on
677 state observations with learned optimal adversary. *arXiv preprint arXiv:2101.08452*, 2021b.

684 A EVALUATION DETAILS

686 A.1 BENCHMARK DESCRIPTIONS

687 Let us briefly describe the specifications of the remaining benchmarks in this section.

688 *ID Quadrotor*: The state is $s = [z \quad \dot{z}]^\top$, with altitude z , velocity \dot{z} , and dynamics (Yuan et al.,
689 2022):

$$690 \dot{s} = \begin{bmatrix} \dot{z} & \frac{a+1}{2m} - g \end{bmatrix}^\top, \quad (29)$$

691 with action space $a \in [-1, 1]$, gravity $g = 9.81$, and mass $m = 0.05$. Starting from initial
692 states $s_0 \in [[-4 \quad 0]^\top, [4 \quad 0]^\top]$, the Quadrotor is stabilized at $s^* = \mathbf{0}$; the reward function is
693 $r(s_t, a_t) = -[1 \quad 0.01] |s_{t+1} - s^*|$ for a time horizon of $3s$.

694 *Navigation Task*: We use a unicycle model with states $s = [x \quad y \quad \theta \quad v]^\top$ and dynamics (Lopez
695 et al., 2022):

$$696 \dot{s} = \begin{bmatrix} v \cos(\theta) & v \sin(\theta) & a_{(1)} & a_{(2)} \end{bmatrix}^\top, \quad (30)$$

697 with action space $a \in [-1, 1]$. From starting point $s_0 = [3 \quad 3 \quad 0 \quad 0]^\top$, the task is to navigate
698 to the goal $s^* = \mathbf{0}$ without colliding with an obstacle $\mathcal{O} = [1, 2 \cdot \mathbf{1}]$. The reward function is

Table 2: Training parameters for *PA-PC*, *SA-PC* and *SA-SC*.

Parameter	DDPG	TD3
Actor learning rate	$1 \cdot 10^{-4}$	$1 \cdot 10^{-4}$
Critic learning rate	$1 \cdot 10^{-3}$	$1 \cdot 10^{-3}$
Critic L_2 weight regularization λ_Q	0.01	0
Discount factor γ	0.99	0.99
Target update factor τ	0.05	0.05
Exploration noise std. deviation σ	0.1	0.1, 0.2
Batchsize	64	64
Buffersize	$1 \cdot 10^6$	$1 \cdot 10^6$
Episodes	2000	2000
Perturbation radius ϵ	0.1	0.1
Actor weighting factor η_μ	0.1	0.1
Critic weighting factor η_Q	0.01	0.01

$r(s_t, a_t) = -[1 \ 1 \ 0 \ 0] |s_{t+1} - s^*| - c$, with $c = 1$ if $s_{t+1} \in \mathcal{O}$ and otherwise $c = 0$. We consider a time horizon of $8s$.

Inverted Pendulum: The state is $s = [\theta \ \dot{\theta}]^\top$, with angle θ , angular velocity $\dot{\theta}$, dynamics (Krasowski et al., 2023):

$$\dot{s} = \left[\dot{\theta} \quad \frac{g}{l} \sin(\theta) + \frac{1}{ml^2} a \right]^\top, \quad (31)$$

with action space $a \in [-15, 15]$, gravity $g = 9.81$, mass $m = 1$, and length $l = 1$. The goal to stabilize the pendulum in the upright position $s^* = \mathbf{0}$; the reward function is $r(s_t, a_t) = -[1 \ 0.01] |(s_{t+1} - s^*)|$ for a time horizon of $3s$.

2D Quadrotor: The state of the system is defined as $s = [x \ \dot{x} \ z \ \dot{z} \ \theta \ \dot{\theta}]^\top$, with horizontal displacement x , horizontal velocity \dot{x} , altitude z , vertical velocity \dot{z} , angle θ , angular velocity $\dot{\theta}$ and dynamics (Yuan et al., 2022):

$$\dot{s} = \begin{bmatrix} \dot{x} \\ \sin(\theta) \frac{\tilde{a}_{(1)} + \tilde{a}_{(2)}}{m} \\ \dot{z} \\ \cos(\theta) \frac{\tilde{a}_{(1)} + \tilde{a}_{(2)}}{m} - g \\ \dot{\theta} \\ \frac{l(\tilde{a}_{(2)} - \tilde{a}_{(1)})}{\sqrt{2}J_y} \end{bmatrix}, \quad (32)$$

where $\tilde{a} = (1 + \frac{1}{2}a) \frac{mg}{2}$ with action $a \in [-1, 1] \subset \mathbb{R}^2$. The constant $g = 9.81$ defines gravity, $m = 0.027$ is the mass of the Quadrotor, $l = 0.0397$ and $J_y = 1.4 \cdot 10^{-4}$ defines the arm length of the propeller mount and the moment of inertia around the y axis. The reward function is given by $r(s_t, a_t) = -[1 \ 0.01 \ 1 \ 0.01 \ 0 \ 0] |s_{t+1} - s^*|$ for a time horizon of $3s$, with the goal to stabilize the Quadrotor at $s^* = \mathbf{0}$.

A.2 HYPERPARAMETERS

We give the hyperparameters used to train the networks in our evaluation in Tab. 2.

A.3 ADDITIONAL EXPERIMENTS

Quadrotor 2D: We provide additional experiments on the *2D Quadrotor* and the *MuJoCo Hopper-v2* benchmark (Todorov et al., 2012). We again average the results over the last five agents in each training run and compute the mean and a 95% confidence interval across five independent random seeds. Fig. 8 compares the lower bounds $\underline{V}_\mu(s_0)$ of the different training algorithms.

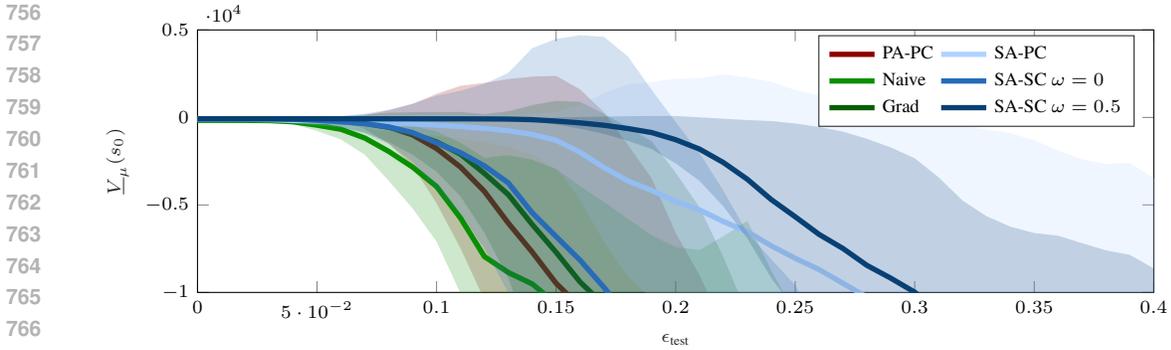


Figure 8: Comparison of $V_{\mu}(s_0)$ for Quadrotor 2D.

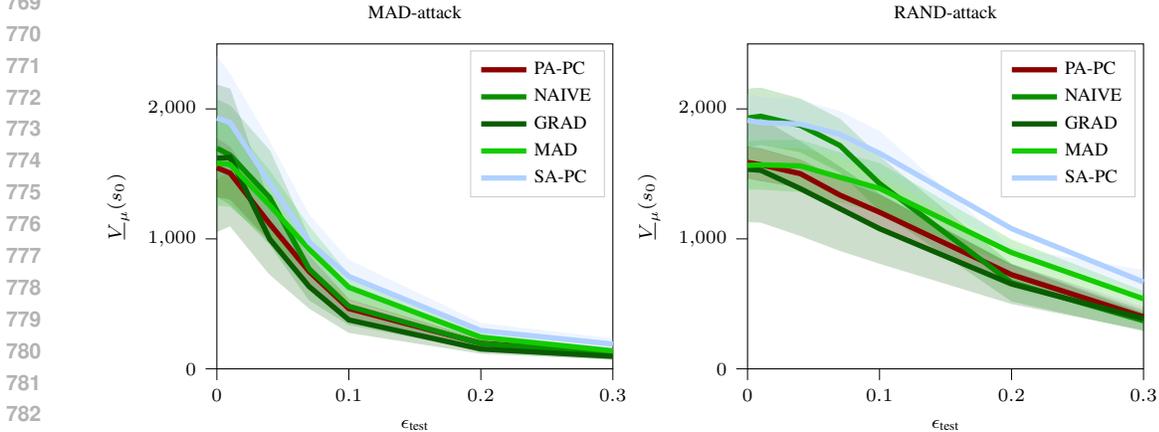


Figure 9: Comparison of $V_{\mu}(s_0)$ for *Hopper-v2* approximated by MAD (Zhang et al., 2020). Figure 10: Comparison of $V_{\mu}(s_0)$ for *Hopper-v2* approximated with uniform random noise.

Locomotion: Since, to our best knowledge, no existing verification toolbox can compute reachable sets for locomotion benchmarks, we are unable to provide a formal lower bound for *Hopper-v2*. However, we want to stress this is not a limitation of our training approach but rather of the subsequent verification step. Thus, to demonstrate the scalability of our approach to such benchmarks, we approximate $V_{\mu}(s_0)$ using 50 trajectories under the MAD-attack (Zhang et al., 2020) and 200 trajectories perturbed with uniform random noise sampled from the ℓ_{∞} ball with radius ϵ_{test} . The corresponding worst-case returns are shown in Figure 9 (MAD attack) and Figure 10 (random noise). For *Hopper-v2*, these empirical estimates provide an upper bound on the worst-case performance and demonstrate the scalability of our method. Notably, even under noise-free conditions, the worst returns of SA-PC trained agents, evaluated over randomly initialized trajectories, consistently exceed those of agents trained with a point-wise robustness criterion. We additionally provide videos illustrating agent behaviors under the uniform-random² and MAD³ attacks using the first random seed, with $\epsilon_{\text{test}} = 0.1$ and $\epsilon_{\text{test}} = 0.075$. To better visualize the agents’ failure modes, we disable early termination in these demonstrations. The results clearly show that the MAD attack is more effective at degrading agent performance compared to the uniform-random baseline. Notably, across both attack types, our SA-PC agent consistently demonstrates superior robustness and overall performance.

Hyperparameter Discussion: Set-based reinforcement learning introduces the additional parameters η_{μ} and η_Q , which appear in the set-based policy gradient Def. 3.2 and the set-based regression loss Prop. 3.1. In Tab. 2, we list the hyperparameters used in our benchmarks. Notably η_{μ} and η_Q are fixed in all experiments. We ease tuning these parameters by choosing β and α in Eq. (20) and (25) as proposed by Koller et al. (2025). For a given perturbation radius ϵ_{test} , we study for SA-PC on the 1D

²Video: <https://t1p.de/ar23e>

³Video: <https://t1p.de/mye61>

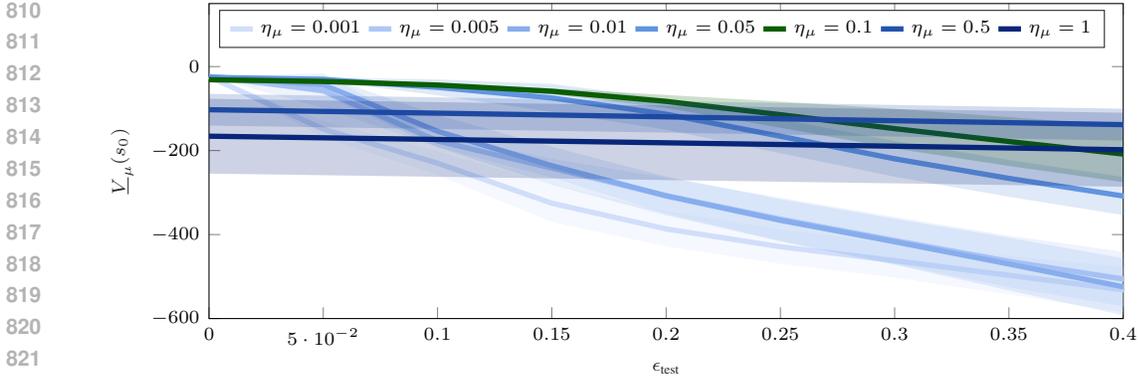


Figure 11: Tradeoff in verified performance for hyper parameter ablation η_μ .

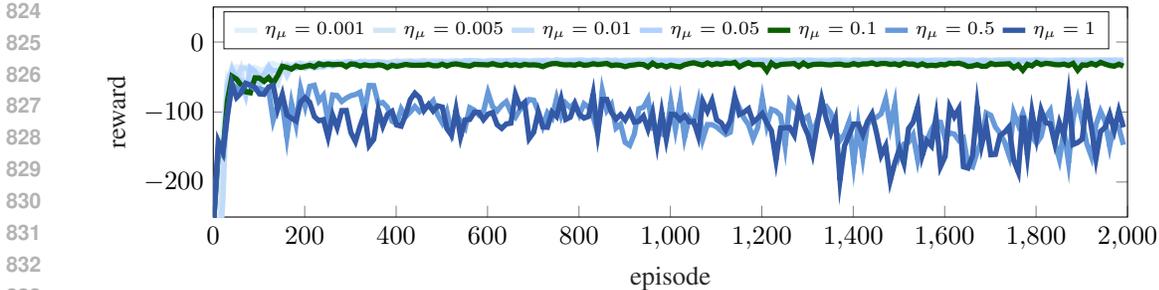


Figure 12: Learning behavior for hyper parameter ablation η_μ .

Quadrotor benchmark, the effect of varying η_μ . The hyperparameter η_μ directly scales the set-based gradient of the zonotope generators. As shown in Fig. 11, increasing η_μ , improves robustness: verified performance near the noise-free case $\epsilon_{\text{test}} = 0$ decreases slightly for larger η_μ , whereas verified performance at larger ϵ_{test} increases. This matches the intended trade-off between nominal (zero-perturbation) performance and robustness to larger perturbations. We therefore conclude that the choice of $\eta_\mu = 0.1$ (green) is considerably good, since it still achieves high verified rewards for large ϵ_{test} , while it preserves good rewards for the noise-free case $\epsilon_{\text{test}} = 0$. We also report the learning history for different hyperparameter settings in Fig. 12. Moderate values of η_μ lead to stable convergence, while very large η_μ slow or prevent convergence to the optimum.

We next conduct an ablation study on the hyperparameter η_Q for SA-SC, keeping $\eta_\mu = 0.1$ fixed. As shown in Fig. 13, large values of η_Q (e.g. $\eta_Q = 1$) make the value function highly contractive, which leads to reduced performance for small perturbation radii ϵ_{test} and to convergence difficulties, as illustrated in Fig. 14. As η_Q decreases, convergence improves and the value function captures more of the variability induced by observation noise. Consequently, the actor learns a more conservative policy, resulting in higher verified performance for larger ϵ_{test} . We find that $\eta_Q = 0.01$ (green) provides the best verified performance across all perturbation radii. Further reducing η_Q , causes verified performance to decline again.

Learning History: In Fig. 16, we present the full learning curves, which show that set-based reinforcement learning exhibits convergence behavior comparable to the adversarial baselines. For clarity, we report the rewards evaluated without observation noise during training. Additionally, we report the learning curves for the Locomotion benchmark (Hopper-v2) from Fig. 9 and 10 in Fig. 15. For this analysis, we evaluate both algorithms using the reward without observation noise. We find that, also in the

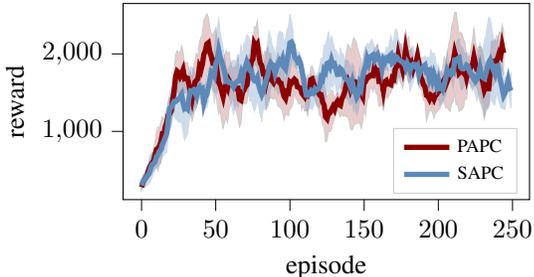


Figure 15: Full learning history for Hopper-v2.

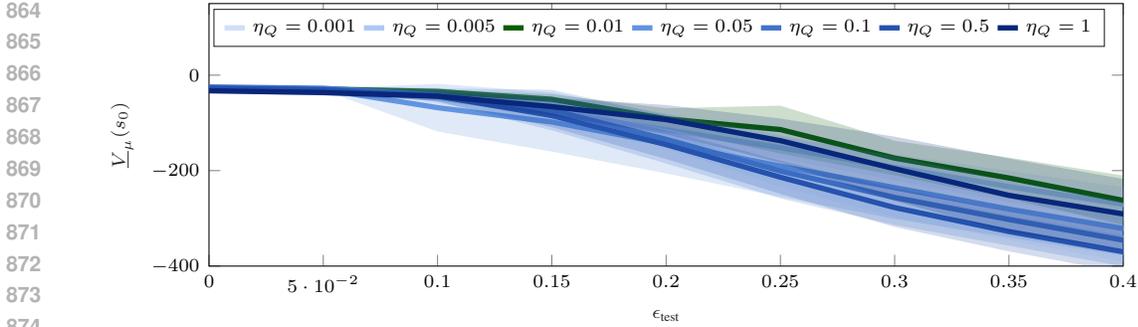


Figure 13: Hyper parameter ablation η_Q .

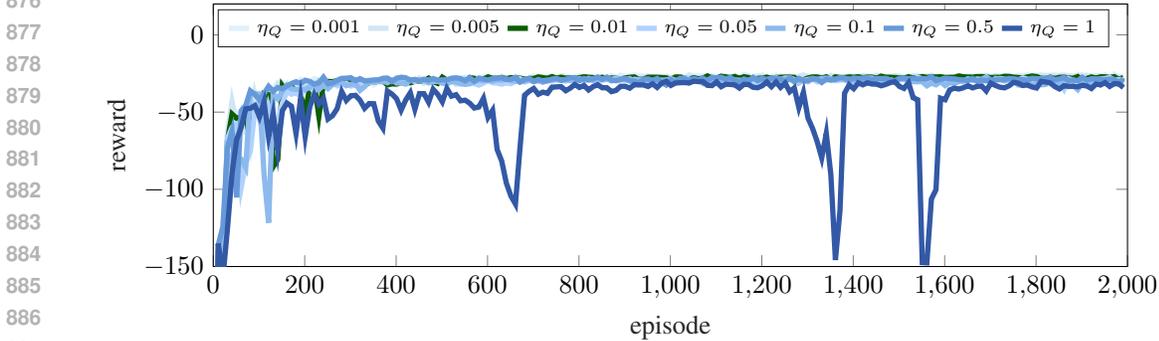


Figure 14: Learning behavior for hyper parameter ablation η_Q .

Locomotion setting, our set-based method exhibits learning dynamics similar to those of the vanilla DDPG Lillicrap et al. (2016) implementation. Based on the hyperparameter analysis in Fig. 11 and 12, we also select $\eta_\mu = 0.1$ for the *Hopper-v2* benchmark.

A.4 RUNTIME COMPLEXITY

Building on set-based neural network training (Koller et al., 2025), our algorithm has polynomial time complexity. Based on the analysis in Koller et al. (2025, Prop. 17), we observe that *SA-PC* has runtime complexity $\mathcal{O}(n_{max,\mu}^2 q_\mu \kappa_\mu)$, where $n_{max,\mu}$ is the maximum number of neurons per layer for the actor. The sum $q = n_0 + \sum_{k \in \kappa_\mu} n_k$ includes the number of initial independent noise generators n_0 and the total number of neurons in the network. For *SA-SC*, we propagate the zonotopes through both networks, actor and critic. Therefore, the time complexity is hence given by

$$\mathcal{O}(n_{max,\mu+Q}^2 q_{\mu+Q} \kappa_{\mu+Q}), \tag{33}$$

where $n_{max,\mu+Q}$ is the maximum number of neurons per layer in either of actor or critic. Second, the total number of layers $\kappa_{\mu+Q} = \kappa_\mu + \kappa_Q$ is the sum of layers in the actor κ_μ and the critic κ_Q . Finally, $q_{\mu+Q} = n_0 + \sum_{k \in \kappa_\mu} n_k + \sum_{k \in \kappa_Q} n_k$ is now defined as the sum of the initial independent noise generators and the total number of neurons in both networks. We additionally remark that consequently our algorithm depends on n_0 in q , and thereby grows with the state dimension, if the noise is considered to be independent per dimension.

Set-based reinforcement learning can be efficiently computed batch-wise on a GPU, but the memory load remains challenging. Especially for *SA-SC*, storing entire action sets with many generators in \mathcal{B} is memory-consuming and computationally more expensive. The different runtimes for 10 learning epochs are listed in Tab. 3 and were run on a server with two *AMD EPYC 7763* 64 core processors, 2 TB RAM, and an *NVIDIA A100-PCIE* 40 GB GPU.

A.5 LIMITATIONS

We do not encode safety as an explicit constraint during training of the *Navigation Task* benchmark; instead, we incorporate safety considerations indirectly by subtracting a penalty from the reward

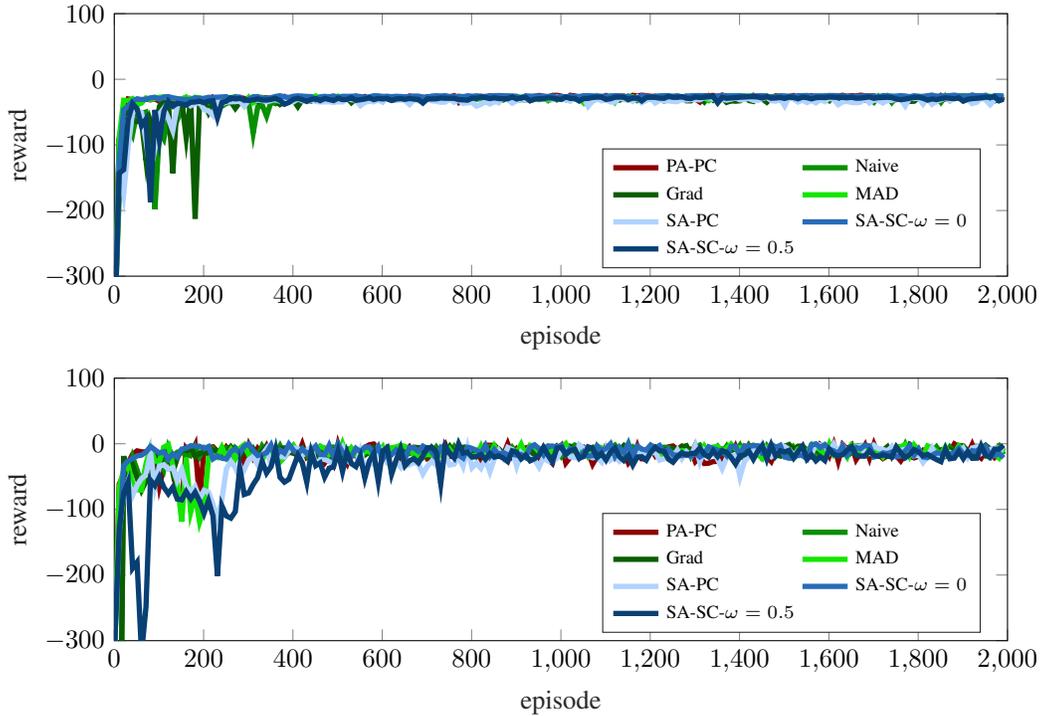


Figure 16: Full learning history for (top) 1D Quadrotor and (bottom) Inverted pendulum.

Table 3: Training times [s/10 epochs]

Benchmark	PA-PC	Naive	Grad	SA-PC	SA-SC
1D Quad.	1.58	2.10	1.98	2.77	7.92
Pendulum	1.82	2.04	2.02	2.87	7.66
Nav. Task	2.35	2.89	2.84	4.35	12.43

972 signal. Consequently, any claim of safe behavior at deployment cannot rest solely on the provable
 973 lower bound of the cumulative reward, this bound guarantees only worst-case performance, not safety
 974 itself. To establish formal safety guarantees, the trained policy must therefore undergo separate,
 975 formal verification rather than relying on its reward lower bound.
 976

977 B PROOFS

979 **Proposition 3.1.** *Given output set $\mathcal{Q}_i = \langle c_{\mathcal{Q}_i}, G_{\mathcal{Q}_i} \rangle_Z \subset \mathbb{R}$ and target $y_i \in \mathbb{R}$, the set-based
 980 regression loss is defined as*

$$981 L_{Reg}(y_i, \mathcal{Q}_i) = \underbrace{1/2(c_{\mathcal{Q}_i} - y_i)^2}_{\text{standard training loss}} + \underbrace{\eta_Q/\epsilon \ln \text{Dia}(G_{\mathcal{Q}_i})}_{\text{robustness}},$$

984 with weighting factor $\eta_Q \in \mathbb{R}_+$ and perturbation radius $\epsilon \in \mathbb{R}_+$. The gradient of L_{Reg} w.r.t. \mathcal{Q}_i is:

$$985 \nabla_{\mathcal{Q}_i} L_{Reg}(y_i, \mathcal{Q}_i) = \langle c - y_i, \eta_Q/\epsilon \ln \text{Dia}'(G_{\mathcal{Q}_i}) \rangle_Z.$$

988 *Proof.* The gradient w.r.t. a zonotope is represented as a zonotope as well, consisting of the gradient
 989 w.r.t. the center and the gradient w.r.t. the generator matrix (Koller et al., 2025, Def. 8). Hence, the
 990 gradient of the set-based regression is:

$$\begin{aligned} 991 \nabla_{\mathcal{Q}_i} L_{Reg}(y_i, \mathcal{Q}_i) &= \langle \nabla_{c_{\mathcal{Q}_i}} L_{Reg}(y_i, \mathcal{Q}_i), \nabla_{G_{\mathcal{Q}_i}} L_{Reg}(y_i, \mathcal{Q}_i) \rangle_Z \\ 992 &= \left\langle c_{\mathcal{Q}_i} - y_i, \frac{\eta_Q}{\epsilon} \text{diag}(|G_{\mathcal{Q}_i}| \mathbf{1})^{-1} \text{sign } G_{\mathcal{Q}_i} \right\rangle_Z \\ 993 &= \left\langle c_{\mathcal{Q}_i} - y_i, \frac{\eta_Q}{\epsilon} \ln \text{Dia}'(G_{\mathcal{Q}_i}) \right\rangle_Z. \quad \square \end{aligned}$$

997 **Proposition 4.1.** *Given a neural network with ReLU-activations and an input set \mathcal{H}_{k-1} with the
 998 enclosing interval $[l_{k-1}, u_{k-1}] \supseteq \mathcal{H}_{k-1}$, the expected value of the enclosure of the k -th layer is*

$$999 \mathbb{E}_{h_k \sim \mathcal{H}_k} [h_k] = \mathbb{E}_{h_{k-1} \sim \mathcal{Z}(l_{k-1}, u_{k-1})} [L_k(h_{k-1})],$$

1000 with $\mathcal{H}_k = \text{enclose}(L_k, \mathcal{H}_{k-1})$ having minimal approximation errors.

1002 *Proof.* We split cases on the type of layer L_k .

1004 *Case (1): Linear layer* The expected value is preserved by linearity of the expectation:

$$\begin{aligned} 1005 \mathbb{E}_{h_k \sim \mathcal{H}_k} [h_k] &\stackrel{(15)}{=} c_k = W_k c_{k-1} + b_k = W_k \mathbb{E}_{h_{k-1} \sim \mathcal{Z}(l_{k-1}, u_{k-1})} [h_{k-1}] + b_k \\ 1006 &= \mathbb{E}_{h_{k-1} \sim \mathcal{Z}(l_{k-1}, u_{k-1})} [W_k h_{k-1} + b_k] = \mathbb{E}_{h_{k-1} \sim \mathcal{Z}(l_{k-1}, u_{k-1})} [L_k(h_{k-1})]. \end{aligned}$$

1009 *Case (2): ReLU layer* Activation functions are applied element-wise, thus we consider each dimension
 1010 individually; to avoid clutter, we drop the dimension index $x_{(i)}$. We distinguish between three
 1011 cases: (2a) $l_{k-1}, u_{k-1} \leq 0$, (2b) $l_{k-1}, u_{k-1} \geq 0$, (2c) $l_{k-1} < 0 < u_{k-1}$. For cases (i) and (ii),
 1012 ReLU is linear, thus by linearity of the expectation the expected value is preserved. For case (iii),
 1013 we approximate the ReLU activation function with the affine map $m_k x + t_k$. We first derive a
 1014 condition to ensure preserving the expected value. After that we show the slope that minimizes the
 1015 approximation errors is:

$$1016 m_k = \frac{\text{ReLU}(u_{k-1}) - \text{ReLU}(l_{k-1})}{u_{k-1} - l_{k-1}} = \frac{u_{k-1}}{u_{k-1} - l_{k-1}}. \quad (34)$$

1018 The expected value is preserved if the offset t_k to satisfies the condition:

$$\begin{aligned} 1020 \mathbb{E}_{h_k \sim \mathcal{H}_k} [h_k] &\stackrel{!}{=} \mathbb{E}[\text{ReLU}(h_{k-1})] \iff c_k = \mathbb{E}[\text{ReLU}(h_{k-1})] \\ 1021 &\iff m_k c_{k-1} + t_k = \mathbb{E}[\text{ReLU}(h_{k-1})] \\ 1022 &\iff t_k = \mathbb{E}[\text{ReLU}(h_{k-1})] - m_k c_{k-1}, \end{aligned} \quad (35)$$

1024 with $h_{k-1} \sim \mathcal{Z}(l_{k-1}, u_{k-1})$. Hence, we fix the offset t_k w.r.t. the slope m_k . Moreover, to
 1025 find the optimal slope m_k , we compute the expected value $\mathbb{E}[\text{ReLU}(h_{k-1})]$ using the probability
 density function $f_{\mathcal{H}_k}$ for the distribution of $\text{ReLU}(\mathcal{Z}(l_{k-1}, u_{k-1}))$. Therefore, we first compute the

probability mass below 0 using the cumulative distribution function (Hinton & Ghahramani, 1997; Soggi et al., 1997):

$$\begin{aligned} F_{\mathcal{W}(l_{k-1}, u_{k-1})}(0) &= \int_{-\infty}^0 f_{\mathcal{W}(l_{k-1}, u_{k-1})}(h_{k-1}) dh_{k-1} \\ &= \int_{l_{k-1}}^0 \frac{1}{u_{k-1} - l_{k-1}} dh_{k-1} = \frac{-l_{k-1}}{u_{k-1} - l_{k-1}}. \end{aligned}$$

The probability mass is concentrated in a peak at zero using the Dirac delta $\delta(x)$ (Au & Tam, 1999). Hence, we obtain for a compactly supported function $h(x)$, with respect to the measure δ the Lebesgue integral:

$$\delta(x) = \begin{cases} \infty & x = 0, \\ 0 & \text{else,} \end{cases} \quad \int_{-\infty}^{\infty} h(x) \delta(x) dx = h(0).$$

Thus, the probability density function for the post-activation is

$$f_{\mathcal{H}_k}(h_k) = \begin{cases} \frac{1-l_{k-1}}{u_{k-1}-l_{k-1}} \delta(h_k), & 0 \leq h_k < u_{k-1}, \\ 0 & \text{otherwise.} \end{cases}$$

The resulting density is composed of the uniform distribution for the support $h_k > 0$ and the aggregated probability mass for the negative support h_{k-1} in the Dirac spike at $h_k = 0$. This follows immediately from the definition of $\text{ReLU}(h_{k-1})$. Thus, the expected value of the transformed distribution is given by:

$$\begin{aligned} \mathbb{E}_{h_{k-1} \sim \mathcal{W}(l_{k-1}, u_{k-1})}[\text{ReLU}(h_{k-1})] &= \int_{-\infty}^{\infty} \text{ReLU}(h_{k-1}) f_{\mathcal{H}_k}(\text{ReLU}(h_{k-1})) dh_{k-1} \\ &= \int_0^{u_{k-1}} h_k f_{\mathcal{H}_k}(h_k) dh_k \\ &= \int_0^{u_{k-1}} h_k \frac{1-l_{k-1}}{u_{k-1}-l_{k-1}} \delta(h_k) dh_k \\ &= \frac{h_k^2}{2(u_{k-1}-l_{k-1})} \Big|_0^{u_{k-1}} = \frac{u_{k-1}^2}{2(u_{k-1}-l_{k-1})}. \end{aligned}$$

Hence,

$$t_k \stackrel{(35)}{=} \frac{u_{k-1}^2}{2(u_{k-1}-l_{k-1})} - m_k c_{k-1} \stackrel{(6)}{=} \frac{1}{2} \left(\frac{u_{k-1}^2}{(u_{k-1}-l_{k-1})} - m_k (u_{k-1} + l_{k-1}) \right). \quad (36)$$

We now find the slope m_k that minimizes the approximation errors; by satisfying (35) we ensure preserving the expected value. For concise notation, we abbreviate the approximation error at x with slope m_k by

$$d_x(m_k) := |(m_k x + t_k) - \text{ReLU}(x)|. \quad (37)$$

We optimize the slope m_k for minimal approximation errors:

$$\min_{m_k} \max_{x \in [l_{k-1}, u_{k-1}]} d_x(m_k).$$

With (Koller et al., 2025, Prop. 7), we know that the approximation error are located at $x \in \{l_{k-1}, 0, u_{k-1}\}$ and rewrite:

$$\min_{m_k} \max_{x \in \{l_{k-1}, 0, u_{k-1}\}} d_x(m_k).$$

From $l_{k-1} < 0 < u_{k-1}$, we can simplify

$$\begin{aligned} d_{l_{k-1}}(m_k) &= |m_k l_{k-1} + t_k|, \\ d_0(m_k) &= |t_k|, \\ d_{u_{k-1}}(m_k) &= |m_k u_{k-1} + t_k - u_{k-1}|. \end{aligned} \quad (38)$$

1080 Moreover, we know that at least one approximation error is greater 0:
1081

$$1082 \quad d_{l_{k-1}}(m_k) > 0 \vee d_0(m_k) > 0 \vee d_{u_{k-1}}(m_k) > 0.$$

1083 Hence, the optimal slope m_k is located at an intersection of two error functions $d_{x_1}(m_k) =$
1084 $d_{x_2}(m_k) \geq d_{x_3}(m_k)$ with $\{x_1, x_2, x_3\} = \{l_{k-1}, 0, u_{k-1}\}$. Thus, for each intersection we find
1085 the optimal slope:
1086

1087 *Case (2c.i):* $d_{l_{k-1}}(m_k) = d_0(m_k)$

$$1088 \quad d_0(m_k) = d_{l_{k-1}}(m_k)$$

$$1089 \quad \stackrel{(38)}{\iff} |t_k| = |m_k l_{k-1} + t_k|$$

$$1090 \quad \iff t_k^2 = (m_k l_{k-1} + t_k)^2$$

$$1091 \quad \iff 0 = (m_k l_{k-1})^2 + 2 m_k l_{k-1} t_k$$

$$1092 \quad \iff 0 = m_k^2 - \frac{u_{k-1}}{u_{k-1} - l_{k-1}} m_k$$

$$1093 \quad \stackrel{(36)}{\iff} m_k = 0 \vee m_k = \frac{u_{k-1}}{u_{k-1} - l_{k-1}}.$$

$$1094 \quad \iff$$

$$1095$$

$$1096$$

$$1097$$

1098 *Case (2c.ii):* $d_{l_{k-1}}(m_k) = d_{u_{k-1}}(m_k)$

$$1099 \quad d_{l_{k-1}}(m_k) = d_{u_{k-1}}(m_k)$$

$$1100 \quad \stackrel{(38)}{\iff} |m_k l_{k-1} + t_k| = |m_k u_{k-1} + t_k - u_{k-1}|$$

$$1101 \quad \iff (m_k l_{k-1} + t_k)^2 = (m_k u_{k-1} + t_k - u_{k-1})^2$$

$$1102 \quad \stackrel{(36)}{\iff} m_k u_{k-1} l_{k-1} = u_{k-1} \frac{u_{k-1}^2}{u_{k-1} - l_{k-1}} - u_{k-1}^2$$

$$1103 \quad \iff m_k = \frac{u_{k-1}}{u_{k-1} - l_{k-1}}.$$

$$1104$$

$$1105$$

$$1106$$

$$1107$$

$$1108$$

$$1109$$

1110 *Case (2c.iii):* $d_0(m_k) = d_{u_{k-1}}(m_k)$

$$1111 \quad d_0(m_k) = d_{u_{k-1}}(m_k)$$

$$1112 \quad \stackrel{(38)}{\iff} |t_k| = |m_k u_{k-1} + t_k - u_{k-1}|$$

$$1113 \quad \iff t_k^2 = (m_k u_{k-1} + t_k - u_{k-1})^2$$

$$1114 \quad \stackrel{(36)}{\iff} -m_k^2 l_{k-1} + m_k \left(\frac{l_{k-1} (2 u_{k-1} - l_{k-1})}{u_{k-1} - l_{k-1}} \right)$$

$$1115 \quad \iff = \frac{u_{k-1}^2}{u_{k-1} - l_{k-1}} + u_{k-1}$$

$$1116 \quad \iff m_k = 1 \vee m_k = \frac{u_{k-1}}{u_{k-1} - l_{k-1}}.$$

$$1117$$

$$1118$$

$$1119$$

$$1120$$

$$1121$$

$$1122$$

1123 In each case, we show that the proposed slope in (34) is optimal and minimizes the approximation
1124 errors $\underline{d}_k, \bar{d}_k \in \{d_{l_{k-1}}, d_0, d_{u_{k-1}}\}$.
1125

□

1126
1127
1128
1129
1130
1131
1132
1133

1134 C DISCLOSURE: USAGE OF LARGE LANGUAGE MODELS (LLMs)
1135

1136 We used a large language model (LLM) as a general-purpose assist tool to aid in polishing the writing
1137 and improving clarity of expression. The research ideas, methodology, analysis, and conclusions are
1138 the authors' own. The LLM did not contribute to the ideation, design, or execution of the research.
1139

1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187