# Bridging Dimensions: Confident Reachability for High-Dimensional Controllers

Yuang Geng, Jake Baldauf, Souradeep Dutta, Chao Huang, and Ivan Ruchkin

[1] University of Florida, FL, USA, `yuang.geng@ufl.edu`
[2] University of Florida, FL, USA, `jakebaldauf@ufl.edu`
[3] University of Pennsylvania, PA, USA, `duttaso@seas.upenn.edu`
[4] University of Southampton, UK, `chao.huang@soton.ac.uk`
[5] University of Florida, FL, USA, `iruchkin@ece.ufl.edu`

**Abstract.** Autonomous systems are increasingly implemented using end-to-end learning-based controllers. Such controllers make decisions that are executed on the real system, with images as one of the primary sensing modalities. Deep neural networks form a fundamental building block of such controllers. Unfortunately, the existing neural-network verification tools do not scale to inputs with thousands of dimensions — especially when the individual inputs (such as pixels) are devoid of clear physical meaning. This paper takes a step towards connecting exhaustive closed-loop verification with high-dimensional controllers. Our key insight is that the behavior of a high-dimensional controller can be approximated with several low-dimensional controllers. To balance the approximation accuracy and verifiability of our low-dimensional controllers, we leverage the latest verification-aware knowledge distillation. Then, we inflate low-dimensional reachability results with statistical approximation errors, yielding a high-confidence reachability guarantee for the high-dimensional controller. We investigate two inflation techniques — based on trajectories and control actions — both of which show convincing performance in three OpenAI gym benchmarks.

**Keywords:** reachability, neural-network control, conformal prediction

## 1 Introduction

End-to-end deep neural network controllers have been extensively used in executing complex and safety-critical autonomous systems in recent years [13,40,51,52]. In particular, *high-dimensional controllers* (HDCs) based on images and other high-dimensional inputs have been applied in areas such as autonomous car navigation [49,61] and aircraft landing guidance [47]. For example, recent work has shown the high performance of controlling aircraft to land on the runway with a vision-based controller [65]. For such critical applications, it is important to develop techniques with strong safety guarantees for HDC-controlled systems.

However, due to the high-dimensional nature of the input space, modern verification cannot be applied directly to systems controlled by HDCs [2, 43]. Current closed-loop verification tools, such as NNV [54], Verisig [30], Sherlock [18], and ReachNN* [28], are capable of combining a dynamical system and a *low-dimensional controller* (LDC) to verify a safety property starting from an initial
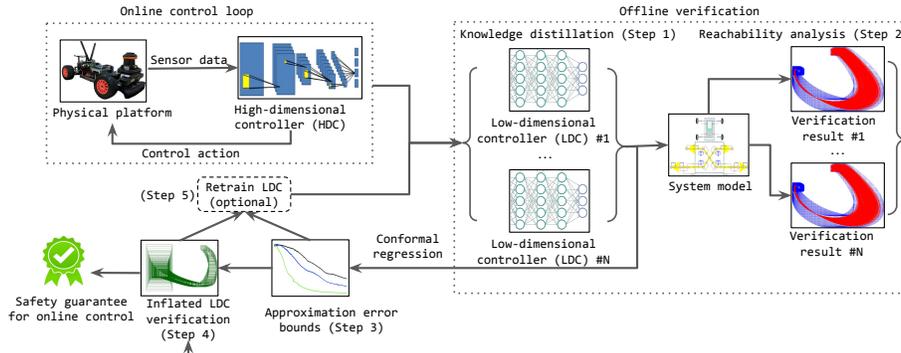
**Fig. 1.** Our verification approach for systems with high-dimensional controllers.

region of the low-dimensional input space, such as position-velocity states of a car. DeepReach [5] has pushed the boundary of applying Hamilton-Jacobi (HJ) reachability to systems with tens of state dimensions. However, such verification tools fail to scale for an input with thousands of dimensions (e.g., an image). One issue is that the dynamics of these dimensions are impractical to describe. Furthermore, the structure of an HDC is usually more complicated than that of an LDC, with convolution and pooling layers. For example, an image-based HDC may have hundreds of layers with thousands of neurons, whereas an LDC usually contains several layers with dozens of neurons, making HDC verification difficult.

To deal with these challenges, researchers have built abstractions of high-dimensional perception into the verification process. One work [31] verified a generative adversarial network (GAN) that creates images from states. Such methods cannot guarantee the GAN's accuracy or relation to reality, which becomes a major falsifiable assumption of their verification outcomes. Another work [47] built a precise mathematical model capturing the exact relationship between states and image pixels to verify the image-based controller, which is effortful and needs to be redone for each system. Inspired by previous work on decreasing the dimensions, we skillfully create verifiable low-dimensional controllers from high-dimensional ones.

This paper proposes an **end-to-end methodology** to verify systems with HDCs by employing the steps displayed in Fig. 1. Instead of verifying an HDC's safety directly over a complicated input space, our key idea is to approximate it with several LDCs so that we can reduce the HDC reachability problem to several LDC reachability problems. A crucial step is to upper-bound the difference between LDC and HDC, which we do statistically. Finally, we extend the reachable sets with the statistical bounds to obtain a safety guarantee for the HDC.

Since the input space and structure of the HDC are too complex to verify, we leverage *knowledge distillation* [25]—a model compression method—to train simplified "student models" (LDC) based on the information from the sophisticated "teacher model" (HDC). This training produces an LDC that is lightweight and amenable to closed-loop verification because it operates on dynamical states, not images. Moreover, due to the importance of the Lipschitz to minimizing the

overapproximation error [28, 29], our methodology adopts *two-objective gradient descent* [21], decreasing both the approximation error and Lipschitz constant.

After training the LDCs, we calculate the statistical upper bound of the discrepancy between the two controllers, since obtaining the true discrepancy is impractical. To this end, we rely on *conformal prediction* [22, 45, 48], one of the cutting-edge statistical methods to provide a lower bound of the confidence interval for prediction residuals without distributional assumptions or explicit dependency on the sample count. We propose *two conformal techniques* to quantify the difference between HDC- and LDC-controlled systems, by bounding: (i) the distance between their trajectories, and (ii) the difference between the actions produced by the HDC and LDC. We inflate reachable sets of the LDC system based on both bounds to obtain safety guarantees on the HDC system.

We evaluate our methodology on three popular verification case studies in OpenAI Gym [7]: inverted pendulum, mountain car, and cartpole. While our results are promising, verifying HDCs remains a challenging task with many improvement opportunities. The contributions of our work are three-fold:

1. Two verification approaches for high-dimensional controllers that combine reachability analysis and statistical inference to provide a safety guarantee for systems controlled by neural networks with thousands of inputs.
2. A novel neural-network approximation technique for training multiple LDCs that collectively mimic an HDC and reduce overapproximation error.
3. An implementation and evaluation of our verification approaches on three case studies: inverted pendulum, mountain car, and cartpole.

Sec. 2 introduces the background and defines our verification problem. Sec. 3 describes the details of our verification approach, which is evaluated in Sec. 4. Finally, we review the related work in Sec. 5 and conclude the paper in Sec. 6.

## 2 Background and Problem Setting

**High- and low-dimensional systems.** The original *high-dimensional closed-loop system* is a tuple $M_{hd} = (S, Z, U, s_0, f, c_{hd}, g)$. Here, the $S$ is the state space, $Z$ is the high-dimensional sensor space of so-called "images" (e.g., camera images or LIDAR scans), and the $U$ is the control action space, $s_0$ is the initial state, $f : S \times U \to S$ is the dynamics, and $c_{hd} : Z \times S \to U$ is the HDC. Note that $c_{hd}$ only uses a *subset* of state dimensions as input (e.g., a convolutional neural network with image and velocity inputs, but not position), getting the rest of the information from the image.

For mathematical convenience, we also define an (unknown) deterministic state-to-image generator as $g : S \to Z$ and the role and assumptions of generator $g$ are stated below. As a verifiable approximation of $M_{hd}$, our *low-dimensional closed-loop system* is defined as $M_{ld} = (S, U, s_0, f, c_{ld})$. Both $M_{hd}$ and $M_{ld}$ have the same state space and action space. The only difference is that the $M_{ld}$ has a low-dimensional controller $c_{ld} : S \to U$, which operates on the exact states.

**System execution.** The execution of $M_{hd}$ starts from the initial state $s_0$. Next, an image $z$ can be generated by image generator $g$ from that state. Then it is

fed into $c_{hd}$ to obtain a corresponding control action $u = c_{hd}(z)$, which is used to update the state via dynamics $f$. For $M_{ld}$, the execution proceeds similarly, except that the current state $s$ directly results in a control action $u = c_{ld}(s)$. Thus, we denote the *state at time $t$* starting from $s_0$ executed by $M_{hd}$ or $M_{ld}$ as $\varphi_{hd}(s_0, t)$ and $\varphi_{ld}(s_0, t)$ respectively. The *trajectory* of $M_{hd}$ is defined as a state sequence: $\tau_{hd}(s_0, T) = [s_0, \varphi_{hd}(s_0, 1), \ldots, \varphi_{hd}(s_0, T)]$, and similarly for $\tau_{hd}$.

Based on previous background, we define reachable sets and tubes:

**Definition 1 (Reachable set).** *Given an initial set $S_0$ and an integer time $t$, a* reachable set $\mathsf{rs}_M(S_0, t)$ *for (either) system $M$ contains all the states that can be reached from $S_0$ in $t$ steps:* $\mathsf{rs}_M(S_0, t) = \{\varphi_M(s_0, t) \mid \forall s_0 \in S_0\}$.

**Definition 2 (Reachable tube).** *Given an initial set $S_0$ and time horizon $T$, a* reachable tube $\mathsf{rt}_M(S_0, T)$ *for (either) system $M$ is a sequence of all the reachable sets from $S_0$ until time $T$:* $\mathsf{rt}_M(S_0, T) = [S_0, \mathsf{rs}_M(S_0, 1), ..., \mathsf{rs}_M(S_0, T)]$.

**Assumptions on image-state mapping $g$.** Our key challenge is establishing a mapping between the high-dimensional image space $Z$ and the low-dimensional state space $S$. Our verification methodology is based on the existence of a deterministic image generator $g$ that is part of $M_{hd}$. This generator is the true and *unknown* mechanism that creates images from states (e.g., a camera system). We do *not* assume or use an analyzable closed-form description of $g$. We also do not assume or verify any perception model (which obtains states from images).

Our only use of $g$ is in the training process for a limited dataset: we assume that we can "lab-study" an instrumented system $M_{hd}$ (e.g., with extra positioning sensors or human annotations) to label each image $z$ with a corresponding low-dimensional state $s$. In practice, this data collection would correspond to executing an autonomous system inside a controlled test facility. It's crucial to clarify that we are not using $g$ to generate any images from any states. Instead, we solely produce a restricted state-image paired dataset for training LDC.

To check the robustness of the above assumption in realistic scenarios, we will perform a sensitivity analysis by adding zero-mean Gaussian noise to the state-image mapping. For simplicity, this noise is added to the true state before generating an image from it, instead of adding pixel noise that may result in unrealistic images. The results of this evaluation will be discussed in Sec. 4.

**Verification problem.** Our problem is to guarantee that the high-dimensional system $M_{hd}$ reaches the goal set $G$ from an initial set $S_0$ within time $T$. To this end, we aim to compute reachable sets of the high-dimensional system $M_{hd}$ and intersect them with the goal set to obtain the verification verdict. Set $G$ is specified in low dimensions (i.e., using physical variables); however, the $M_{hd}$ behavior is determined by the images from generator $g$ and the HDC's response to them.

Thus, given initial set $S_0$, goal set $G$, system $M_{hd}$, and time horizon $T$, <u>our goal is to verify this assertion</u>:

$$\forall s_0 \in S_0 \cdot \mathsf{rs}_{M_{hd}}(S_0, T) \subseteq G \tag{1}$$

This problem can be divided into two parts : (a) approximating $M_{hd}$ with low-dimensional systems $M_{ld}^1, \ldots, M_{ld}^n$ and verifying them; (b) combining these

reachability results based on the approximation error bounds into a reachability verdict to solve the above $M_{hd}$ problem with statistical confidence.

## 3 Verification of High-Dimensional Systems

Considering the challenges of complex structure and dynamics of high-dimensional systems, and the difficulties of defining safety in high dimensions, our end-to-end approach is structured in five steps: (1) train low-dimensional controller(s), (2) perform reachability analysis on them, (3) compute statistical discrepancy bounds between high- and low-dimensional controllers, (4) inflate the reachable tubes from low-dimensional verification with these bounds, and (5) combine the verification results and repeat the process as if needed on different states/LDCs.

### Step 1: Training low-dimensional controllers

Given the aforementioned challenges of directly verifying $M_{hd}$, we plan to first verify the behavior in the low dimensions according to $M_{ld}$. Hence, we train a $c_{ld}$ to imitate the performance of $c_{hd}$ starting from a given state region, which serves as an input to Step 1 (our first iteration uses the full initial state region $S_0$ to train one $c_{ld}$). As a start, we collect the training data for $c_{ld}$: given the $c_{hd}$, access to image generator $g$, and the initial state space region $S_0$, we construct a supervised training dataset $\mathcal{D}_{tr} = \left\{ \left( \tau_{hd}(s_i, T), (u_1, ..., u_T)_i \right) \right\}_{i=1}^m$ by sampling the initial states $s_i \sim D_0$ from some given distribution $D_0$ (in practice, $D_0 = \mathrm{Uniform}(S_0)$).

Training a verifiable LDC has two conflicting objectives. On the one hand, we want to approximate the given $c_{hd}$ with minimal Mean Squared Error (MSE) on $\mathcal{D}_{tr}$. On the other hand, neural networks with smaller Lipschitz constants (Def. 11 in the Appendix) are more predictable and verifiable [15, 23, 50].

We balance the ability of the $c_{ld}$ to mimic the $c_{hd}$ and the verifiability of $c_{ld}$ by using a recent *verification-aware knowledge distillation technique* [21]. Originally, this method was developed to compress low-dimensional neural networks for better verifiability — and we extend it to approximate an HDC with LDCs using the supervised dataset $\mathcal{D}_{tr}$. Specifically, we implement knowledge distillation with *two-objective gradient descent*, which aims to optimize the MSE loss function $L_{mse}$ and Lipschitz constant loss function $L_{lip}$. First, it computes the directions of two gradients with respect to the $c_{ld}$ parameters $\theta$:

$$d_{L_{mse}} = \frac{\partial L_{mse}}{\partial \theta}, \quad d_{L_{lip}} = \frac{\partial L_{lip}}{\partial \theta} \tag{2}$$

The two-objective descent operates case-by-case to optimize at least one objective as long as possible. If $d_{L_{mse}} \cdot d_{L_{lip}} > 0$, the objectives can be optimized simultaneously by following the direction of the angular bisector of the two gradients. If $d_{L_{mse}} \cdot d_{L_{lip}} < 0$, then it is impossible to improve both objectives. Then, weights are updated along the vector of $d_{L_{mse}}$ (the higher priority) projected onto the hyperplane perpendicular to $d_{L_{lip}}$. The thresholds for MSE and Lipschitz constants in our system $M_{ld}$ are denoted as $\epsilon$ and $\lambda$ respectively. The stopping condition is met when both loss functions are below their thresholds or the training time exceeds the limit. Later on, Step 1 will be referred to with function TRAINLDC, and our way of tuning $\epsilon$ and $\lambda$ will be described later in Step 5.

**Step 2: Reachability analysis in low dimensions**

After training LDCs $\{c_{ld}^1, ..., c_{ld}^m\}$, we construct overapproximate reachable tubes for each. We perform reachability analysis for systems $M_{ld}^1, \ldots, M_{ld}^m$ with the respective controllers and the initial set $S_0$ specified in the original verification problem. This will result in a set of reachable tubes $\mathsf{rt}_{M_{ld}^1}(S_0, T), \ldots, \mathsf{rt}_{M_{ld}^m}(S_0, T)$.

To carry out the reachability analysis, we use the $POLAR^6$ *toolbox*, version of December 2022 [27,62] — an implementation that computes univariate Bernstein polynomials to overapproximate activation functions in $c_{ld}$, and then selectively uses Taylor or Bernstein polynomials for tight overapproximation of $c_{ld}$. For dynamics reachability, which alternates with neural-network overapproximation, the POLAR toolbox relies on the mature Flow* tool with Taylor model approximations [9]. The latest experimental results [62] show that POLAR outperforms other neural-network verification tools in terms of both computational efficiency and tightness. The verification details are formalized in Algorithm 2 in Step 5.

**Step 3a: Defining discrepancy bounds**

The LDC reachable tubes from Step 2 cannot be used directly to obtain HDC guarantees because of the discrepancy between LDC and HDC behaviors, which inevitably arises when compressing a higher-parameter neural network [24]. Therefore, we will quantify the difference between LDCs and HDCs using *discrepancy functions*, inspired by the prior work on testing hybrid systems [19, 20, 44]. We introduce and investigate two types of discrepancy functions in our setting:

**1. Trajectory-based discrepancy** $\beta$ considers the difference between the HDC and LDC *trajectories* starting from a *matched* state-image pair $(s, z)$, i.e., $z = g(s)$. It is defined as the least upper bound on the maximum L1 distance between two trajectories, i.e., $\|\tau_{hd}(s_0, T) - \tau_{ld}(s_0, T)\|_1$, over time $T$ for all initial states $s_0$ within the initial set $S_0$. Therefore, each initial set $S_0$ gives rise to its trajectory-based discrepancy $\beta(S_0)$.

**2. Action-based discrepancy** $\gamma$ considers the difference between LDC and HDC *actions* on a *matched* state-image pair $(s, z)$, i.e., $z = g(s)$. Similarly to the above, it is defined as the least upper bound on the difference between control actions over time horizon $T$ starting from any initial state $s_0$ within the initial set $S_0$. Note that the control difference, $\|c_{hd}(g(s_{hd}^t)) - c_{ld}(s_{ld}^t)\|_1$, is considered at each time step, where the $s$ is each state in the two trajectories.

The formal Definitions 12 and 13 can be found in the Appendix.

**Step 3b: Computing statistical discrepancy bounds**

Unfortunately, obtaining the true discrepancies is impractical: it would require solving optimization/feasibility problems in high-dimensional image spaces. Instead, we calculate the statistical upper bounds for these discrepancies via *conformal prediction*, which is a distribution-free statistical technique to provide

---

[6] https://github.com/ChaoHuang2018/POLAR_Tool

probabilistically valid uncertainty regions for complex prediction models — without strong assumptions about these models or their error distributions [55].

Below we briefly summarize basic conformal prediction. Consider $k+1$ independent and identically distributed random variables $\Delta, \Delta^1, ..., \Delta^k$, also known as *non-conformity scores*. Conformal prediction computes an uncertainty region for $\Delta$ via a function $\bar{\Delta} : \mathbb{R}^k \to \mathbb{R}$ from the other $k$ values. Given a failure probability $\alpha \in (0, 1)$, conformal prediction provides an uncertainty bound on $\bar{\Delta}$ such that $\Pr(\Delta \leq \bar{\Delta}) \geq 1 - \alpha$. This is performed with a surprisingly simple quantile argument, where the uncertainty bound $\bar{\Delta}$ is calculated as the $(1-\alpha)$-th quantile of the empirical distribution over the values of $\Delta^1, \Delta^2, ..., \Delta^k$, and $\infty$. The guarantee is formalized in the lemma below, and for details see a popular tutorial [48].

**Lemma 1.** *(Lemma 1 in [22]) Let $\Delta, \Delta^1, \Delta^2, ..., \Delta^k$ be $k$+1 independent identically distributed real-valued random variables. Without loss of generality, let $\Delta, \Delta^1, \Delta^2, ..., \Delta^k$ be stored in non-decreasing order and define $\Delta^{k+1} := \infty$. For $\alpha \in (0, 1)$, it holds that $\Pr(\Delta \leq \bar{\Delta}) \geq 1 - \alpha$ where $\bar{\Delta} := \Delta^{(r)}$, which is the $r$-ranked variable with $r = \lceil (k + 1)(1 - \alpha) \rceil$, and $\lceil . \rceil$ is the ceiling function.*

Leveraging conformal prediction, we define the statistical versions of our discrepancy functions. For the trajectory-based one, we define the non-conformity as the maximum L1 distance between states at the same time in two matched trajectories $\tau_{ld}(s_0, T)$ and $\tau_{hd}(s_0, T)$ starting from a random state $s_0 \sim D_0$ sampled independently and identically distributed (i.i.d.) from a given distribution $D_0$ over the initial region $S_0$, similar to recent works [11, 44]. This leads to a trajectory dataset $\mathcal{D}_{tb}$, from which $k$ non-conformity scores are calculated.

**Definition 3 (Statistical trajectory-based discrepancy).** *Given distribution $D_0$ over $S_0$, confidence $\alpha \in (0, 1)$, and state functions $\varphi_{hd}(s, t)$ and $\varphi_{ld}(s, t)$ for systems $M_{hd}$ and $M_{ld}$, a statistical trajectory-based discrepancy $\bar{\beta}(D_0)$ is an $\alpha$-confident upper bound on the max trajectory distance starting from $s_0 \sim D_0$:*

$$\Pr_{s_0 \sim D_0} \left[ \max_{t=0..T} \|\varphi_{hd}(s_0, t) - \varphi_{ld}(s_0, t)\|_1 \leq \bar{\beta}(D_0) \right] \geq 1 - \alpha$$

To obtain this bound $\bar{\beta}(D_0)$, we leverage conformal prediction as follows. Dataset $\mathcal{D}_{tb}$ contains i.i.d. samples $s_1, s_2, ..., s_k$ from our chosen distribution $D_0$. In practice, we choose the uniform distribution, namely $s \sim \text{Uniform}(S)$, because we value the safety of each state equally. We compute the corresponding non-conformity scores $\delta^1, \delta^2, ..., \delta^k, \delta^{k+1}$ as the maximum L1 distances between the same-time states in the two trajectories over all times $t \in [0..T]$:

$$\delta^i = \max_{t=0..T} \|\varphi_{hd}(s_i, t) - \varphi_{ld}(s_i, t)\|_1 \text{ for } i = 1 \ldots k; \text{ and } \delta^{k+1} = \infty$$

We sort the scores in the increasing order and set $\bar{\beta}(S_0)$ to the $r$-th quantile:

$$\bar{\beta}(D_0) := \delta^{(r)} \text{ with } r = \lceil (k+1)(1-\alpha) \rceil \tag{3}$$

We follow a similar procedure for the statistical action-based discrepancy, except that now the non-conformity scores are defined as the maximum differences between actions at the same time in two paired trajectories.

**Definition 4 (Statistical action-based discrepancy).** *Given confidence* $\alpha \in (0, 1)$, *distribution* $D_0$ *over* $S_0$, *and systems* $M_{ld}$ *and* $M_{hd}$, *a statistical action-based discrepancy* $\bar{\gamma}(D_0)$ *is an* $\alpha$-*confident upper bound on maximum action discrepancy in two trajectories starting from* $s_0 \sim D_0$:

$$\mathrm{Pr}_{D(S_0)} \left[ \max_{t=0..T} \| c_{hd}\big(g(\varphi_{hd}(s_0, t))\big) - c_{ld}\big(\varphi_{ld}(s_0, t)\big) \|_1 \leq \bar{\gamma}(D_0) \right] \geq 1 - \alpha$$

To implement this statistical action-based discrepancy function, we sample initial states $s_1, s_2, ..., s_k$ from a given set $S_0$ following the distribution $D_0$ (in practice, uniform) and obtain the corresponding low-dimensional trajectories. Then we generate with $g$ the corresponding images matched to each state in each trajectory — and these pairs form our action-based dataset $\mathcal{D}_{ab}$. The corresponding nonconformity scores $\delta^1, \delta^2, ..., \delta^k, \delta^{k+1}$ are maximum action differences:

$$\delta^i = \max_{t=0..T} \| c_{hd}(g(\varphi_{hd}(s_0, t))) - c_{ld}(\varphi_{ld}(s_0, t)) \|_1 \text{ for } i = 1 \dots k; \delta^{k+1} = \infty.$$

Then we sort these non-conformity scores in the non-decreasing order and determine the statistical bound for the action-based discrepancy as:

$$\bar{\gamma}(D_0) \coloneqq \delta^{(r)} \text{ with } r = \lceil (k+1)(1-\alpha) \rceil \tag{4}$$

**Step 4: Inflating reachability with discrepancies**

This step combines low-dimensional reachable tubes (Step 2) with statistical discrepancies (Step 3b) to provide a safety guarantee on the high-dimensional system. Thus, we inflate the original LDC reach tubes with either trajectory or action discrepancy to contain the (unknown) true HDC tube with chance $1 - \alpha$.
**Trajectory-based inflation.** The trajectory-based approach inflates the LDC reachable set starting in region $S_0$ with the statistical trajectory-based discrepancy $\bar{\beta}(D_0)$. Since the final reachable tube for a given initial set of $c_{ld}$ is represented as a sequence of discrete state polytopes calculated by concretizing the Taylor model with interval arithmetic on the initial set [27], we inflate these polygons by adding $\bar{\beta}(D_0)$ to their boundaries.

**Definition 5 (Trajectory-inflated reachable set).** *Given a distribution* $D_0$ *over initial set* $S_0$ *that is controlled by LDC* $c_{ld}$, *reachable set* $\mathsf{rs}(S_0, t)$, *and its trajectory discrepancy* $\bar{\beta}(D_0)$, *a* trajectory-inflated reachable set *is defined as:*

$$\mathsf{irs}(S_0, t, \bar{\beta}(D_0)) = \big\{ s \in S \mid \exists s' \in \mathsf{rs}(S_0, t) \cdot \| s - s' \|_1 \leq \bar{\beta}(D_0) \big\}$$

**Definition 6 (Trajectory-inflated reachable tube).** *Given a distribution* $D_0$ *over initial set* $S_0$ *that is controlled by LDC* $c_{ld}$, *a reachable tube* $\mathsf{rt}(S_0, t) = \big[ S_0, \mathsf{rs}(S_0, 1), \dots, \mathsf{rs}(S_0, T) \big]$ *over time horizon* $T$, *and its trajectory discrepancy* $\bar{\beta}(D_0)$ *over the initial set* $S_0$, *a* trajectory-inflated reachable tube $\mathsf{irt}(S_0, \bar{\beta}(D_0))$ *is defined as:*

$$\mathsf{irt}(S_0, \bar{\beta}(D_0)) = \big[ \mathsf{irs}(S_0, 0, \bar{\beta}(D_0)), \mathsf{irs}(S_0, 1, \bar{\beta}(D_0)), \dots, \mathsf{irs}(S_0, T, \bar{\beta}(D_0)) \big].$$

Based on Defs 5 and 6, we establish Theorem 1 that the trajectory-inflated LDC reachable tube contains the HDC reachable tube with at least $1 - \alpha$ probability.

**Theorem 1 (Confident trajectory-based overapproximation).** *Consider distribution $D_0$ over initial set $S_0$, confidence $\alpha$, a high-dimensional system $M_{hd}$, approximated with a low-dimensional system controlled by $c_{ld}$ with an $\alpha$-confident statistical trajectory-based discrepancy function $\bar{\beta}(S_0)$. Then the trajectory-inflated low-dimensional tube $\mathsf{irt}_{M_{ld}}(S_0, \bar{\beta}(D_0))$ contains the high-dimensional reachable tube $\mathsf{rt}_{M_{hd}}(S_0)$ with probability $1 - \alpha$:*

$$\Pr_{D_0}\left[\mathsf{rt}_{M_{hd}}(S_0) \subseteq \mathsf{irt}_{M_{ld}}(S_0, \bar{\beta}(S_0))\right] \geq 1 - \alpha$$

*Proof.* All the proofs are found in the Appendix.

Defs. 5 and 6 and Thm. 1 describe inflation and guarantees with a *single LDC*. However, one LDC usually cannot mimic the behavior of the HDC accurately. Therefore, we train several LDCs $\{c_{ld}^1, c_{ld}^2, \ldots, c_{ld}^m\}$, one for each subregion of initial set $\{S_1, S_2, \ldots, S_m\}$ with respective distributions $D_0 = \{D_1, D_2, \ldots, D_m\}$. Subsequently, the trajectory-inflated tube with multiple LDCs can be represented as a union of all the single trajectory-inflated tube $\mathsf{irt}(S_0, \bar{\beta}(D_0)) := \bigcup_{i=1}^{m} \mathsf{irt}(S_i, \bar{\beta}(D_i))$. The definitions and the multi-LDC version of Thm. 1 are in the Appendix under Defs. 14, 15 and Thm. 4.

**Action-based inflation.** Action-based inflation is less direct than with trajectories: we inflate the neural network's *output set* that is represented by a *Taylor model* $\mathrm{TM}(p(S_0), I)$ [27], where $p(S_0)$ is a polynomial representing order-$k$ Taylor series expansion of the $c_{ld}$ activation functions in region $S_0$, and the remainder interval $I$ ensures that Taylor model overapproximates the neural network's output. In this context, we widen the bounds of the remainder interval $I$ in the last layer of the $c_{ld}$ by our statistical action-based discrepancy $\bar{\gamma}(D_0)$, ensuring that the potential outputs of $c_{hd}$ are contained in the resulting Taylor model.

**Definition 7 (Action-inflated reachable set).** *Given distribution $D_0$ over set $S_0$ that is controlled by LDC $c_{ld}$, statistical action-based discrepancy $\bar{\gamma}(D_0)$, and low-dimensional control bounds $[u_{min}(t), u_{max}(t)] \supseteq c_{ld}(S_0)$, the action-inflated reachable set contains states reachable by inflating the action bounds:*

$$\mathsf{irs}(S_0, \bar{\gamma}(D_0)) = \left\{ f(s, u) \mid s \in S_0, u \in \left[u_{min}(t) - \bar{\gamma}(D_0), u_{max}(t) + \bar{\gamma}(D_0)\right] \right\}$$

**Definition 8 (Action-inflated reachable tube).** *Given an distribution $D_0$ over initial set $S_0$ that is controlled by LDC $c_{ld}$, dynamics $f$, time horizon $T$, and action-based discrepancy functions $\bar{\gamma}(D_0)$, the action-inflated reachable tube is a recursive sequence of inflated action-based reachable sets:*

$$\mathsf{irt}(S_0, \bar{\gamma}(D_0)) = \left[S_0, \mathsf{irs}_1(S_0, \bar{\gamma}(D_0)), \mathsf{irs}_2(\mathsf{irs}_1, \bar{\gamma}(D_0)), \ldots, \mathsf{irs}_T(\mathsf{irs}_{T-1}, \bar{\gamma}(D_0))\right].$$

Based on Defs. 7 and 8, we put forward Thm. 2 below for the lower probability bound of the action-inflated LDC tube containing the true HDC tube.

---

**Algorithm 1** Iterative LDC training for the action-based approach

---

**function** IterativeTrainingAB(HDC $c_{hd}$, image generator $g$, sample count $N$,
initial state space $S_0$, confidence $\alpha$, discrepancy thresh. $\xi$, time steps $T$, goal set $G$)
    $\lambda, \epsilon \leftarrow$ initial values
    $\mathbf{S} \leftarrow$ initial gridding of $S_0 : S_1, S_2, \ldots$
    **while** Computing resources last **do**
        **for** $i = 1$ to $|\mathbf{S}|$ **do**
            $c_{ld}^i \leftarrow$ TrainLDC($c_{hd}$, g, $S_i$, $\lambda$, $\epsilon$)
            $\delta^i \leftarrow$ ComputeActionDiscr($c_{ld}^i$, $c_{hd}$, g, $S_i$, $\alpha$, $N$)
            **if** $\delta^i > \xi$ **then**
                $\epsilon \leftarrow \epsilon/2$                           ▷ Reduce MSE threshold
            **end if**
        **end for**
        **if** $\hat{\delta} > \xi$ in some sub-region $\hat{\mathbf{S}} \subseteq \mathbf{S}$ **then**         ▷ Too much discrepancy
            $\mathbf{S}' \leftarrow \mathbf{S}$ with refined re-gridding of $\hat{\mathbf{S}}$
        **end if**
        **if** $\hat{\delta} \leq \xi \wedge \mathsf{rs}_{M_{ld}}(\hat{\mathbf{S}}, T) \not\subseteq G$ in some sub-region $\hat{\mathbf{S}} \subseteq \mathbf{S}$ **then**
            $\lambda \leftarrow \lambda/2$ and keep the same $\epsilon$ in $\hat{\mathbf{S}}$      ▷ Reduce Lipschitz threshold
        **end if**
        $\mathbf{S} \leftarrow \mathbf{S}'$                                    ▷ Use the updated grid
    **end while**
    $\bar{\gamma} \leftarrow \delta^1, \delta^2, \ldots$
    **return** $c_{ld}^1, c_{ld}^2, \ldots, \bar{\gamma}$
**end function**

---

**Theorem 2 (Confident action-based overapproximation).** *Consider distribution $D_0$ over initial set $S_0$, high-dimensional system $M_{hd}$ with controller $c_{hd}$, approximated by low-dimensional system $M_{ld}$ controlled by $c_{ld}$ with $\alpha$-confident statistical action-based discrepancies $\bar{\gamma}(S_0)$. Then the action-inflated low-dimensional tube $\mathsf{irt}_{M_{ld}}(S_0, \bar{\gamma}(S_0))$ contains the high-dimensional tube $\mathsf{rt}_{M_{hd}}(S_0)$ with probability $1 - \alpha$:*

$$\mathrm{Pr}_{D_0}\left[\mathsf{rt}_{M_{hd}}(S_0) \subseteq \mathsf{irt}_{M_{ld}}(S_0, \bar{\gamma}(S_0))\right] \geq 1 - \alpha$$

Defs. 7 and 8 describe inflation with a *single LDC*, which we extend to multiple LDCs by taking the union of all the LDCs' inflated tubes. Given a partitioned initial set $S_0 = \{S_1, ..., S_m\}$ with respective controllers $\{c_{ld}^1, \ldots, c_{ld}^m\}$ and distributions $D_0 = \{D_1, ..., D_m\}$, the multiple LDCs action-inflated reachable tube is $\mathsf{irt}(S_0, \bar{\gamma}(D_0)) := \bigcup_{i=1}^m \mathsf{irt}(S_i, \bar{\gamma}(D_i))$. As it turns out, this reachable tube also contains the HDC tube with at least $1-\alpha$ chance (See Thm. 5 in the Appendix).

**Step 5: Iterative retraining and re-gridding**

Once the inflated reachable tubes are obtained in Step 4, we focus on the regions of the current subset of the initial set where HDC simulations succeed — yet safety verification fails. This can happen for two reasons: (i) overly high overapproximation error in the LDC reachability, or (ii) overly high conformal discrepancy bounds from $\bar{\beta}$ or $\bar{\gamma}$. We handle these issues as follows.

---

**Algorithm 2** End-to-end reachability verification of an HDC

---

**function** ENDTOENDVERIFICATION(HDC $c_{hd}$, generator $g$, sample count $N$, state
space $S$, initial set $S_0$, confidence $\alpha$, discrepancy threshold $\xi$, time horizon $T$, goal
set $G$, approach selection $J \in \{$ trajectory-based, action-based $\})$

    **if** $J =$ trajectory-based **then**

        $c_{ld}^1 \ldots c_{ld}^n, \bar{\beta} \leftarrow$ ITERATIVETRAININGTB$(c_{hd}, g, N, S_0, \alpha, \xi, T)$

        $X \leftarrow \bar{\beta}$                            ▷ Store the trajectory discrepancies

    **else**

        $c_{ld}^1 \ldots c_{ld}^n, \bar{\gamma} \leftarrow$ ITERATIVETRAININGAB$(c_{hd}, g, N, S, \alpha, \xi, T, G)$

        $X \leftarrow \bar{\gamma}$                         ▷ Store the action discrepancies

    **end if**

    $\mathbf{S_{ver}} \leftarrow$ split $S_0$ into regions: $S_0^1, S_0^2, \ldots$      ▷ Gridding for parallel verification

    $S_{safe}, S_{unsafe} \leftarrow \emptyset$                    ▷ Initialize safe and unsafe regions

    **for** $j = 1$ to $|\mathbf{S_{ver}}|$ **do**

        $\mathsf{irs}(S_0^j, X, T) \leftarrow$ REACH$(c_{ld}^1, \ldots, c_{ld}^n, S_0^j, X, T)$

        **if** $\mathsf{irs}(S_0^j, X, T) \subseteq G$ **then**

            $S_{safe} \leftarrow S_{safe} \cup S_0^j$

        **else**

            $S_{unsafe} \leftarrow S_{unsafe} \cup S_0^j$

        **end if**

    **end for**

    **return** $S_{safe}, S_{unsafe}$

**end function**

---

**Reducing reachability overapproximation error.** We lower the threshold for the Lipschitz constant $\lambda$ to retrain the respective LDCs in Step 1. In our experience, this almost always reduces the overapproximation in the LDC analysis and makes low-dimensional reachable tubes tighter — but may result in higher statistical discrepancy bounds, addressed below.

**Reducing conformal discrepancy bounds.** When these bounds are high, our LDC is not sufficient to imitate the HDC performance in a certain region of the state space. Given a desired discrepancy bound $\xi$, when this bound $\xi$ is exceeded in a state-space region, we split it into subregions by taking its midpoints in each dimension, leading to an updated state-space grid $\mathbf{S}'$. Then in each sub-region, we retrain an LDC as per Step 1 with a reduced acceptable MSE threshold $\epsilon$ and re-compute its bounds as per Step 3b. leading to tighter statistical overapproximations of HDC reachable tubes.

To summarize, Alg. 1 shows our iterative training procedure for the action-based approach (its trajectory-based counterpart proceeds analogously, except for computing the discrepancies over trajectories, see Alg. 6 in the Appendix).

Combining all the five steps together, we present Alg. 2 that displays our end-to-end verification of a given HDC with either trajectory-based or action-based discrepancies. The LDCs and their discrepancies are input into the reachability analysis, implemented with the function REACH, to calculate the inflated reachable tubes (using the POLAR toolbox in practice). Our end-to-end algorithm

guarantees that an affirmative answer to our verification problem is correct with at least $1 - \alpha$ probability, as per Thm. 3.

**Theorem 3 (Confident guarantee of HDC safety).** *Consider a partitioned initial set grid $S_0 = \{S_1, \ldots, S_m\}$, a set of corresponding distributions $\{D_1, \ldots D_m\}$, a high-dimensional system $M_{hd}$ with controller $c_{hd}$, and a set of low-dimensional systems $M_{ld}^1, \ldots, M_{ld}^m$ with respective controllers $c_{ld}^1, \ldots, c_{ld}^n$ that approximate $c_{hd}$ with either an $\alpha$-confident trajectory discrepancy or action discrepancy, the probability that HDC safe set $S_{safe}$ calculated by Alg. 2 with either discrepancy belongs to ground truth safe set $S_{safe}^*$ is at least $(1 - \alpha)$:*

$$\Pr_{D_1 \ldots D_m} \left[ S_{safe} \subseteq S_{safe}^* \right] \geq (1 - \alpha)$$

## 4   Experimental Evaluation

**Benchmark systems and controllers.** We evaluate our approach on three benchmarks from OpenAI Gym [7]: two two-dimensional case studies — an *inverted pendulum* (IP) with angle $\theta$ and angular velocity $\dot{\theta}$; a *mountain car* (MC) with position $x$ and velocity $v$, and a four-dimensional case study — a *cart pole* (CP) with cart position $x$, cart velocity $v$, angle $\theta$, and angular velocity $\dot{\theta}$. Our selection of case studies is limited because of the engineering challenge of setting up *both* vision-based control and low-dimensional verification for the same system. Our continuous-action, convolutional HDCs $c_{hd}$ for these systems were trained with deep deterministic policy gradient (DDPG) [36]. To imitate the performance of $c_{hd}$, we train simpler feedforward neural networks $c_{ld}$ with only low-dimensional state inputs. See the Appendix for their architecture and dynamics, and our code can be accessed from GitHub [7]

**Experimental procedure.** Our verification's goal is to check whether the system will stay inside the specified goal set $G$ after $T$ time steps (e.g., the mountain car's position must stay within the target set $[0.45, \infty]$ after 60 steps). The verification returns "safe" if the inflated reachable set for $t = T$ lies entirely in $G$ — and "unsafe" otherwise. The details are found in the Appendix.

For both approaches, we calculate the discrepancies in 0.25-sized state squares within the initial set in IP, hence creating $8 \times 8 = 64$ regions (MC has $8 \times 9 = 72$ regions; CP has $5 \times 5 \times 5 \times 5 = 625$). In each, we sample 60 trajectories to compute both trajectory-based discrepancies $\bar{\beta}$ and action-based discrepancies $\bar{\gamma}$ because it is a relatively small sample count that avoids the highest non-conformity score or the infinity as the conformal bound. We also implement a *pure conformal prediction baseline* and, for a fair comparison, give it the same data/regions. This results in 3840 sampled trajectories in IP, 4320 in MC, and 76800 for CP.

We use closed-loop simulation to obtain the (approximate) ground truth (GT) of safety. For IP and CP, we grid the initial set into squares with an interval of 0.01. For MC, we grid the initial set with the position step 0.01 and velocity step 0.001. Within each grid cell, we uniformly sample 10 initial states

---

[7] https://github.com/yuanggeng/Bridging-dimensions

and simulate a trajectory from each. If all 10 trajectories end in the goal set $G$, we mark this cell as "truly safe", otherwise "truly unsafe". In IP, the truly safe-to-unsafe cell ratio is 0.56, 0.78 in MC, and 0.58 in CP. The verification process uses the same grid cells as its initial state regions, leading to 40k low-dimensional verification runs for IP, 14k for MC, and 50k for CP. The trajectory-based verification time for IP, MC, and CP are 6.2, 5.8, and 6.4 hours respectively; the action-based verification takes 6.3, 6.1, and 6.6 hours respectively.

**Success metrics.** We evaluate verification as a binary classifier of the GT safety, with "safe" being the positive class and "unsafe" being the negative. Our evaluation metrics are the (i) *true positive rate* (TPR, a.k.a. sensitivity and recall), indicating the fraction of truly safe regions that were successfully verified; (ii) *true negative rate* (TNR, a.k.a. specificity), indicating the fraction of truly unsafe regions that failed verification; (iii) *precision*, indicating the fraction of safe verification verdicts that are truly safe (which is essential for safety-critical systems and controlled by rate $\alpha$ as per Thm. 3); and (iv) *F1 score*, which is a harmonic mean of precision and recall to provide a class-balanced assessment of predictions.

**Table 1.** Verification performance ($M = 4$ for IP and CP, $M = 10$ for MC).

| Benchmark | Metrics | Pure conformal prediction | Trajectory-based approach | | Action-based approach | |
|---|---|---|---|---|---|---|
| | | HDC | 1 LDC | $M$ LDCs | 1 LDC | $M$ LDCs |
| Inverted Pendulum (IP) | True positive rate | 0.6564 | 0.4662 | **0.7938** | 0.0603 | 0.4050 |
| | True negative rate | **0.9999** | 0.9976 | **0.9995** | **1.0000** | **0.9999** |
| | Precision | **0.9998** | 0.9880 | 0.9985 | **1.0000** | 0.9997 |
| | F1-score | 0.7925 | 0.6335 | **0.8844** | 0.1137 | 0.5765 |
| Mountain Car (MC) | True positive rate | 0.4686 | **0.7220** | 0.7207 | 0.1050 | 0.2659 |
| | True negative rate | **0.9967** | 0.9693 | 0.9872 | 0.9964 | **1.0000** |
| | Precision | 0.9916 | 0.9621 | 0.9793 | **0.9999** | **1.0000** |
| | F1-score | 0.6364 | 0.8249 | **0.8303** | 0.1900 | 0.4201 |
| Cartpole (CP) | True positive rate | 0.6697 | 0.7225 | **0.7450** | 0.6554 | 0.7238 |
| | True negative rate | **1.0000** | 0.9998 | **1.0000** | **1.0000** | **1.0000** |
| | Precision | **1.0000** | 0.9999 | **1.0000** | **1.0000** | **1.0000** |
| | F1-score | 0.8022 | 0.8389 | **0.8539** | 0.7918 | 0.8398 |

**Table 2.** Verification performance for multiple LDCs with zero-mean Gaussian noise added to true state before image generator $g$.

| Benchmark | Metrics | Trajectory-based method | | Action-based method | |
|---|---|---|---|---|---|
| Inverted Pendulum (IP) | STD of $\theta, \dot{\theta}$ noise | 0.01 | 0.1 | 0.01 | 0.1 |
| | True positive rate | **0.6732** | 0.5272 | 0.3675 | 0.1924 |
| | True negative rate | **1.0000** | **1.0000** | 0.9999 | **1.0000** |
| | Precision | **1.0000** | **1.0000** | 0.9997 | 0.9997 |
| | F1-score | **0.8046** | 0.6904 | 0.5374 | 0.3228 |
| Mountain Car (MC) | STD of $x$ noise | 0.01 | 0.1 | 0.01 | 0.1 |
| | STD of $v$ noise | 0.0001 | 0.003 | 0.0001 | 0.003 |
| | True positive rate | **0.6797** | 0.4189 | 0.1558 | 0.0658 |
| | True negative rate | 0.9878 | 0.9889 | **1.0000** | **1.0000** |
| | Precision | 0.9790 | 0.9753 | **1.0000** | **1.0000** |
| | F1-score | **0.8023** | 0.5861 | 0.2696 | 0.1235 |
| Cartpole (CP) | STD of $x,v,\theta,\dot{\theta}$ noise | 0.03 | 0.1 | 0.03 | 0.1 |
| | True positive rate | **0.7108** | 0.6253 | 0.6724 | 0.6040 |
| | True negative rate | **0.9998** | **0.9998** | 0.9996 | 0.9996 |
| | Precision | **0.9995** | **0.9995** | 0.9990 | 0.9989 |
| | F1-score | **0.8308** | 0.7692 | 0.8038 | 0.7528 |

**Verification results.** The quantitative results of the three case studies are summarized in Tab. 1. Confidence $\alpha$ is set to 0.05 for all methods, which sets the

minimum precision to 0.95, satisfied by all the approaches. The pure conformal prediction baseline shows high precision and TNR, but loses in TPR to our approaches — thus being able to correctly verify a significantly smaller region of the state space. When it comes to well-balanced safety prediction in practice, F1 score shows that our trajectory-based approach outperforms the other two.

Across all case studies, the baseline is significantly more conservative than the requested 95% precision. While this can be an advantage in safety-critical settings, excessive conservatism can also hamper adoption, so the approach should be sensitive to the desired confidence — which our trajectory-based approach demonstrates in the mountain car case study (see Precision in Tab. 1).

Across all case studies, the multi-LDC approaches always match or outperform the one-LDC approaches. This result demonstrates the utility of modularizing the HDC approximation problem. Also, our single-LDC action-based approach successfully verifies relatively few regions, leading to its low TPR. That is because unlike in the case of trajectory discrepancies, only one LDC cannot provide tight statistical upper bounds for control actions, causing large overapproximation in the inflated reachable sets, resulting in false negatives.

**Sensitivity to noisy images.** Despite adding Gaussian noise to generator $g$, our approaches perform similarly to noise-free $g$ when under low noise variance as per Tab. 2, thus showing some robustness. However, we saw a significant decline in the verification coverage (TPR, but not the TNR and $\alpha$-guaranteed precision) under substantial noise variance (up to 0.5, not shown in Tab. 2).

**Limitations.** Our approach relies on statistical inference based on i.i.d. sampling from a fixed distribution, which downgrades the exhaustive guarantees of formal verification. However, it may be possible to exhaustively bridge this gap with neural-network conformance analysis based on satisfiability solving [41]. We also envision relaxing the i.i.d. assumption with time-series conformal prediction [3, 58], as well as uncertainty-guided gridding [37] to reduce our discrepancy bounds.

## 5   Related work

**Low-dimensional verification of closed-loop systems.** Neural-network controlled systems have been used widely [42, 46, 52], which has highlighted the challenges of verifying their correctness within closed-loop systems. Since it's impossible to calculate all the exact states, especially in non-linear systems, current approaches primarily focus on how to make tight overapproximate reachable sets [2, 9, 10]. For sigmoid-based NNCS, Verisig [30] toolbox can transform the neural-network controlled system into a hybrid system, which can be verified by other tools like flow*. NNV [54] performs overapproximation analysis by combining star sets [38, 53] for feed-forward neural networks with zonotopes for non-linear plant dynamics in CORA [2]. POLAR [27] overcame the challenges of non-differentiable activation functions by combining the Bernstein-Bézier Form [28] and the symbolic remainder. This method achieves state-of-the-art performance in both the tightness of reachable tubes and computation times. Another type of verification called *Hamilton-Jacobi* (HJ) reachability [4], is inspired by optimal control. The DeepReach [5] technique can solve the verification problem

with tens of dimensions by leveraging a deep neural network to represent the value function in the HJ reachability analysis. Nonetheless, such methods remain ill-suited for handling inputs with hundreds or thousands of dimensions.

These verification tools cannot deal with complicated neural network controllers. Therefore, an alternative approach is to simplify complex controllers into smaller, verifiable controllers by model reduction techniques [16,33], such as parameter pruning, compact convolution filters, and knowledge distillation [25]. **Abstractions of perception models.** Given the challenge of verifying the image-based closed-loop systems directly, many methods construct abstractions of the perception model to map the relationship between the image and the states for verification [43]. One abstraction approach [31] employs the generative model, especially Generative Adversarial Network (GAN), mapping states to images. The generated images will be put into the controller in the verification phase. Hence, the accuracy of the verification results depends on the quality of the image produced by the generative model. Other researchers [26] construct the exact mathematical formula mapping the real state into the simplified image [47], which can be verified in another neural network checker [32]. One limitation of exact modeling is the effort to generalize for other systems or scenarios. For instance, their implementation may be specific to a proportional controller in the aircraft landing or lane-keeping scenarios, which may not be suitable for the more complicated image-based systems in other cases. **Statistical verification**. Statistical verification draws samples to determine the property satisfaction from a finite number of trajectories [1, 11, 34, 35]. One advantage of such algorithms is that they provide assurance for arbitrarily complex black-box systems, merely requiring the ability to simulate them [59, 60]. Conformal prediction [55], which has been a popular choice for distribution-free uncertainty quantification, has recently been used to provide probabilistic guarantees on the satisfaction of a given STL property [37, 45]. Purely statistical methods come at the price of drawing sufficient samples — and only obtaining the guarantees at some level of statistical confidence, which can be difficult to interpret in the context of a dynamical system. Our work restricts the use of sampling only to the most challenging aspects and leverages exhaustive verification for the rest of the system, thus reducing our reliance on statistical assurance.

## 6    Conclusion

This paper takes a significant step towards addressing the major challenge of verifying end-to-end controllers implemented with high-dimensional neural networks. Our insight is that the behavior of such neural networks can be effectively approximated by several low-dimensional neural networks operating over physically meaningful space. To balance approximation error and verifiability in our low-dimensional controllers, we harness the state-of-the-art knowledge distillation. To close the gap between low- and high-dimensional controllers, we apply conformal prediction and provide a statistical upper bound on their difference either in trajectories or actions. Finally, by inflating the reachable tubes with two discrepancy types, we establish a high-confidence reachability guarantee for high-dimensional controllers. Future work may further reduce the role of sampling.

# References

1. Agha, G., Palmskog, K.: A survey of statistical model checking. ACM Trans. Model. Comput. Simul. **28**(1) (jan 2018). https://doi.org/10.1145/3158668, `https://doi.org/10.1145/3158668`

2. Althoff, M. An introduction to CORA 2015. *Proc. Of The Workshop On Applied Verification For Continuous And Hybrid Systems.* pp. 120-151 (2015)

3. Auer, A., Gauch, M., Klotz, D., Hochreiter, S.: Conformal Prediction for Time Series with Modern Hopfield Networks (Mar 2023). https://doi.org/10.48550/arXiv.2303.12783, `http://arxiv.org/abs/2303.12783`, arXiv:2303.12783 [cs, stat]

4. Bansal, S., Chen, M., Herbert, S.L., Tomlin, C.J.: Hamilton-jacobi reachability: A brief overview and recent advances. 2017 IEEE 56th Annual Conference on Decision and Control (CDC) pp. 2242–2253 (2017), `https://api.semanticscholar.org/CorpusID:35768454`

5. Bansal, S., Tomlin, C.J.: Deepreach: A deep learning approach to high-dimensional reachability. In: 2021 IEEE International Conference on Robotics and Automation (ICRA). pp. 1817–1824. IEEE (2021)

6. Bassan, S., Katz, G.: Towards formal xai: Formally approximate minimal explanations of neural networks. In: International Conference on Tools and Algorithms for the Construction and Analysis of Systems. pp. 187–207. Springer (2023)

7. Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., Zaremba, W.: OpenAI Gym (Jun 2016). https://doi.org/10.48550/arXiv.1606.01540, `http://arxiv.org/abs/1606.01540`, arXiv:1606.01540 [cs]

8. Chakraborty, K., Bansal, S.: Discovering closed-loop failures of vision-based controllers via reachability analysis. IEEE Robotics and Automation Letters **8**(5), 2692–2699 (2023)

9. Chen, X., Ábrahám, E., Sankaranarayanan, S.: Flow*: An analyzer for non-linear hybrid systems. In: International Conference on Computer Aided Verification (2013)

10. Chen, X., Sankaranarayanan, S.: Reachability Analysis for Cyber-Physical Systems: Are We There Yet?, pp. 109–130 (05 2022). https://doi.org/10.1007/978-3-031-06773-0̇6

11. Cleaveland, M., Lee, I., Pappas, G.J., Lindemann, L.: Conformal prediction regions for time series using linear complementarity programming. arXiv preprint arXiv:2304.01075 (2023)

12. Cleaveland, M., Sokolsky, O., Lee, I., Ruchkin, I.: Conservative Safety Monitors of Stochastic Dynamical Systems. In: Proc. of the NASA Formal Methods Conference (May 2023)

13. Codevilla, F., Miiller, M., López, A., Koltun, V., Dosovitskiy, A.: End-to-End Driving Via Conditional Imitation Learning. In: 2018 IEEE International Conference on Robotics and Automation (ICRA). pp. 1–9. IEEE Press, Brisbane, Australia (May 2018). https://doi.org/10.1109/ICRA.2018.8460487, `https://doi.org/10.1109/ICRA.2018.8460487`

14. Cofer, D., Amundson, I., Sattigeri, R., Passi, A., Boggs, C., Smith, E., Gilham, L., Byun, T., Rayadurgam, S.: Run-time assurance for learning-based aircraft taxiing. pp. 1–9 (10 2020). https://doi.org/10.1109/DASC50938.2020.9256581

15. Combettes, P.L., Pesquet, J.C.: Lipschitz Certificates for Layered Network Structures Driven by Averaged Activation Operators. SIAM Journal on Mathematics

of Data Science **2**(2), 529–557 (Jan 2020). https://doi.org/10.1137/19M1272780, `https://epubs.siam.org/doi/10.1137/19M1272780`, publisher: Society for Industrial and Applied Mathematics

16. Deng, L., Li, G., Han, S., Shi, L., Xie, Y.: Model compression and hardware acceleration for neural networks: A comprehensive survey. Proceedings of the IEEE **108**(4), 485–532 (2020)

17. Dutta, S., Caprio, M., Lin, V., Cleaveland, M., Jang, K.J., Ruchkin, I., Sokolsky, O., Lee, I.: Distributionally Robust Statistical Verification with Imprecise Neural Networks (Aug 2023). https://doi.org/10.48550/arXiv.2308.14815, `http://arxiv.org/abs/2308.14815`, arXiv:2308.14815 [cs]

18. Dutta, S., Chen, X., Jha, S., Sankaranarayanan, S., Tiwari, A.: Sherlock - a tool for verification of neural network feedback systems: Demo abstract. HSCC '19, Association for Computing Machinery, New York, NY, USA (2019). https://doi.org/10.1145/3302504.3313351, `https://doi.org/10.1145/3302504.3313351`

19. Fan, C., Mitra, S.: Bounded Verification with On-the-Fly Discrepancy Computation. In: Finkbeiner, B., Pu, G., Zhang, L. (eds.) Automated Technology for Verification and Analysis. pp. 446–463. Lecture Notes in Computer Science, Springer International Publishing, Cham (2015). https://doi.org/10.1007/978-3-319-24953-7˙32

20. Fan, C., Qi, B., Mitra, S., Viswanathan, M.: DryVR: Data-Driven Verification and Compositional Reasoning for Automotive Systems. In: Computer Aided Verification. pp. 441–461. Lecture Notes in Computer Science, Springer International Publishing, Cham (2017). https://doi.org/10.1007/978-3-319-63387-9˙22

21. Fan, J., Huang, C., Li, W., Chen, X., Zhu, Q.: Towards verification-aware knowledge distillation for neural-network controlled systems: Invited paper. 2019 IEEE/ACM International Conference on Computer-Aided Design (ICCAD) pp. 1–8 (2019), `https://api.semanticscholar.org/CorpusID:209497572`

22. Fannjiang, C., Bates, S., Angelopoulos, A.N., Listgarten, J., Jordan, M.I.: Conformal prediction under feedback covariate shift for biomolecular design. Proceedings of the National Academy of Sciences **119**(43), e2204569119 (2022). https://doi.org/10.1073/pnas.2204569119, `https://www.pnas.org/doi/abs/10.1073/pnas.2204569119`

23. Fazlyab, M., Robey, A., Hassani, H., Morari, M., Pappas, G.: Efficient and Accurate Estimation of Lipschitz Constants for Deep Neural Networks. In: Advances in Neural Information Processing Systems. vol. 32. Curran Associates, Inc. (2019), `https://proceedings.neurips.cc/paper_files/paper/2019/hash/95e1533eb1b20a97777749fb94fdb944-Abstract.html`

24. Gou, J., Yu, B., Maybank, S.J., Tao, D.: Knowledge distillation: A survey. International Journal of Computer Vision **129**, 1789–1819 (2021)

25. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015)

26. Hsieh, C., Li, Y., Sun, D., Joshi, K., Misailovic, S., Mitra, S.: Verifying controllers with vision-based perception using safe approximate abstractions. IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems **41**(11), 4205–4216 (2022). https://doi.org/10.1109/TCAD.2022.3197508

27. Huang, C., Fan, J., Chen, X., Li, W., Zhu, Q.: Polar: A polynomial arithmetic framework for verifying neural-network controlled systems. In: International Symposium on Automated Technology for Verification and Analysis. pp. 414–430. Springer (2022)

28. Huang, C., Fan, J., Li, W., Chen, X., Zhu, Q.: Reachnn: Reachability analysis of neural-network controlled systems. ACM Transactions on Embedded Computing Systems (TECS) **18**(5s), 1–22 (2019)

29. Fazlyab, M., Robey, A., Hassani, H., Morari, M. & Pappas, G. Efficient and accurate estimation of lipschitz constants for deep neural networks. *Advances In Neural Information Processing Systems.* **32** (2019)

30. Ivanov, R., Weimer, J., Alur, R., Pappas, G.J., Lee, I.: Verisig: Verifying safety properties of hybrid systems with neural network controllers. p. 169–178. HSCC '19, Association for Computing Machinery, New York, NY, USA (2019). https://doi.org/10.1145/3302504.3311806, `https://doi.org/10.1145/3302504.3311806`

31. Katz, S.M., Corso, A.L., Strong, C.A., Kochenderfer, M.J.: Verification of image-based neural network controllers using generative models. Journal of Aerospace Information Systems **19**(9), 574–584 (2022)

32. Khedr, H., Ferlez, J., Shoukry, Y.: Peregrinn: Penalized-relaxation greedy neural network verifier. In: Computer Aided Verification: 33rd International Conference, CAV 2021, Virtual Event, July 20–23, 2021, Proceedings, Part I. p. 287–300 (2021). https://doi.org/10.1007/978-3-030-81685-8˙13, `https://doi.org/10.1007/978-3-030-81685-8_13`

33. Ladner, T., Althoff, M.: Specification-driven neural network reduction for scalable formal verification. arXiv preprint arXiv:2305.01932 (2023)

34. Larsen, K.G., Legay, A.: Statistical model checking: Past, present, and future. In: Margaria, T., Steffen, B. (eds.) Leveraging Applications of Formal Methods, Verification and Validation: Foundational Techniques. ISoLA 2016. Lecture Notes in Computer Science, vol. 9952. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-47166-2˙1, `https://doi.org/10.1007/978-3-319-47166-2_1`

35. Lew, T., Janson, L., Bonalli, R., Pavone, M.: A simple and efficient sampling-based algorithm for general reachability analysis, `http://arxiv.org/abs/2112.05745`

36. Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D.: Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971 (2015)

37. Lindemann, L., Qin, X., Deshmukh, J.V., Pappas, G.J.: Conformal prediction for stl runtime verification. In: Proceedings of the ACM/IEEE 14th International Conference on Cyber-Physical Systems (with CPS-IoT Week 2023). p. 142–153. ICCPS '23, Association for Computing Machinery, New York, NY, USA (2023). https://doi.org/10.1145/3576841.3585927, `https://doi.org/10.1145/3576841.3585927`

38. Lopez, D.M., Musau, P., Tran, H.D., Johnson, T.T.: Verification of closed-loop systems with neural network controllers. EPiC Series in Computing **61**, 201–210 (2019)

39. Luo, R., Zhao, S., Kuck, J., Ivanovic, B., Savarese, S., Schmerling, E., Pavone, M.: Sample-efficient safety assurances using conformal prediction. In: International Workshop on the Algorithmic Foundations of Robotics. pp. 149–169. Springer (2022)

40. Matsumoto, E., Saito, M., Kume, A., Tan, J.: End-to-End Learning of Object Grasp Poses in the Amazon Robotics Challenge. In: Advances on Robotic Item Picking: Applications in Warehousing & E-Commerce Fulfillment, pp. 63–72. Springer International Publishing, Cham (2020). https://doi.org/10.1007/978-3-030-35679-8˙6, `https://doi.org/10.1007/978-3-030-35679-8_6`

41. Mohammadinejad, S., Paulsen, B., Deshmukh, J.V., Wang, C.: DiffRNN: Differential Verification of Recurrent Neural Networks. In: Dima, C., Shirmohammadi, M. (eds.) Formal Modeling and Analysis of Timed Systems. pp. 117–134. Lecture Notes in Computer Science, Springer International Publishing, Cham (2021). https://doi.org/10.1007/978-3-030-85037-1_8

42. Pan, Y., Cheng, C.A., Saigol, K., Lee, K., Yan, X., Theodorou, E., Boots, B.: Agile autonomous driving using end-to-end deep imitation learning. arXiv preprint arXiv:1709.07174 (2017)

43. Păsăreanu, C.S., Mangal, R., Gopinath, D., Getir Yaman, S., Imrie, C., Calinescu, R., Yu, H.: Closed-loop analysis of vision-based autonomous systems: A case study. In: International Conference on Computer Aided Verification. pp. 289–303. Springer (2023)

44. Qin, X., Hashemi, N., Lindemann, L., Deshmukh, J.V.: Conformance testing for stochastic cyber-physical systems. In: Conference on Formal Methods in Computer-Aided Design–FMCAD 2023. p. 294 (2023)

45. Qin, X., Xia, Y., Zutshi, A., Fan, C., Deshmukh, J.V.: Statistical verification of cyber-physical systems using surrogate models and conformal inference. In: 2022 ACM/IEEE 13th International Conference on Cyber-Physical Systems (ICCPS). pp. 116–126 (2022). https://doi.org/10.1109/ICCPS54341.2022.00017

46. Ruchkin, I., Cleaveland, M., Ivanov, R., Lu, P., Carpenter, T., Sokolsky, O., Lee, I.: Confidence Composition for Monitors of Verification Assumptions. In: ACM/IEEE 13th Intl. Conf. on Cyber-Physical Systems (ICCPS). pp. 1–12 (May 2022). https://doi.org/10.1109/ICCPS54341.2022.00007

47. Santa Cruz, U., Shoukry, Y.: Nnlander-verif: A neural network formal verification framework for vision-based autonomous aircraft landing. Springer-Verlag, Berlin, Heidelberg (2022). https://doi.org/10.1007/978-3-031-06773-0_11, `https://doi.org/10.1007/978-3-031-06773-0_11`

48. Shafer, G., Vovk, V.: A Tutorial on Conformal Prediction. J. Mach. Learn. Res. **9**, 371–421 (Jun 2008), `http://dl.acm.org/citation.cfm?id=1390681.1390693`

49. Stocco, A., Nunes, P.J., D'Amorim, M., Tonella, P.: Thirdeye: Attention maps for safe autonomous driving systems. In: Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering. ASE '22, Association for Computing Machinery, New York, NY, USA (2023). https://doi.org/10.1145/3551349.3556968, `https://doi.org/10.1145/3551349.3556968`

50. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. In: International Conference on Learning Representations (2014)

51. Teeti, I., Khan, S., Shahbaz, A., Bradley, A., Cuzzolin, F.: Vision-based Intention and Trajectory Prediction in Autonomous Vehicles: A Survey. vol. 6, pp. 5630–5637 (Jul 2022). https://doi.org/10.24963/ijcai.2022/785, `https://www.ijcai.org/proceedings/2022/785`, iSSN: 1045-0823

52. Topcu, U., Bliss, N., Cooke, N., Cummings, M., Llorens, A., Shrobe, H., Zuck, L.: Assured Autonomy: Path Toward Living With Autonomous Systems We Can Trust (Oct 2020). https://doi.org/10.48550/arXiv.2010.14443, `http://arxiv.org/abs/2010.14443`, arXiv:2010.14443 [cs]

53. Tran, H.D., Manzanas Lopez, D., Musau, P., Yang, X., Nguyen, L.V., Xiang, W., Johnson, T.T.: Star-based reachability analysis of deep neural networks. In: Formal Methods–The Next 30 Years: Third World Congress, FM 2019, Porto, Portugal, October 7–11, 2019, Proceedings 3. pp. 670–686. Springer (2019)

54. Tran, H.D., Yang, X., Lopez, D.M., Musau, P., Nguyen, L.V., Xiang, W., Bak, S., Johnson, T.T.: Nnv: The neural network verification tool for deep neural networks and learning-enabled cyber-physical systems (2020)
55. Vovk, V., Gammerman, A., Shafer, G.: Algorithmic Learning in a Random World. Springer, New York, 2005 edition edn. (Mar 2005)
56. Xiang, W., Shao, Z.: Approximate bisimulation relations for neural networks and application to assured neural network compression. In: 2022 American Control Conference (ACC). pp. 3248–3253. IEEE (2022)
57. Xiang, W., Shao, Z.: Safety verification of neural network control systems using guaranteed neural network model reduction. In: 2022 IEEE 61st Conference on Decision and Control (CDC). pp. 1521–1526. IEEE (2022)
58. Xu, C., Xie, Y.: Conformal prediction interval for dynamic time-series. In: Proceedings of the 38th International Conference on Machine Learning. pp. 11559–11569. PMLR (Jul 2021), `https://proceedings.mlr.press/v139/xu21h.html`, iSSN: 2640-3498
59. Xue, B., Zhang, M., Easwaran, A., Li, Q.: Pac model checking of black-box continuous-time dynamical systems. IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems **39** (07 2020). https://doi.org/10.1109/TCAD.2020.3012251
60. Zarei, M., Wang, Y., Pajic, M.: Statistical verification of learning-based cyber-physical systems. In: Proceedings of the 23rd International Conference on Hybrid Systems: Computation and Control. HSCC '20, Association for Computing Machinery, New York, NY, USA (2020). https://doi.org/10.1145/3365365.3382209, `https://doi.org/10.1145/3365365.3382209`
61. Zhang, M., Zhang, Y., Zhang, L., Liu, C., Khurshid, S.: Deeproad: Gan-based metamorphic testing and input validation framework for autonomous driving systems. In: Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering. p. 132–142. ASE '18, Association for Computing Machinery, New York, NY, USA (2018). https://doi.org/10.1145/3238147.3238187, `https://doi.org/10.1145/3238147.3238187`
62. Wang, Y., Zhou, W., Fan, J., Wang, Z., Li, J., Chen, X., Huang, C., Li, W. and Zhu, Q.: Polar-express: Efficient and precise formal reachability analysis of neural-network controlled systems. In: IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems. (2023).
63. Xin, L., Tang, Z., Gai, W. & Liu, H. Vision-based autonomous landing for the UAV: A review. *Aerospace.* **9**, 634 (2022)
64. Tang, C. & Lai, Y. Deep Reinforcement Learning Automatic Landing Control of Fixed-Wing Aircraft Using Deep Deterministic Policy Gradient. *2020 International Conference On Unmanned Aircraft Systems (ICUAS).* pp. 1-9 (2020)
65. Oszust, M., Kapuscinski, T., Warchol, D., Wysocki, M., Rogalski, T., Pieniazek, J., Kopecki, G., Ciecinski, P. & Rzucidlo, P. A vision-based method for supporting autonomous aircraft landing. *Aircraft Engineering And Aerospace Technology.* **90**, 973-982 (2018)

# 7    Appendix

Our appendix is organized into five parts:

- Subsection 7.1 describes complementary definitions, which supplement those presented in the main body.
- Subsection 7.2 provides all the proofs for our theorems.
- Subsection 7.3 provides qualitative results.
- Subsection 7.4 shows the details of our case study.
- Subsection 7.5 includes all complementary algorithms related to this paper.

## 7.1    Complementary definitions

This section starts with widely accepted concepts to provide an exhaustive formalization behind Secs. 2 and 3.

**Definition 9 (State function).** *Given initial state $s_0$, the state of (either low- or high-dimensional) system $M$ at discrete time $t$ is $\varphi_M(s_0, t)$, given by* state function $\varphi_M : S \times \mathbb{N}_{\geq 0} \to S$.

**Definition 10 (Trajectory for $M_{ld}$ and $M_{hd}$).** *From state $s_0$, system $M_{ld}$ produces the* low-dimensional trajectory $\tau_{ld}$ *for $T$ steps:* $\tau_{ld}(s_0, T) = [s_0, \varphi_{M_{ld}}(s_0, 1), \ldots, \varphi_{M_{ld}}(s_0, T)]$. *From state $s_0$, system $M_{hd}$ produces the* high-dimensional trajectory $\tau_{hd}$ *for $T$ steps:* $\tau_{hd}(s_0, T) = [s_0, \varphi_{M_{hd}}(s_0, 1), \ldots, \varphi_{M_{hd}}(s_0, T)]$.

**Definition 11 (Lipschitz Constant for LDC).** *An LDC $c_{ld} : S \to U$ is Lipschitz-continuous with* Lipschitz constant $L$ *on state region $\bar{S}$ if there exists $L \geq 0$ such that:* $\forall s_1, s_2 \in \bar{S} \cdot \|c_{ld}(s_1) - c_{ld}(s_2)\|_2 \leq L\|s_1 - s_2\|_2$.

Here we formalize the two notions of discrepancy.

**Definition 12 (Trajectory-based discrepancy).** *A* trajectory-based discrepancy $\beta(S_0)$ *is the supremum of the L1 differences between high- and low-dimensional trajectories starting in initial set $S_0$ at time $t \in [0..T]$:*

$$\beta(S_0) := \sup_{\forall s_0 \in S_0, t \in [0..T]} \|\tau_{hd}(s_0, t) - \tau_{ld}(s_0, t)\|_1.$$

**Definition 13 (Action-based discrepancy).** *An* action-based discrepancy $\gamma(S_0)$ *is the supremum of the control action differences $\delta$ between an HDC and LDC on the initial set $S_0$ for each time until horizon $T$:*

$$\gamma(S_0) := \sup_{\forall s_0 \in S_0, t \in [0..T]} \|c_{hd}\big(g(\varphi_{M_{hd}}(s_0, t))\big) - c_{ld}\big(\varphi_{M_{ld}}(s_0, t)\big)\|_1.$$

The definitions and theorems below straightforwardly extend Defs. 5, 6, 7, 8 and Thms. 1, 2 to the multi-LDC case.

**Definition 14 (Multi-LDC Trajectory-inflated reachable set).** *Given a set of distribution* $\mathbf{D_0} = \{D_1, ..., D_m\}$ *over partitioned initial set* $\mathbf{S_0} = \{S_1, ..., S_m\}$ *that is controlled by* $\{c_{ld}^1, ..., c_{ld}^m\}$, *reachable set* $\mathsf{rs}(S_i, t)$, *and its trajectory discrepancy* $\bar{\boldsymbol{\beta}} = \{\bar{\beta}(D_1), ..., \bar{\beta}(D_m)\}$, *a* Multi-LDC trajectory-inflated reachable set *is defined as:*

$$\mathsf{irs}(\mathbf{S_0}, t, \bar{\boldsymbol{\beta}}) = \bigcup_{i=1}^{m} \mathsf{irs}(S_i, t, \bar{\beta}(D_i)) = \bigcup_{i=1}^{m} \{s \in S \mid \exists s' \in \mathsf{rs}(S_i, t) \cdot \|s - s'\|_1 \leq \bar{\beta}(D_i)\}.$$

**Definition 15 (Multi-LDC Trajectory-inflated reachable tube).** *Given a set of distribution* $\mathbf{D_0} = \{D_1, ..., D_m\}$ *over partitioned initial set* $\mathbf{S_0} = \{S_1, ..., S_m\}$ *that is controlled by* $\{c_{ld}^1, c_{ld}^2, ..., c_{ld}^m\}$, *reachable set* $\mathsf{rs}(S_i, t)$, *and its trajectory discrepancy* $\bar{\boldsymbol{\beta}} = \{\bar{\beta}(D_1), ..., \bar{\beta}(D_m)\}$, *a* Multi-LDC trajectory-inflated reachable tube *is defined as:*

$$\begin{aligned} \mathsf{irt}(\mathbf{S_0}, \bar{\boldsymbol{\beta}}) &= \bigcup_{i=1}^{m} \mathsf{irt}(S_i, \bar{\beta}(D_i)) \\ &= \bigcup_{i=1}^{m} \left[ \mathsf{irs}(S_i, 0, \bar{\beta}(D_i)), \mathsf{irs}(S_0, 1, \bar{\beta}(D_i)), ..., \mathsf{irs}(S_i, T, \bar{\beta}(D_i)) \right]. \end{aligned}$$

**Theorem 4 (Confident multi-LDCs trajectory-based overapproximation).** *Consider a set of distribution* $\mathbf{D_0} = \{D_1, ..., D_m\}$ *over split initial set* $\mathbf{S_0} = \{S_1, ..., S_m\}$, *a high-dimensional system* $M_{hd}$, *and a set of low-dimensional systems with a set of controllers* $c_{ld} = \{c_{ld}^1, ..., c_{ld}^m\}$ *that approximates it with an* $\alpha$-*confident statistical trajectory-based discrepancy function* $\bar{\boldsymbol{\beta}} = \{\bar{\beta}(D_1), ... \bar{\beta}(D_m)\}$. *Then the trajectory-inflated low-dimensional tube* $\mathsf{irt}_{M_{ld}}(\mathbf{S_0}, \bar{\beta}(D_0))$ *contains the high-dimensional reachable tube* $\mathsf{rt}_{M_{hd}}(S_0)$ *with probability* $1 - \alpha$ *over the time* $T$:

$$\Pr_{D_1...D_m} \left[ \mathsf{rt}_{M_{hd}}(S_0) \subseteq \mathsf{irt}_{M_{ld}}(\mathbf{S_0}, \bar{\boldsymbol{\beta}}) \right] \geq 1 - \alpha$$

**Definition 16 (Multi-LDC action-inflated reachable set).** *Given distributions* $\{D_1, ..., D_m\}$ *over each region of grid* $\mathbf{S_0} = \{S_1, ..., S_m\}$ *of initial set* $S_0$, *an LDC for each region* $\{c_{ld}^1, c_{ld}^2, ..., c_{ld}^m\}$, *low-dimensional reachable sets* $\mathsf{rs}(S_i, t)$, *dynamics* $f$, *and action-based discrepancy functions for each LDC and state region* $\bar{\gamma} = \{\bar{\gamma}(D_i)\}_{i=1}^{m}$, *low-dimensional control bounds* $[u_{min}^i(t), u_{max}^i(t)] \supseteq c_{ld}^i(\mathsf{rs}(S_i, t))$, *the* multi-LDC action-inflated reachable set *contains a union of states reachable by inflating the actions bounds of the individual controllers:*

$$\mathsf{irs}(\mathbf{S_0}, t+1, \bar{\gamma}) = \bigcup_{i=1}^{m} \mathsf{irs}(S_i, t, \bar{\gamma}(D_i))$$

$$= \bigcup_{i=1}^{m} \left\{ f(s_i, u_i) \mid s_i \in \mathsf{rs}(S_i, t), u_i \in \left[ u_{min}^i(t) - \bar{\gamma}(D_i), u_{max}^i(t) + \bar{\gamma}(D_i) \right] \right\}$$

**Definition 17 (Multi-LDC action-inflated reachable tube).** *Given an initial set* $S_0$, *an initial state-space grid* $\mathbf{S_0} = \{S_1, ..., S_m\}$, *multiple LDCs*

$\{c_{ld}^1, c_{ld}^2, \ldots, c_{ld}^m\}$, *dynamics $f$, and action-based discrepancy functions for each LDC and state region $\bar{\gamma} = \{\bar{\gamma}(S_i)\}_{i=1}^m$, the* multi-LDC action-inflated reachable tube *is a recursive sequence of inflated action-based reachable sets:*

$$\mathsf{irt}(S_0, \bar{\gamma}) = \left[S_0, \mathsf{irs}_1(S_0, S_0, \bar{\gamma}), \mathsf{irs}_2(S_0, \mathsf{irs}_1, \bar{\gamma}), \ldots, \mathsf{irs}_T(S_0, \mathsf{irs}_{T-1}, \bar{\gamma})\right].$$

**Theorem 5 (Confident multi-LDCs action-based overapproximation).**
*Consider an initial set $S_0$ that can be partitioned into $m$ subsets, $\mathbf{S_0} = \{S_1, ..., S_m\}$, a high-dim. system $M_{hd}$ with controller $c_{hd}$, and low-dimensional systems $M_{ld}^1, \ldots, M_{ld}^m$ with respective controllers $c_{ld}^1, \ldots, c_{ld}^m$ that approximate $c_{hd}$ with $\alpha$-confident statistical action-based discrepancies $\bar{\gamma} = \{\bar{\gamma}(D_i)\}_{i=1}^m$. Then the action-inflated low-dimensional tube $\mathsf{irt}_{M_{ld}}(S_0, \bar{\gamma})$ contains the high-dimensional tube $\mathsf{rt}_{M_{hd}}(S_0)$ with probability $1 - \alpha$ over the time $T$:*

$$\mathrm{Pr}_{D_1 \ldots D_m} \left[\mathsf{rt}_{M_{hd}}(S_0) \subseteq \mathsf{irt}_{M_{ld}}(\mathbf{S_0}, \bar{\gamma})\right] \geq 1 - \alpha$$

## 7.2   Proofs of theorems

This section contains theorem proofs that could not be included in the main body of the paper due to the page limit.

### Proof of Theorem 1 (Confident trajectory-based overapproximation)

*Proof.* Based on Lemma 1, after calculating the conformal bound $\bar{\beta}(D_0)$ (in practice, $D_0 = \mathrm{Uniform}(S_0)$) with specified miscoverage $\alpha$ for the statistical trajectory-based discrepancy, for any new initial point $s_0$ sample uniformly and i.i.d. from $S_0$, the following holds:

$$\forall t \in [0..T] \cdot \mathrm{Pr}_{D_0} \left[\|\varphi_{hd}(s_0, t) - \varphi_{ld}(s_0, t)\|_1 \leq \bar{\beta}(D_0)\right] \geq 1 - \alpha$$

The low-dimensional reachable set $\mathsf{rs}(S_0, t)$ is guaranteed to contain all trajectories starting from the initial set $S_0$ in system $M_{ld}$; The $\bar{\beta}(D_0)$ is the confidence upper bound of the maximum L1 norm over the whole trajectories starting from the initial set. Then, with at least $1 - \alpha$ probability, the statistical lower bound of the trajectory-inflated reachable set is:

$$\mathsf{irs}_{low}(S_0, t) = \mathsf{rs}_{low}(S_0, t) - \bar{\beta}(D_0)$$

Similarly, the statistical upper bound of the trajectory-inflated reachable set is:

$$\mathsf{irs}_{up}(S_0, t) = \mathsf{rs}_{up}(S_0, t) + \bar{\beta}(D_0)$$

Then if any trajectory differences starting from this initial set satisfy the specified confidence, the reachable set containing all trajectories will satisfy this confidence, too. Since the whole reachable tube contains all the reachable sets, this inflated reachable tube also satisfies this probability bound:

$$\mathrm{Pr}_{D_0} \left[\mathsf{rt}_{M_{hd}}(S_0) \subseteq \mathsf{irt}_{M_{ld}}(S_0, \bar{\beta}(D_0))\right] \geq 1 - \alpha.$$

**Proof of Theorem 2 (Confident action-based overapproximation)**

*Proof.* Consider a low-dimensional reachable tube inflated with action-based discrepancy bounds. Given the statistical action-based upper bound $\bar{\gamma} = \bar{\gamma}(D_0)$ (in practice, $D_0 = \mathrm{Uniform}(S_0)$) with one LDC over the initial set $S_0$, starting from any states in $S_0$ sample uniformly, in the next $T$ time steps, control action difference between LDC and HDC is within this conformal bound with at least $1 - \alpha$ confidence,

$$\mathrm{Pr}_{s_0 \sim \mathrm{Uniform}(S_0)} \left[ \max_{t=0..T} \| c_{hd}\big(g(\varphi_{M_{hd}}(s_0, t))\big) - c_{ld}\big(\varphi_{M_{ld}}(s_0, t)\big) \|_1 \leq \bar{\gamma} \right] \geq 1 - \alpha.$$

Therefore, the $(1 - \alpha)$-confident statistical upper and lower bounds on the $c_{hd}$ output are respectively $[u_{min}^{c_{ld}} - \bar{\gamma}(D_0), u_{max}^{c_{ld}} + \bar{\gamma}(D_0)]$ for each time step in the whole time horizon $T$, where $u_{min}^{c_{ld}}$ and $u_{max}^{c_{ld}}$ are the bounds on $c_{ld}$ output based on the reachability analysis in Step 2 and the $\bar{\gamma}(D_0)$ has been calculated. Through the reachability analysis, we calculated each inflated $\gamma$-based reachable set given inflated control bounds in each time step. We assert that,

$$\mathrm{Pr}_{D_0} \left[ \forall t \in [0, T], \mathsf{rs}_{M_{hd}}(S_0, t) \subseteq \mathsf{irs}_{M_{ld}}(S_0, t, \bar{\gamma}(D_0)) \right] \geq 1 - \alpha$$

Since the action-inflated reachable tube contains all the action-inflated reachable sets, the reachable tube of $c_{hd}$ also obtains the probabilistic containment in the inflated reachable tube, namely:

$$\mathrm{Pr} \left[ \mathsf{rt}_{M_{hd}}(S_0) \subseteq \mathsf{irt}_{M_{ld}}(S_0, \bar{\gamma}(D_0)) \right] \geq 1 - \alpha.$$

**Proof of Theorem 3 (Confident guarantee of HDC safety)**

*Proof.* With the help of Thms. 1 and 2, given all the trained $c_{ld}$ and their corresponding $\alpha$-confident statistical action-based discrepancy functions $\bar{\gamma} = \{\bar{\gamma}(S_i)\}_{i=1}^n$ or $\alpha$-confident statistical trajectory-based discrepancy function $\bar{\beta} = \{\bar{\beta}(S_i)\}_{i=1}^n$, both action-based and trajectory-based inflated reachable tubes contain the ground truth HDC reachable tube with a lower probability bound:

$$\mathrm{Pr} \left[ \mathsf{rt}_{M_{hd}}(S_0) \subseteq \mathsf{irt}_{M_{ld}}(S_0, \bar{\beta}) \right] \geq 1 - \alpha.$$

$$\mathrm{Pr} \left[ \mathsf{rt}_{M_{hd}}(S_0) \subseteq \mathsf{irt}_{M_{ld}}(S_0, \bar{\gamma}) \right] \geq 1 - \alpha.$$

We check the safety of the HDC by checking that $\mathsf{irs}(S_0^j, T) \subseteq G$. Since the target set and goal set are fixed, the confidence of safety depends on the confidence of the reachable tube and set starting from partitioned initial regions $\mathbf{S_{ver}} = S_1, \cdots, S_n$. Since the verification regions are independent, the probability that the HDC safety set $S_{safe}$ being contained in the ground truth safety set $S_{safe}^*$ is at least $(1 - \alpha)$:

$$\mathrm{Pr} \left[ S_{safe} \subseteq S_{safe}^* \right] \geq 1 - \alpha$$

**Proof of Theorem 4 (Confident multi-LDCs trajectory-based overapproximation)**

*Proof.* This proof is based on Thm. 1. Given the statistical action-based upper bound $\bar{\boldsymbol{\beta}} = \{\bar{\beta}(D_i)\}_{i=1}^m$ for LDCs ($D_i = \text{Uniform}(S_i)$) over the partitioned initial set $\mathbf{S_0} = \{S_i\}_{i=1}^m$, for each subset $S_i$, Thm. 1 gives us that

$$\Pr\left[\mathsf{rt}_{M_{hd}}(S_i) \subseteq \mathsf{irt}_{M_{ld}}(S_i, \bar{\beta}(D_i))\right] \geq 1 - \alpha.$$

Then we iterate each subset of $S_0$ to get the probability containment,

$$\Pr\left[\bigcup_{i=1}^m \left(\mathsf{rt}_{M_{hd}}(S_i) \subseteq \mathsf{irt}_{M_{ld}}(S_i, \bar{\beta}(D_i))\right)\right] \geq 1 - \alpha.$$

Therefore, we add all the subregions together and conclude that

$$\Pr_{D_1...D_m}\left[\mathsf{rt}_{M_{hd}}(S_0) \subseteq \mathsf{irt}_{M_{ld}}(\mathbf{S_0}, \bar{\boldsymbol{\beta}})\right] \geq 1 - \alpha$$

**Proof of Theorem 5 (Confident multi-LDCs action-based overapproximation)**

*Proof.* This proof is based on Thm. 2. Given the statistical action-based upper bound $\bar{\boldsymbol{\gamma}} = \{\bar{\gamma}(D_i)\}_{i=1}^m$ for LDCs ($D_i = \text{Uniform}(S_i)$) over the initial set $S_0 = \{S_i\}_{i=1}^m$, for each subset $S_i$, Thm. 2 gives us that

$$\Pr\left[\mathsf{rt}_{M_{hd}}(S_i) \subseteq \mathsf{irt}_{M_{ld}}(S_i, \bar{\gamma}(D_i))\right] \geq 1 - \alpha.$$

And then we iterate each subset of $S_0$ to get the probability containment,

$$\Pr\left[\bigcup_{i=1}^m \left(\mathsf{rt}_{M_{hd}}(S_i) \subseteq \mathsf{irt}_{M_{ld}}(S_i, \bar{\gamma}(D_i))\right)\right] \geq 1 - \alpha.$$

Therefore, we add all the subregions together, we can conclude that

$$\Pr_{D_1...D_m}\left[\mathsf{rt}_{M_{hd}}(S_0) \subseteq \mathsf{irt}_{M_{ld}}(\mathbf{S_0}, \bar{\gamma})\right] \geq 1 - \alpha$$

### 7.3  Complementary qualitative illustration

For most truly safe partitions of the initial state set, both of our approaches successfully verified safety. A successful example of the action-based approach is shown on the left of Fig. 2. Given the narrow margins of safety in our benchmarks, even small overapproximation or discrepancy errors can lead to losing the safety guarantees, resulting in a false negative verification outcome. Such an example is shown on the right of Fig. 2: the reach set is not fully contained in the goal set (even though all simulations from the initial set are contained and therefore safe).

These figures display only action-inflated reachable tubes: low-dimensional tubes and trajectory-inflated tubes are visually indistinguishable due to the low statistical discrepancies between the HDC and trained LDCs. Our reachable tube computations enable solving reach-avoid problems, which specify an avoid set in addition to a reach set.
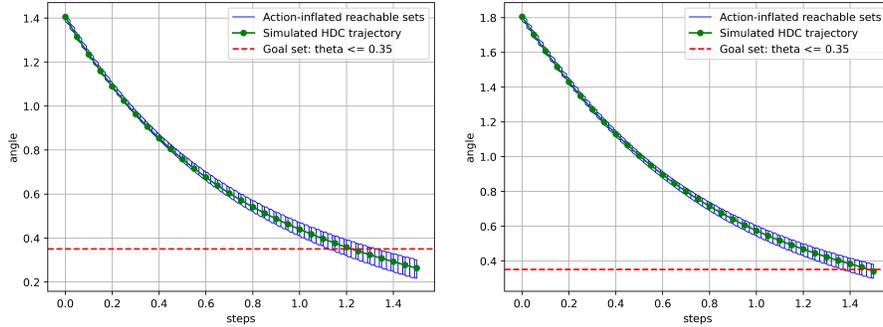


**Fig. 2.** Reachability examples of the action-based approach on the inverted pendulum. Left: safe, right: unsafe.
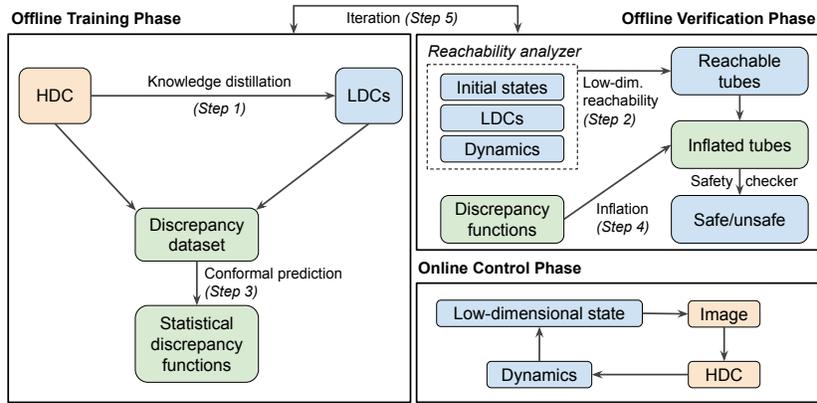


**Fig. 3.** Three phases of our approach: training, verification, and control. Orange shows high-dimensional elements, and green shows novel contributions.

## 7.4   Details of experimental setup

**Approach and initial set**  We present a detailed summary of our approach, illustrated step-by-step in Fig. 3. Our methodology begins with matching the
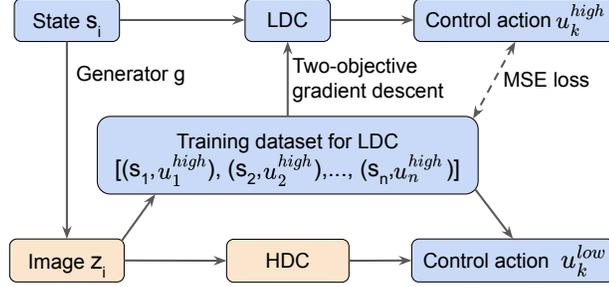
**Fig. 4.** Training an LDC with supervision from an HDC.

initial state and image, where the initial set $S_0$ can be construed as a quasi-inverse set $g^{-1}(z_0)$ of an initial image $z_0$. This initial set $S_0$ is readily available to us due to our assumption that the initial conditions are sufficiently instrumented, in contrast to the subsequent states. Moreover, our approach is not confined to HDCs with solely high-dimensional inputs. For instance, in all of our case studies, the HDCs accommodate velocities alongside the images. Furthermore, our approach could be extended to HDCs that process sequences of images (e.g., to extract velocity information).

**Benchmark systems** The equations below describe the discrete-time, friction-free, transition dynamics of the systems. We have omitted the time step $\Delta t$ for simplicity, assuming a constant time step, where the effects of the changes in all relevant variables are implicitly assumed to occur over a fixed interval.

The inverted pendulum dynamics are as follows with the fixed parameters — the mass of the rod ($m = 0.1kg$), rod inertia ($I = 0.125kg \cdot m^2$), rod length ($L = 0.25m$), and damping ratio ($c = 0kg^{-1} \cdot m^{-2}$):

$$\theta_{t+1} = \theta_t + \dot{\theta}_t, \qquad \dot{\theta}_{t+1} = \dot{\theta}_t + \frac{1}{I}(u_t - c\dot{\theta}_t + mgL\sin(\theta_t)) \tag{5}$$

The mountain car dynamics are as follows with the fixed parameters — the mass of the car ($m = 2.5 \times 10^{-4}kg$) and the force produced by the car's engine ($F = 3.75 \times 10^{-7}N$):

$$x_{t+1} = x_t + v_t, \qquad v_{t+1} = v_t + \frac{F}{m}u_t - mg\cos(3x_t) \tag{6}$$

The cartpole dynamics are as follows with the fixed parameters — the mass of the cart ($m_c = 1kg$), the mass of the pole ($m_p = 0.1kg$) and the length of the pole ($l = 0.5m$):

$$x_{t+1} = x_t + v_t, \qquad \dot{x}_{t+1} = \dot{x}_t + \frac{u + m_p l(\dot{\theta}^2 \sin\theta - \ddot{\theta}\cos\theta)}{m_c + m_p} \tag{7}$$

$$\theta_{t+1} = \theta_t + \dot{\theta}_t, \qquad \dot{\theta}_{t+1} = \dot{\theta} + \frac{g\sin\theta + \cos\theta(\frac{-u - m_p l\dot{\theta}^2 \sin\theta}{m_c + m_p})}{l(\frac{4}{3} - \frac{m_p \cos^2\theta}{m_c + m_p})} \tag{8}$$

**Controller details** Our HDCs $c_{hd}$ were trained in continuous action spaces with the deep deterministic policy gradient (DDPG) [36], which is an off-policy, actor-critic algorithm for deep reinforcement learning. For the pendulum, the $c_{hd}$ input is a $64 \times 64$ image and an angular velocity $\dot{\theta}$. The structure of $c_{hd}$ consists of convolutional layers (10 for IP, 12 for MC, 12 for CP) containing 400 neurons, fully connected layers, pooling layers, and flattened layers, along with $ReLU$ activation functions. for the mountain car, we have the same $c_{hd}$ structure and input dimensions (with velocity $v$ instead of angular velocity $\dot{\theta}$). Although these controllers perform well, they are impractical to verify directly due to their complexity.

To approximate high-dimensional controllers, we train feedforward neural networks $c_{ld}$ with only low-dimensional state inputs to imitate the performance of the $c_{hd}$. Step 1 in the training process of $c_{ld}$ is illustrated in Fig. 4 and formalized in Alg. 3. The structure of $c_{ld}$ is simpler to enable exhaustive verification: 2 layers with 20 neurons each and $Sigmoid$ and $Tanh$ activation functions.

**Experimental verification details**

- In the IP case, the initial set we considered is $S_0^{ip} = \{(\theta_0, \dot{\theta}_0) \in [0, 2] \times [-2, 0]\}$. In $T = 30$ steps, we checked whether the rod stays in the target set $G_{ip} = \{\theta \in [0, 0.35]\}$.
- In the MC case, the initial set is $S_0^{mc} = \{(x_0, v_0) \in [-0.6, -0.4] \times [-0.02, 0.05]\}$. Given $T = 60$ steps, we checked whether the mountain car will stay in the target set $G^{ip} = \{x \in [0.45, \infty]\}$.
- In the CP case, the initial set is $S_0^{cp} = \{(x_0, v_0, \theta_0, \dot{\theta}_0) \in [0, 0.1] \times [0, 0.1] \times [0.05, 0.15] \times [-0.4, -0.35]\}$. Given $T = 20$ steps, we checked whether the cartpole would stay in the target set $G^{ip} = \{\theta \times x \in [-0.2, 0.2] \times [0, 0.2]\}$.

## 7.5  Complementary algorithms

This section provides the algorithms discussed in previous sections. Algorithm 3 refers to the first step—knowledge distillation LDC training. Algorithms 4 and 5 represent Step 3 in our end-to-end method—computation of action- and trajectory-based discrepancies. Algorithm 6 is the counterpart of Algorithm 1.

---
**Algorithm 3** Training an LDC with HDC supervision

---
**function** TRAINLDC(HDC $c_{hd}$, image generator $g$, initial state region $\bar{S}$, time horizon $T$, threshold for Lipschitz constant $\lambda$, threshold for MSE $\epsilon$)

    $X_{init} \leftarrow s_1, s_2, \ldots, s_n \sim \text{Uniform}(\bar{S})$

    **for** $i = 1$ to $n$ **do**

        $X \leftarrow \tau_{hd}(s_i, T)$            ▷ Simulate HDC system to get low-dimen. state

        **for** $j = 1$ to $T$ **do**

            $s_j \leftarrow \varphi_{M_{hd}}(s_i, j)$

            $z_j \leftarrow g(s_j)$

            $Y \leftarrow c_{hd}(z_i)$            ▷ Store HDC control action dataset

        **end for**

    **end for**

    Training dataset $\mathcal{D}_{tr} \leftarrow (X, Y)$

    $c_{ld} \leftarrow$ two-objective gradient descent$(\mathcal{D}_{tr}, \lambda, \epsilon)$

    **return** $c_{ld}$

**end function**

---

---
**Algorithm 4** Computation of trajectory-based discrepancy

---
**function** COMPUTETRAJDISCR(LDC $c_{ld}$, HDC $c_{hd}$, image generator $g$, state region $\bar{S}$, confidence $\alpha$, sample count $N$, time steps $T$)

    $s_0, s_1, \ldots, s_N \sim \text{Uniform}(\bar{S})$

    **for** $i = 1$ to $N$ **do**

        $\delta_i \leftarrow \max_{t=1,\ldots,T} \left\| \varphi_{M_{hd}}(s_i, t) - \varphi_{M_{ld}}(s_i, t) \right\|_1$    ▷ non-conformity scores

    **end for**

    $r \leftarrow \lceil (N+1)(1-\alpha) \rceil$            ▷ conformal quantile

    **return** $r$-th smallest value among $[\delta_1, \ldots, \delta_N, \infty]$

**end function**

---

---
**Algorithm 5** Computation of action-based discrepancy

---
**function** COMPUTEACTIONDISCR(LDC $c_{ld}$, HDC $c_{hd}$, image generator $g$, state region $\bar{S}$, confidence $\alpha$, sample count $N$)

    $s_0, s_1, \ldots, s_N \sim \text{Uniform}(\bar{S})$

    **for** $i = 1$ to $N$ **do**

        $\delta_i \leftarrow \max_{t=0..T} \left\| c_{hd}(g(\varphi_{M_{hd}}(s_0, t))) - c_{ld}(\varphi_{M_{hd}}(s_0, t)) \right\|_1$    ▷
non-conformity scores

    **end for**

    $r \leftarrow \lceil (N+1)(1-\alpha) \rceil$            ▷ conformal quantile

    **return** $r$-th smallest value among $[\delta_1, \ldots, \delta_N, \infty]$

**end function**

---

---

**Algorithm 6** Iterative LDC training for the trajectory-based approach

---

**function** IterativeTrainingTB(HDC $c_{hd}$, image generator $g$, sample count $N$, initial state space $S_0$, confidence $\alpha$, discrepancy threshold $\xi$, time steps $T$)

    $\lambda, \epsilon \leftarrow$ initial values

    $\mathbf{S_0} \leftarrow$ initial gridding of $S_0 : S_1, S_2, \ldots$

    **while** Computing resources last **do**

        **for** $i = 1$ to $|\mathbf{S_0}|$ **do**

            $c_{ld}^i \leftarrow$ TrainLDC($c_{hd}$, g, $S_i$, $\lambda$, $\epsilon$)

            $\delta^i \leftarrow$ ComputeTrajDiscr($c_{ld}^i$, $c_{hd}$, $g$, $S_i$, $\alpha$, $N$, $T$)

            **if** $\delta^i > \xi$ **then**

                $\epsilon \leftarrow \epsilon/2$                           ▷ Reduce MSE threshold

            **end if**

        **end for**

        **if** $\hat{\delta} > \epsilon$ in some sub-regions $\hat{\mathbf{S}} \subseteq \mathbf{S_0}$ **then**

            $\mathbf{S_0} \leftarrow \mathbf{S_0'}$ with refined re-gridding of $\hat{S}$        ▷ Reduce Lipschitz threshold

        **end if**

        **if** $\hat{\delta} \leq \xi \wedge \mathsf{rs}_{M_{ld}}(\hat{\mathbf{S}}, T) \not\subseteq G$ in some sub-regions $\hat{\mathbf{S}} \subseteq \mathbf{S_0}$ **then**

            $\lambda \leftarrow \lambda/2$ and keep the same $\epsilon$ in $\hat{\mathbf{S_0}}$     ▷ Reduce Lipschitz threshold

        **end if**

        $\mathbf{S_0} \leftarrow \mathbf{S_0'}$                               ▷ Use the updated grid

    **end while**

    $\bar{\beta} \leftarrow \delta^1, \delta^2, \ldots$

    **return** $c_{ld}^1, c_{ld}^2, \ldots, c_{ld}^n, \bar{\beta}$

**end function**

---