

# A Distance-based Anomaly Detection Framework for Deep Reinforcement Learning

Anonymous authors

Paper under double-blind review

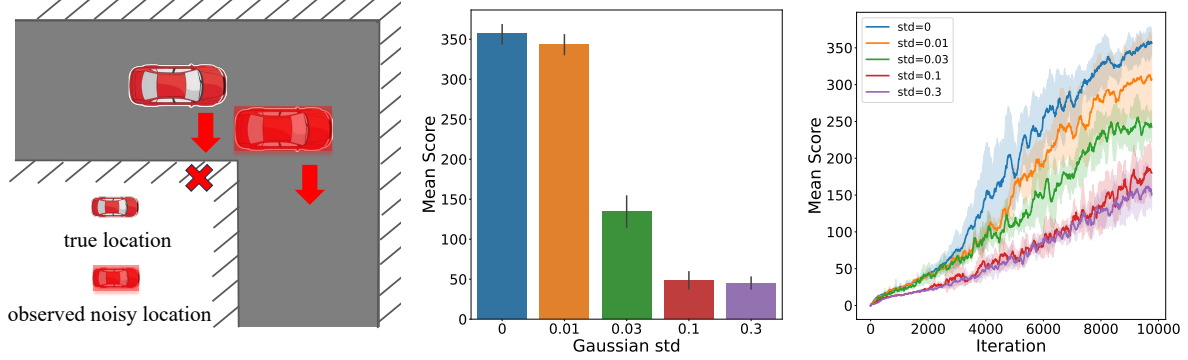
## Abstract

In deep reinforcement learning (RL) systems, abnormal states pose significant risks by potentially triggering unpredictable behaviors and unsafe actions, thus impeding the deployment of RL systems in real-world scenarios. It is crucial for reliable decision-making systems to have the capability to cast an alert whenever they encounter unfamiliar observations that they are not equipped to handle. In this paper, we propose a novel Mahalanobis distance-based (MD) anomaly detection framework, called *MDX*, for deep RL algorithms. MDX simultaneously addresses random, adversarial, and out-of-distribution (OOD) state outliers in both offline and online settings. It utilizes Mahalanobis distance within class-conditional distributions for each action and operates within a statistical hypothesis testing framework under the Gaussian assumption. We further extend it to robust and distribution-free versions by incorporating Robust MD and conformal inference techniques. Through extensive experiments on both Atari games and autonomous driving scenarios, we demonstrate the effectiveness of our MD-based detection framework. MDX offers a simple, unified, and practical tool for enhancing the safety and reliability of RL systems in real-world applications.

## 1 Introduction

Deep reinforcement learning (RL) algorithms vary considerably in their performance and are highly sensitive to a wide range of factors, including the environment, state observations, and hyper-parameters (Jordan et al., 2020; Patterson et al., 2020). The lack of robustness of RL algorithms hinders their deployment in real-world scenarios, particularly in safety-critical applications, such as autonomous driving (Kiran et al., 2021). Recently, the reliability of RL algorithms has garnered substantial attention (Chan et al., 2020; Gu et al., 2024), emphasizing the need for anomaly detection-based strategies to build trustworthy RL systems (Haider et al., 2023; Danesh & Fern, 2021; Sedlmeier et al., 2020).

**Practical Scenarios.** Observed states often contain natural measurement errors (random noises), adversarial perturbations, and out-of-distribution (OOD) observations. For instance, consider an autonomous vehicle with malfunctioning or unreliable sensors or cameras. Under such circumstances, the collected data, such as the vehicle’s observed location, can be contaminated by random measurement errors. Furthermore, an autonomous car can encounter sensory inputs that have been adversarially manipulated regarding traffic signs. For example, a stop sign maliciously altered to be misclassified as a speed limit sign (Chen et al., 2019), increases the risk of traffic accidents. Regarding OOD samples, an RL policy trained to drive only on sunny days will struggle with observations from rainy days, which are beyond its trained experience. Such OOD observations can lead to safety violations, performance degradation, and potentially catastrophic failures. All these scenarios highlight the necessity of detecting inaccurate sensor signals from noisy state observations to ensure a vehicle’s accurate and reliable operation. Beyond autonomous driving, anomaly detection is critical in many other applications involving sequential decision-making. In healthcare, the RL agent might adjust treatment recommendations if it detects a sudden change in the patient’s health condition (Hu et al., 2022). Similarly, detecting fraud and anomalous market states in financial systems is becoming increasingly instrumental in preventing substantial financial losses from market manipulation and fraudulent activities (Hilal et al., 2022).



(a) Unsafe behavior in autonomous driving under noisy sensor signals. (b) Performance degradation when noises injected in policy deployment. (c) Performance degradation when noises injected during policy learning.

Figure 1: (a) An autonomous car navigates using location data observed from sensors such as GPS. Without an effective anomaly detection mechanism, inaccuracies or malfunctions in these sensors can cause the car to prematurely turn right, leading to a collision. (b) and (c): Performance degradation occurs when noisy states are observed in the Breakout environment. Gaussian noises with increasing standard deviations are injected into the state observations during policy deployment (b) and policy learning (c).

**Motivating Examples.** Figure 1(a) illustrates a potential collision scenario where an autonomous car, relying on noisy location data in the red region (such as GPS coordinate errors), turns right prematurely, risking an accident. Without anomaly detection, the car reacts incorrectly due to the location error. Figure 1(b) highlights how increasing measurement errors, represented by the standard deviation of Gaussian noises, dramatically degrade policy performance. For instance, autonomous cars with RL systems may take sub-optimal or unsafe actions when processing noisy sensory signals in deployment. In addition, incorporating excessive noise during online training (Figure 1(c)) can severely impair policy learning and diminish performance. These motivating examples underscore the importance of detecting different types of abnormal states for developing trustworthy RL systems in real-world applications.

Our research aims to provide a general framework for applying anomaly detection in deep RL problems, including problem formulation, detection algorithms, and evaluation scenarios. Specifically, we strive to develop an effective and unified anomaly detection framework for deep RL in *both offline and online settings*.

1. **Offline Setting.** In this setting, a dataset is fixed without additional online data collection. Given a pre-trained policy, our objective is to utilize a fixed dataset to develop a distance-based anomaly detector tailored for a pre-trained policy. This detector aims to effectively identify whether a state is an outlier <sup>1</sup>.
2. **Online Setting.** In this setting, the RL agent interacts with a noisy environment and continuously updates its policy. Our goal is to develop a detection strategy that identifies state outliers, which are outside the RL system’s training experience. Removing these outliers can prevent them from interfering with policy training.

Methodologically, we first design an RL outlier detection approach using Mahalanobis Distance (MD) (De Maesschalck et al., 2000) within a statistical hypothesis test framework and extend it to a robust MD version (Butler et al., 1993). These strategies are applied *in a parametric manner* under the Gaussian assumption for state features in each class, which may not always be accurate in practice. To address this limitation, we introduce a *non-parametric conformal version* of MD detection to relax the Gaussian assumption. We empirically investigate the effectiveness of these proposed detection approaches in both offline and online settings across a representative set of RL environments, including Atari games and autonomous driving. Our contributions can be summarized as follows:

- Our primary technical contribution is the design of RL outlier detection strategies based on the concepts of Mahalanobis Distance (MD), robust MD, and conformal inference. The anomaly detection

<sup>1</sup>Compared with the classical tasks of policy evaluation and learning in offline RL, our offline setting also utilizes a fixed dataset but specifically focuses on developing detection methods given a fixed policy.

strategies are specially developed for deep RL within a hypothesis test framework, accommodating both parametric (Gaussian assumption) and non-parametric (conformal calibration) approaches.

- Secondly, in our online setting, our anomaly detection can be applied to a dynamic dataset, where the RL policy continually improves when interacting with the environment. This dynamic setting contrasts with the simpler anomaly detection in supervised learning with a static dataset. To address this challenge, we particularly develop *moving window estimation* and *double self-supervised detectors* for anomaly detection in the online RL setting.
- To our best knowledge, we are the first to conduct a comprehensive study on distance-based anomaly detection in deep RL, covering all typical types of outliers. Our anomaly detectors can simultaneously identify random, adversarial, and out-of-distribution state outliers. We perform extensive experiments to verify the effectiveness of our proposed methods in both offline and online settings.

## 2 Related Work

**Anomaly Detection in Reinforcement Learning.** Anomaly detection has yet to be extensively explored in RL. The connection between anomaly detection and RL was first established in (Müller et al., 2022); however, their work is mainly conceptual and does not propose practical detection algorithms. Change point detection has been investigated in the tabular setting of RL, particularly in environments described as doubly inhomogeneous under temporal non-stationarity and subject heterogeneity (Hu et al., 2022). They focus on identifying “best data chunks” within the environment that exhibit similar dynamics for policy learning, while our detection focuses on anomaly detection in *deep* RL scenarios. Prior studies have also probed anomaly detection in specific RL contexts, such as the offline imitation learning with a transformer-based policy network (Wang et al., 2024) and detecting adversarial attacks within cooperative multi-agent RL (Kazari et al., 2023). However, these studies are limited to specific scenarios that do not address general anomaly detection, even in single-agent RL. Haider et al. (2023) proposed a model-based method using probabilistic dynamics models and bootstrapped ensembles, but this approach is computationally expensive. Our research aims to develop a unified and practical anomaly detection framework that applies to general RL scenarios.

**Distance-based Anomaly Detection.** Recently, there has been a growth of interest in developing anomaly detection strategies in deep learning scenarios (Pang et al., 2021; Elmrabit et al., 2020). In image classification, Mahalanobis distance (MD) was effectively applied by (Lee et al., 2018), who constructed a Mahalanobis confidence score by training a logistic regression detector using validation samples. This score was evaluated in a supervised way, relying on a *validation set*, and thus it is unsuitable for the RL setting. The “tied” covariance assumption used by (Lee et al., 2018), where class-conditional distributions of pre-trained features share the same covariance, was criticized as implausible by (Kamoi & Kobayashi, 2020) based on Gaussian discriminant analysis (Klecka et al., 1980). In contrast, our detection framework MDX avoids the unrealistic “tied covariance” assumption by estimating variance for each class using quadratic discriminant analysis. This approach extends linear boundaries to quadratic ones between classes, offering a more flexible and accurate detection (Hastie et al., 2009).

**Robust Statistics for RL.** Deep RL algorithms inherently face challenges related to instability and divergence due to the use of function approximation, bootstrapping, and off-policy learning (Sutton & Barto, 2018). Employing Mahalanobis distance (MD) for anomaly detection can be particularly sensitive during unstable learning phases. The computation of MD is based on Maximum Likelihood Estimate (MLE), which is susceptible to outliers or noisy data (Rousseeuw & Van Zomeren, 1990). Robust statistics (Huber, 2004) have been developed to address these robustness problems, especially leveraging robust estimation techniques that are not unduly affected by outliers. For example, Robust MD is a robust version of MD that employs robust estimators, e.g., Minimum Covariance Determinant (MCD) (Rousseeuw, 1984; Grübel, 1988), for location and covariance estimation (Maronna & Yohai, 2014).

**Conformal Prediction and Conformal Anomaly Detection.** Conformal anomaly detection (Ishimtsev et al., 2017; Laxhammar & Falkman, 2011) is based on the conformal prediction (Teng et al., 2023; Angelopoulos et al., 2021), a popular, modern technique for providing valid prediction intervals for arbitrar-

ily machine learning models. Conformal prediction has garnered increasing attention as it can provide a simple, distribution-free, and computationally effective way of tuning the distribution threshold. Its validity relies on the data exchangeability condition (Shafer & Vovk, 2008), where different orderings of samples are equally likely, but recent studies have verified its applicability in scenarios involving distribution shift (Tibshirani et al., 2019; Barber et al., 2023) and off-policy evaluation (Zhang et al., 2023). These examples justify the potential of using conformal inference to detect outliers in the context of RL.

### 3 Background

**Markov Decision Process.** The interaction of an agent with its environment can be modeled as a Markov Decision Process (MDP), a 5-tuple  $(\mathcal{S}, \mathcal{A}, R, P, \gamma)$ .  $\mathcal{S}$  and  $\mathcal{A}$  are the state and action spaces,  $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is the environment transition dynamics,  $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is the reward function and  $\gamma \in (0, 1)$  is the discount factor. The policy  $\pi$  is continually updated in this online interaction paradigm. Compared to the online setting, a recent popular paradigm for reinforcement learning is offline RL (Levine et al., 2020). In the offline setting, RL algorithms utilize previously collected data to extract policies without additional online data collection.

**Proximal Policy Optimization (PPO).** The policy gradient algorithm of Proximal Policy Optimization (PPO) (Schulman et al., 2017) has achieved state-of-the-art or competitive performance on Atari games (Bellemare et al., 2013) and MuJoCo robotic tasks (Todorov et al., 2012). Typical policy gradient algorithms optimize the expected reward function  $\rho(\theta, s_0) = \mathbb{E}_{\pi_\theta} [\sum_{t=0}^{\infty} \gamma^t r(s_t) \mid s_0]$  by using the policy gradient theorem (Sutton & Barto, 2018). Here  $\pi_\theta$  is the  $\theta$ -parameterized policy function. Trust Region Policy Optimization (TRPO) (Schulman et al., 2015) and PPO (Schulman et al., 2017) utilize constraints and advantage estimation to perform the update by reformulating the original optimization problem with the surrogate loss  $L(\theta)$  as:

$$L(\theta) = \mathbb{E}_t \left[ \frac{\pi_\theta(s_t, a_t)}{\pi_{\theta_{\text{old}}}(s_t, a_t)} A_{\pi_{\theta_{\text{old}}}}(s_t, a_t) \right], \quad (1)$$

where  $A_{\pi_{\theta_{\text{old}}}}$  is the generalized advantage function (GAE) (Schulman et al., 2018). PPO introduces clipping in the objective function in order to penalize changes to the policy that make  $\pi_\theta$  vastly different from  $\pi_{\theta_{\text{old}}}$ :

$$L^{\text{CLIP}}(\theta) = \mathbb{E}_t \left[ \min \left( \frac{\pi_\theta(s_t, a_t)}{\pi_{\theta_{\text{old}}}(s_t, a_t)} A_{\pi_{\theta_{\text{old}}}}(s_t, a_t), \text{clip} \left( \frac{\pi_\theta(s_t, a_t)}{\pi_{\theta_{\text{old}}}(s_t, a_t)}, 1 - \epsilon, 1 + \epsilon \right) A_{\pi_{\theta_{\text{old}}}}(s_t, a_t) \right) \right], \quad (2)$$

where  $\epsilon$  is a hyperparameter. We use PPO as the algorithm testbed to examine the efficacy of our anomaly detection framework. However, our detection methods are general and can be easily applied to other RL algorithms (Zhang & Yu, 2020) such as DQN (Mnih et al., 2015), A3C (Mnih et al., 2016), and DDPG (Lillicrap et al., 2016).

**Conformal Prediction.** Conformal anomaly detection (Laxhammar & Falkman, 2011; Ishimtsev et al., 2017) is grounded in conformal prediction (Shafer & Vovk, 2008; Angelopoulos et al., 2021), which aims to construct a confidence band  $\mathcal{C}_{1-\alpha}(X)$  for  $Y$  given a random data pair  $(X, Y) \sim \mathcal{P}$  and a confidence level  $1 - \alpha$ . Suppose we have a pre-trained model  $\hat{\mu}$  and a calibration dataset  $(X_1, Y_1), \dots, (X_n, Y_n)$  for conformal prediction. We can then compute a predictive interval for the new sample  $X_{n+1}$  to cover the unseen response  $Y_{n+1}$  by leveraging the empirical quantiles of the residuals  $|Y_i - \hat{\mu}(X_i)|$  on the calibration dataset. This further leads to valid prediction intervals such that:

$$\mathbb{P}(Y_{n+1} \in \mathcal{C}_{1-\alpha}(X_{n+1})) \geq 1 - \alpha, \quad (3)$$

where the confidence band is expected to be as small as possible while maintaining the desired coverage. A fundamental quantity in conformal prediction is the *non-conformity measure*, e.g., the residual  $|Y_i - \hat{\mu}(X_i)|$ , which measures how “different” an example is relative to a set of examples (Vovk et al., 2005).

### 4 Anomaly Detection in the Offline RL Setting

In this section, we design our MD-based detection framework MDX in the offline setting, where a fixed dataset collected from the environment and a pre-trained RL policy are provided. Based on the Gaussian

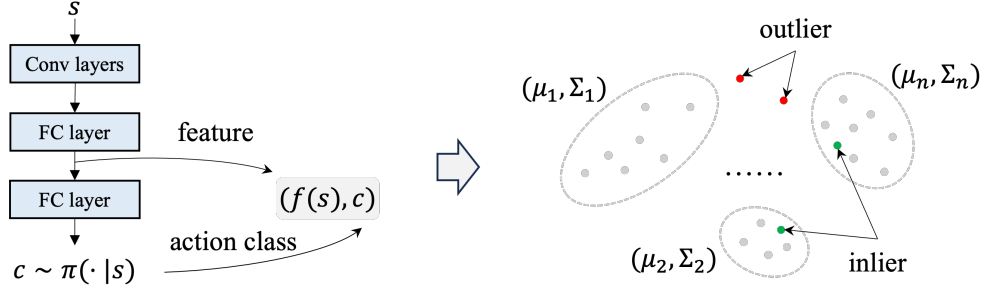


Figure 2: **The detection pipeline of MDX.** We feed the state into the policy network to extract the feature vector and identify its class. For each class, we estimate  $(\mu, \Sigma)$  and establish a detection threshold depicted as a dashed ellipse. To determine whether a new state is an outlier, we evaluate its features and compute the distance to the class centers. If the distance falls below the set threshold, the state is classified as an inlier (green points). Conversely, the state is marked as an outlier (red points).

assumption, we introduce the basic Mahalanobis Distance (MD) detection strategy. We then extend it to the robust MD and conformal MD-based detection methods. Finally, we present the deployment of MDX in a potentially noisy environment.

**Description of Detection Framework.** Our detection framework is structured around two core components: feature extraction and detector estimation. The process begins by assessing whether a state is anomalous, crucially dependent on the associated policy. A state that prompts the policy to initiate a potentially unsafe action is labeled an outlier. Specifically, we input the state into the policy network and access the feature vector extracted from the penultimate layer of this network. We categorize states according to the actions determined by the policy, where the underlying intuition is that states associated with the same action share similar features. To ascertain whether a new state is an outlier, we compute its distance from the established class centroids based on its feature vector. A state is deemed an outlier if the distance surpasses a set threshold. Figure 2 illustrates the operational flow of MDX. By ensuring that only states within the policy’s capability are considered valid, MDX thereby enhances the safety and reliability of the RL system.

#### 4.1 Mahalanobis Distance (MD)-based Detection

**Gaussian Assumption.** The given pre-trained parameterized RL policy  $\pi_\theta$  is a discriminative softmax classifier,  $\pi(a_t = c | s_t) = \exp(\mathbf{w}_c^\top f(s_t) + b_c) / \sum_{c'} \exp(\mathbf{w}_{c'}^\top f(s_t) + b_{c'})$ , where  $\mathbf{w}_c$  and  $b_c$  are the weight and bias of the policy classifier for action class  $c$ . The function  $f(\cdot)$  represents the output of the penultimate layer of the policy network  $\pi_\theta$ , serving as the state feature vector. Here,  $C = |A|$  is the size of the action space, and  $\mu_c$  is the mean vector of  $f(s)$  corresponding to the action class  $c$ <sup>2</sup>. If we assume that the class-conditional distribution follows a multivariate Gaussian distribution sharing a single covariance  $\Sigma$  (tied covariance) in a generative classifier, i.e.,  $\pi(f(s) | a = c) = \mathcal{N}(f(s) | \mu_c, \Sigma)$ , then the posterior distribution of  $f(s)$  matches the form of a discriminative softmax classifier (Lee et al., 2018). This equivalence implies that  $f(s)$  fit a Gaussian distribution under  $\pi_\theta$ . We approximate state feature vectors with a class-conditional Gaussian distribution with  $\mu_c$  and  $\Sigma_c$  for each action class, rather than using a single "tied" covariance  $\Sigma$  across all action classes (Kamoi & Kobayashi, 2020).

An MD-based detection based on Gaussian assumption can be immediately developed based on the mean vectors  $\mu_c$  and the covariance matrix  $\Sigma_c$  calculated from  $f(s)$  for each action class  $c$ . We first collect  $N_c$  state action pairs  $\{(s_i, a_i)\}$ , separately for each action class  $c$ , and compute the empirical class mean and covariance of  $c$ :

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i: a_i=c} f(s_i), \quad \hat{\Sigma}_c = \frac{1}{N_c} \sum_{i: a_i=c} (f(s_i) - \hat{\mu}_c)(f(s_i) - \hat{\mu}_c)^\top. \quad (4)$$

In distance-based detection, a straightforward metric is Euclidean distance (ED). However, MD generally outperforms ED in many tasks (Lee et al., 2018; Ren et al., 2021; Kamoi & Kobayashi, 2020), as it incorpo-

<sup>2</sup>For continuous action spaces, we can discretize the actions into several bins and then follow the same detection pipeline.

rates the additional data covariance information to normalize the distance scales. Following the estimation in Eq. 4, we derive the class-conditional Gaussian distribution to characterize the data structure within the state representation space for each action class. For each state  $s$  observed by the agent, we compute its *Detection Mahalanobis Distance*  $M(s)$  between  $s$  and the nearest class-conditional Gaussian distribution by:

$$M(s) = \min_c (f(s) - \hat{\mu}_c)^\top \hat{\Sigma}_c^{-1} (f(s) - \hat{\mu}_c). \quad (5)$$

Unlike the previous work Lee et al. (2018), which defined a Mahalanobis confidence score based on a binary classifier in a validation dataset, we utilize  $M(s)$  as the detection metric within a statistical hypothesis test framework. Proposition 1 demonstrates that  $M(s)$  follows a Chi-squared distribution under the Gaussian assumption.

**Proposition 1.** (*Test Distribution of Detection Mahalanobis distance  $M(s)$* ) Let  $f(\mathbf{s})$  be the  $p$ -dimensional state random vector for action class  $c$ . Under the Gaussian assumption  $P(f(\mathbf{s})|a = c) = \mathcal{N}(f(\mathbf{s}) | \mu_c, \Sigma_c)$ , the Detection Mahalanobis Distance  $M(\mathbf{s})$  in Eq. 5 is Chi-Square distributed:  $M(\mathbf{s}) \sim \chi_p^2$ .

Please refer to Appendix A for the proof. Based on Proposition 1, we can define a threshold  $\Theta = \chi_p^2(1 - \alpha)$  by selecting a  $\alpha$  value from the specified Chi-Squared distribution to distinguish normal states from outliers. Given a new state observation  $s$  and a confidence level  $1 - \alpha$ , if  $M(s) > \Theta$ ,  $s$  is detected as an outlier.

## 4.2 Robust MD-based Detection

**Motivation.** The estimation of  $\mu_c$  and  $\Sigma_c$  in Eq. 4 relies on Maximum Likelihood Estimate (MLE), which is sensitive to the presence of outliers in the dataset (Rousseeuw & Van Zomeren, 1990). As the offline data collected from the environment tends to be noisy, directly introducing MD for outlier detection in RL easily results in a less statistically effective estimation of  $\mu_c$  and  $\Sigma_c$ , thus undermining the detection accuracy for outliers. This vulnerability of the MD-based detector against noisy states prompts us to instantiate MDX with a more robust estimator (Huber, 2004).

To this end, we apply the Minimum Covariance Determinant (MCD) estimator (Hubert & Debruyne, 2010) to estimate  $\mu_c$  and  $\Sigma_c$  by only using a subset of all collected samples. It only uses the observations where the determinant of the covariance matrix is as small as possible. Concretely, MCD determines the subset  $J$  of observations with a size  $h$ , while minimizing the determinant of the sample covariance matrix calculated solely from these  $h$  points. The choice of  $h$  determines the trade-off between the robustness and efficiency of the estimator. The robust MCD mean vector  $\hat{\mu}_c^{\text{rob}}$  and covariance matrix  $\hat{\Sigma}_c^{\text{rob}}$  in the action class  $c$  are computed as

$$\hat{\mu}_c^{\text{rob}} = \frac{1}{h} \sum_{i:i \in J, a_i = c} f(s_i), \quad J = \left\{ \text{set of } h \text{ points} : \left| \hat{\Sigma}_J \right| \leq \left| \hat{\Sigma}_K \right| \text{ for all subsets } K \right\}, \quad (6)$$

where we set  $h$  as  $(\text{number\_of\_samples} + \text{number\_of\_features} + 1)/2$  (Rousseeuw, 1984).  $K$  represents the total number of subsets that contain  $h$  points. In practice, the MCD estimator can be efficiently solved by the FAST-MCD algorithm (Hubert & Debruyne, 2010) instead of performing a brute-force search over

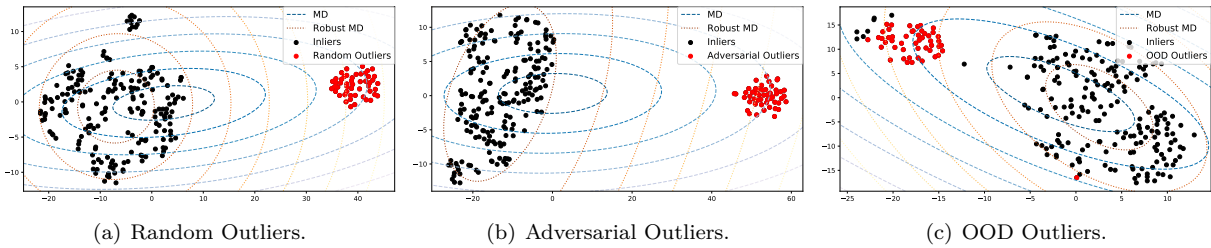


Figure 3: Contours under the estimation based on MD and Robust MD across different outlier types on Breakout. Black and red points denote inliers and outliers, respectively. The dimension of state feature vectors after a pre-trained PPO policy is reduced by t-SNE (Van der Maaten & Hinton, 2008).

all possible subsets. Akin to Mahalanobis Distance, we define the *Detection Robust Mahalanobis Distance*  $M_{\text{rob}}(s)$  as robust detection metric:

$$M_{\text{rob}}(s) = \min_c (f(s) - \hat{\mu}_c^{\text{rob}})^\top \hat{\Sigma}_c^{\text{rob}-1} (f(s) - \hat{\mu}_c^{\text{rob}}). \quad (7)$$

Since the robust Mahalanobis distance can still approximate the true Chi-squared distribution (Hardin & Rocke, 2005), we still employ the threshold value  $\Theta = \chi_p^2(1 - \alpha)$  for detecting outliers as in the MD case.

**Potential Advantages of Robust MD on Real Data.** As a motivating example, Figure 3 displays contours computed by both MD and Robust MD detection methods for state feature vectors in the Breakout game from the popular Atari benchmark (Bellemare et al., 2013; Brockman et al., 2016) with different types of outliers. These results demonstrate that estimation based on Robust MD is less vulnerable to outlying states (red points) and better fits inliers (black points) than MD. This robust parameter estimation highlights the potential advantage of Robust MD for RL outlier detection, where the data used for estimation tends to be noisy.

### 4.3 Conformal MD-based Detection

**Motivation.** Although robust MD-based detection is less vulnerable to noise in RL environments, both MD and robust MD strategies heavily rely on the Gaussian assumption to construct the detection thresholds based on Proposition 1. This distribution assumption is often violated in practice, diminishing the effectiveness of MD and robust MD. In contrast, conformal prediction offers a mathematical framework that provides valid and rigorous prediction distribution without assuming a specific underlying data distribution. The resulting conformal anomaly detection circumvents the limitation of the distribution assumption, potentially improving the detection efficacy.

In the context of RL, conformal anomaly detection evaluates how a state conforms to a model’s current prediction distribution, thereby discriminating abnormal states. As a distribution-free detection approach, conformal anomaly detection can enhance the distance-based detectors by additionally tuning the anomaly threshold in the calibration dataset. To design the conformal anomaly detection method, we leverage the Detection Mahalanobis Distance  $M(s)$  as the *non-conformity score*, which measures how dissimilar a state is from the instances in the calibration set. Following split conformal inference (Papadopoulos et al., 2002; Shafer & Vovk, 2008), we split the the previously collected offline dataset into the the calibration set  $\mathcal{D}_{\text{cal}}$  and the evaluation set. A simple way is to evaluate the quantiles of the resulting empirical distribution to create the corresponding confidence band. Using the calibration set  $\mathcal{D}_{\text{cal}}$ , we define the fitted quantiles  $\hat{Q}_{1-\alpha}^c$  of the conformity scores for the action class  $c$  as follows:

$$\hat{Q}_{1-\alpha}^c = \inf \left\{ q : \left( \frac{1}{N_c} \sum_{s_i \in \mathcal{D}_{\text{cal}}, a_i = c} \mathbf{1}_{\{M^c(s_i) \leq q\}} \right) \geq 1 - \alpha \right\}, \quad (8)$$

where each  $(s_i, a_i)$  is drawn from the calibration set  $\mathcal{D}_{\text{cal}}$  and  $c$  is calculated by  $c = \arg \min M^c(s_i)$  in  $M^c(s_i)$  among all action classes. Finally, we use the class-dependent and well-calibrated detection thresholding  $\Theta = \hat{Q}_{1-\alpha}^c$  in conformal MD-based detection instead of  $\chi_p^2(1 - \alpha)$  used in MD and Robust MD strategies.

### 4.4 MD-based Detection Algorithm in the Offline Setting

Algorithm 1 summarizes the instantiation of MDX in the offline setting. We compute the (robust) mean vector and covariance matrix among the state feature vectors in the penultimate layer of  $\pi_\theta$  for each action class. Next, given a state observation  $s$ , we compute the detection Mahalanobis distance  $d = M(s)$  or  $d = M_{\text{rob}}(s)$ . And compare it with the threshold  $\Theta = \chi_p^2(1 - \alpha)$  under the Gaussian assumption or  $\Theta = \hat{Q}_{1-\alpha}^c$  from distribution-free conformal quantiles. If  $d > \Theta$ ,  $s$  is detected as an outlier. Conversely, if  $d \leq \Theta$ ,  $s$  is deemed as an inlier.

**Algorithm 1** MDX Detection Framework in the Offline Setting

---

```

1: Input: The given policy  $\pi_\theta$ , the dimension of state feature vectors  $p$ , and a confidence level  $1 - \alpha$ .
2: Output: Detection labels  $\{y_s\}$  for each  $s$  in the evaluation trajectory.
3: / * Step 1: Detection Design by Estimating Mean and Covariance * /
4: Given state action pairs  $\{(s_i, a_i)\}$  where  $a_i \sim \pi_\theta(\cdot|s_i)$ .
5: for each action class  $c$  do
6:   if we choose MD detection then
7:     Estimate  $\hat{\mu}_c$  and  $\hat{\Sigma}_c$  via Eq. (4). / * Approach 1: MD Detection * /
8:   else if we choose Robust MD detection then
9:     Estimate  $\hat{\mu}_c^{\text{rob}}$  and  $\hat{\Sigma}_c^{\text{rob}}$  via Eqs. (6) and (7). / * Approach 2: Robust MD Detection * /
10:  else
11:    Estimate  $\hat{\mu}_c, \hat{\Sigma}_c$  via Eq. (4), calibrate  $\hat{Q}_{1-\alpha}^c$  via Eq. (8) / * Approach 3: Conformal MD-based Detection * /
12:  end if
13: end for
14: / * Step 2: Detection Deployment * /
15: for  $s$  in the noisy environment do
16:   Compute distance  $d = M(s)$  or  $d = M_{\text{rob}}(s)$ , and threshold  $\Theta = \chi_p^2(1 - \alpha)$  or  $\Theta = \hat{Q}_{1-\alpha}^c$ .
17:   Set Detection label  $y_s = 1$  if  $d > \Theta$  else  $y_s = -1$ .
18: end for

```

---

## 5 Anomaly Detection in the Online RL Setting

In the online RL setting (Sutton & Barto, 2018; Dong et al., 2020), a policy is updated continuously, unlike the fixed pre-trained policy used in our offline setting. Robust policy training with noisy states is crucial in safe RL, as the agents are more likely to encounter state outliers during training. In this section, we extend MDX to the online RL training scenario. Unlike the offline setting, the challenge here stems from the dynamic nature of policy updates, requiring our detector to adapt to the evolving distribution of feature vector outputs. The complexity increases when the improved policy starts gathering new samples through exploration, posing a fundamental challenge in an online RL framework. An effective detection system must differentiate between actual noisy observations and newly collected data through exploration. Training the RL agent and estimating the detector are interleaved in a noisy online environment. Various options for managing detected outliers during training include removing or denoising the outlier states. In our detection framework, we focus on direct removal and assess the resulting learning curves in the presence of noisy states during the training process. To address the challenges in detecting abnormal states in the online training setting, we propose *Moving Window Estimation* and *Double Self-supervised Detectors*, both of which are pivotal for the empirical success of our anomaly detection approach.

**Moving Window Estimation.** In the online setting, improving the policy  $\pi_\theta$  causes a shift in the data distribution within the replay buffer as the agent interacts with the environment (Rolnick et al., 2019; Xiao et al., 2019). To effectively utilize information from the updated data distribution, we maintain a moving window to store experiences throughout the interaction steps. The window size can be adjusted to either prioritize a long historical context with a larger window size or more recent experiences with a smaller one. Based on the constantly updated state feature vectors,  $\mu_c$  and  $\Sigma_c$  ( $\mu_c^{\text{rob}}$  and  $\Sigma_c^{\text{rob}}$ ) are continually estimated. This continuous updating allows us to accurately track the state feature distribution, ensuring that our detector remains sensitive to recent and historical data shifts.

**Double Self-Supervised Detectors.** Our current detector is continually refined using self-detected inliers, while any detected outliers are promptly discarded. However, a more practical approach is to leverage these outliers to create a complementary detector for outliers. This secondary self-supervised detector validates the detection results from the primary detector. For example, if the primary detector classifies a state as an inlier and the secondary detector agrees that it is not an outlier, the state is confidently classified as such. Conversely, if there is a difference between the discrimination of the two detectors, the state is randomly classified as either an outlier or an inlier. In the event of disagreement, this random classification



**Algorithm 2** MDX Detection Framework in the Online Setting, PPO Style

---

```

1: Initialize policy network  $\pi_\theta$  and estimator  $\hat{\mu}_c$  and  $\hat{\Sigma}_c$  (or  $\hat{\mu}_c^{\text{rob}}$  and  $\hat{\Sigma}_c^{\text{rob}}$ ).
2: Initialize confidence level  $1 - \alpha$ , the window size  $m$ , inlier and outlier buffers  $\mathcal{B}_I$ ,  $\mathcal{B}_O$ .
3: for iteration = 1, 2, ...,  $K$  do
4:   for actor = 1, 2, ...,  $N$  do
5:     Run policy  $\pi_\theta$  in environment for  $T$  timesteps.
6:     Compute distance  $d = M(s)$  or  $d = M_{\text{rob}}(s)$ , and threshold  $\Theta = \chi_p^2(1 - \alpha)$  or  $\Theta = \hat{Q}_{1-\alpha}^c$ .
7:     if  $d \leq \Theta$  then
8:       Add it to  $\mathcal{B}_I$ .
9:     else
10:      Add it to  $\mathcal{B}_O$ .
11:    end if
12:  end for
13:  Optimize policy  $\pi_\theta$  using inlier trajectories.
14:  Update  $\hat{\mu}_c$  and  $\hat{\Sigma}_c$  (or  $\hat{\mu}_c^{\text{rob}}$  and  $\hat{\Sigma}_c^{\text{rob}}$ ) of the two detectors based on  $\mathcal{B}_I$  and  $\mathcal{B}_O$  respectively every  $N_c$  samples.
15: end for

```

---

is motivated by the need to avoid systematic bias that could arise from consistently favoring one detector’s output over the other. By introducing randomness, we ensure the system remains fair and does not overly rely on potentially flawed outputs from either detector. This approach also preserves the system’s ability to learn and adapt over time, preventing the reinforcement of incorrect classifications. The double-detector system thus enhances the robustness and reliability of the detection process, ensuring more accurate and consistent identification of abnormal states.

**MD-based Detection Algorithm in the Online Setting.** Algorithm 2 outlines our MD-based detection procedure for online RL, incorporating both moving window estimation and double self-supervised detectors. To update our double detectors, inliers and outliers are stored in buffers  $\mathcal{B}_I$  and  $\mathcal{B}_O$ , respectively. For each class, a window size  $m$  is specified. Within each class, the state-action pairs in the window are used to estimate  $\hat{\mu}_c$  and  $\hat{\Sigma}_c$  ( $\hat{\mu}_c^{\text{rob}}$  and  $\hat{\Sigma}_c^{\text{rob}}$ ). These parameters are updated after every  $N_c$  newly collected data points in the window for action class  $c$ . This adaptive updating mechanism ensures that the detectors remain responsive to evolving data distributions.

**Online Anomaly Detection Procedure.** In a real-time scenario like a recommendation system, we typically first deploy a pre-trained policy as a warm start in the online system to provide initial recommendations for each user. Feedback from users, such as the click-through rate (CTR), is then observed to update the online policy iteratively. Similarly, within our online detection algorithm, we pre-train a policy as a warm start. After pre-training, the policy is introduced to the noisy environment for further online learning. Throughout this process, our MDX framework is used to identify outliers in the subsequent training phases. We then evaluate the training performance of algorithms equipped with these detection mechanisms. This systematic approach facilitates the gradual refinement of the policy while concurrently integrating outlier detection to enhance robustness in real-world settings.

## 6 Experiments

We first conduct experiments on six typical Atari games (Bellemare et al., 2013) to verify the effectiveness of our MDX framework in both offline and online settings. The Atari games are divided into two different groups. The first group includes Breakout, Asterix, and SpaceInvaders, which feature nearly static backgrounds. Enduro, FishingDerby, and Tutankham in the second group have time-changing or dramatically different backgrounds, presenting more challenging scenarios. We further conduct experiments on autonomous driving environments (Dosovitskiy et al., 2017) as one potential application. We select Proximal Policy Optimization (PPO) (Schulman et al., 2017) as our baseline RL algorithm.

**Three Types of Outliers.** (1) **Random Outliers.** We generate random outliers by adding Gaussian noise with zero mean and different standard deviations on state observations, simulating natural measurement errors. (2) **Adversarial Outliers.** We perform white-box adversarial perturbations (Goodfellow et al., 2014b; Szegedy et al., 2013; Cao et al., 2020) on state observations for the current policy, following the strategy proposed in (Huang et al., 2017; Pattanaik et al., 2017). Particularly, we denote  $a_w^t$  as the "worst" action, with the lowest probability from the current policy  $\pi_t(a|s)$ . The optimal adversarial perturbation  $\eta_t$ , constrained in an  $\epsilon$ -ball, can be derived by minimizing the objective function  $J: \min_{\eta} J(s_t + \eta, \pi_t) = -\sum_{i=1}^n p_i^t \log \pi_t(a_i|s_t + \eta), s.t. \|\eta\| \leq \epsilon$ , where  $p_w^t = 1$  and  $p_i^t = 0$  for  $i \neq w$ . We solve this minimization problem with the Fast Gradient Sign Method (FGSM) (Goodfellow et al., 2014b), a typical adversarial attack method in the deep learning literature. The resulting adversarial outliers  $s_t + \eta_t^*$  force the policy to choose  $a_w^t$ . (3) **Out-of-Distribution (OOD) outliers.** OOD outliers arise from the disparity in data distribution across different environments. To simulate them, we randomly select states from other environments and introduce them to the current environment. In our experiments, we select images from other Atari games to serve as Out-of-Distribution (OOD) outliers within the considered environment. In the autonomous driving scenario, we designate rainy and nighttime observations as OOD outliers for the primary daytime setting on a sunny day. This deliberate selection of diverse outlier examples enables comprehensive testing of our method’s robustness across varied environments.

**Baseline Methods.** A fundamental obstacle in assessing the anomaly detection strategies in RL lies in the scarcity of baselines in deep RL settings as introduced in Section 2. To rigorously substantiate the effectiveness of MDX, we initiate our evaluation by comparing them with the foundational baselines we have developed ourselves. (1) **Euclidean distance (ED)** assumes that all features are independent under the Gaussian assumption with one standard deviation, which can be considered as a simplified version of our MD method with an identity covariance matrix. (2) **MD with Tied covariance (TMD)** follows the tied covariance assumption in (Lee et al., 2018), where features among all action classes share a single covariance matrix estimation. (3) **MD** is our first proposed method with class-conditional Gaussian assumption. (4) **Robust MD (RMD)** is the robust variant of MD under the Gaussian assumption. (5) **MD+C** uses well-

Detection Accuracy (%)	Outliers	ED	TMD	MD	RMD	MD+C
Breakout	Random	53.2	60.0	64.0	<b>71.2</b>	62.8
	Adversarial	83.8	89.1	91.0	80.4	<b>92.3</b>
	OOD	50.0	47.6	50.5	<b>78.7</b>	51.5
Asterix	Random	44.3	46.0	59.6	<b>71.2</b>	54.8
	Adversarial	84.2	85.5	91.3	75.8	<b>93.7</b>
	OOD	40.1	40.8	45.9	<b>65.2</b>	49.7
SpaceInvader	Random	52.1	66.2	72.3	<b>79.2</b>	70.4
	Adversarial	72.4	91.2	95.9	83.4	<b>96.4</b>
	OOD	45.2	56.6	51.4	<b>83.2</b>	50.2
Enduro	Random	49.0	51.6	60.2	<b>78.5</b>	51.6
	Adversarial	93.9	90.8	95.2	80.4	<b>97.5</b>
	OOD	57.0	62.8	69.8	<b>80.3</b>	53.2
FishingDerby	Random	49.1	66.3	69.2	<b>85.6</b>	65.3
	Adversarial	85.3	92.9	<b>97.5</b>	87.4	97.4
	OOD	51.1	55.9	59.2	<b>75.7</b>	57.9
Tutankham	Random	50.0	47.9	49.2	<b>77.0</b>	52.2
	Adversarial	61.1	88.9	<b>94.2</b>	78.7	87.2
	OOD	55.0	86.6	<b>92.0</b>	78.7	77.8
Average	Random	49.6	56.3	62.4	<b>77.1</b>	59.5
	Adversarial	80.1	89.7	<b>94.2</b>	81.0	94.0
	OOD	49.7	58.4	61.5	<b>76.9</b>	56.7
	Average	59.8	68.1	72.7	<b>78.3</b>	70.1

Table 1: **Average detection accuracy** of MD, RMD, and MD+C compared with baselines across different outlier types in all six environments in the **offline** setting. The averages are computed across environments and outlier types. Accuracy is determined by applying detection techniques to the balanced data composed equally of clean and noisy states.

calibrated conformality scores to construct a valid empirical distance distribution instead of relying on the Chi-Squared distribution established upon the Gaussian assumption.

### 6.1 Anomaly Detection in the Offline Setting

In the offline setting, we randomly split the states from the given dataset into calibration and evaluation sets, each containing 50%. The calibration set is used to construct our detectors, and the evaluation set is for testing. We first use PCA to reduce the state feature vectors into a 50-dimensional space. We then apply MD or robust MD to estimate mean vectors and covariances and calibrate the conformality score based on the calibration dataset. Finally, we add the three types of noises to the evaluation dataset and combine them with the clean evaluation dataset. We assess the performance of our detection methods on the entire evaluation dataset.

**Main results.** Table 1 shows the detection accuracy of MDX instantiated with MD, robust MD, and conformal MD with  $\alpha = 0.05$  across a wide range of outlier types on each game. A higher accuracy indicates a more successful identification of anomalies for the evaluated detection method. Detailed results for settings with different noise types are provided in Table 4 of Appendix B.1. We conclude that: (1) All MD-based methods, i.e., TMD, MD, RMD, and MD+C, outperform ED, confirming the usefulness of covariance matrix information in RL outlier detection. (2) Robust MD generally performs the best, significantly surpassing MD and other methods in detecting random and OOD outliers. Nonetheless, robust MD is not effective enough to detect adversarial outliers. (3) MD+C excels in identifying adversarial outliers and performs similarly to MD in other scenarios.

**Sensitivity Analysis on Feature Dimension Reduction.** We provide a sensitivity analysis regarding the number of feature dimensions reduced by PCA, showing that the detection accuracy for all considered outliers tends to improve as the number of principal components increases. This indicates that better detection performance can be achieved with higher feature dimensions. The detailed results are presented in Appendix B.4.

**Effectiveness of Robust MD.** In robust MD analysis, it is typically concluded that outlier states are more distinctly separated from inlier states. By comparing the Mahalanobis distance distributions between inliers and outliers under both MD and Robust MD, we show that this conclusion also applies to the RL anomaly detection scenario. This effect explains the detection advantage of robust MD in RL. Detailed results are provided in Appendix B.3.

### 6.2 Anomaly Detection in the Online Setting

The PPO agent, utilizing multi-processes as detailed in the original PPO algorithm (Schulman et al., 2017), runs eight independent environments in parallel, and we introduce state outliers into four of these environments. For random and adversarial outliers, actions are determined based on the PPO policy network  $\pi_\theta$ . For OOD outliers, due to the potential differences in action spaces between the original environment and the OOD environment, we select OOD states from the OOD environment by taking random actions within its own action space. For the Robust MD method, we use PCA to reduce state feature vectors into a 50-dimensional space due to the expensive computation of the robust MD method. For the other methods, we use the original feature vectors output from the penultimate layer of  $\pi_\theta$ . Results are averaged over three seeds with hyperparameters given in Table 5 of Appendix C.1. When our detectors identify an outlier, it is removed from training. We compare the resulting learning curves for different detection methods.

**Additional Baselines.** We add another two baselines as performance upper bound and lower bound. (1) For an ideal baseline, the method **Auto** automatically deletes true state outliers, showing the optimal training performance of algorithms *without the interruption from outliers*. (2) At the other extreme, **Random** uses a totally random detector that detects a state as an inlier or outlier with a probability of 0.5.

**Main Results.** Figure 4 presents learning curves of cumulative rewards (first row) based on the PPO algorithm and the corresponding detection F1 Score (second row) for all tested detection methods across three types of outliers in Tutankham games. To better highlight their differences, we omit the confidence bands in Figure 4, while providing similar results on all six Atari games with confidence bands in Appendix C.1

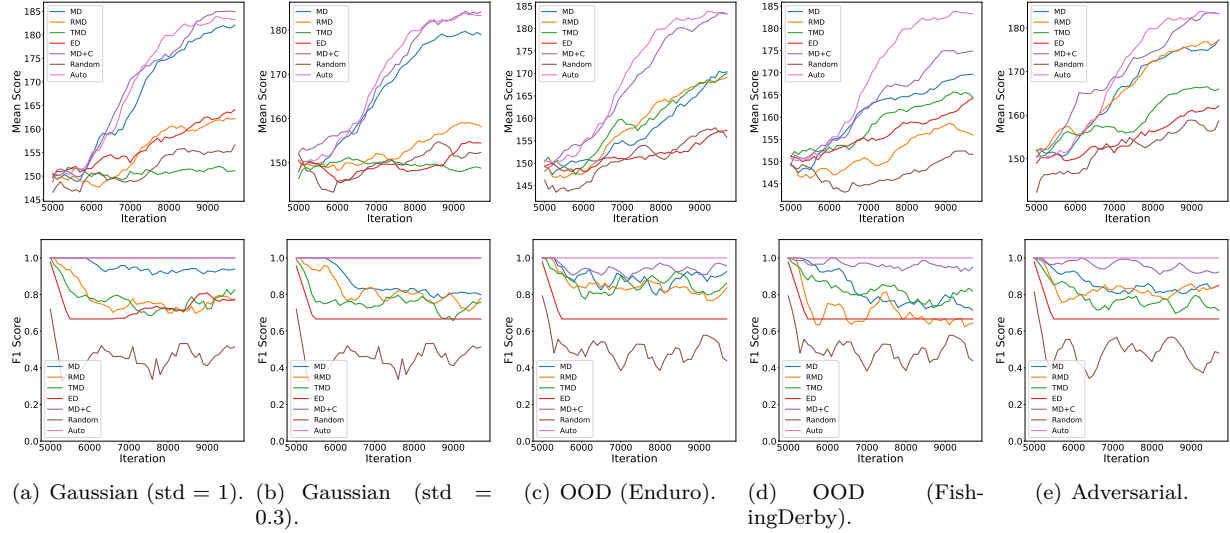


Figure 4: **Learning curves and detection performance** across various state outliers in online learning on Tutankham. "Mean Score" in the first row indicates the cumulative rewards, and "F1 Score" in the second row evaluates the detection performance during training. We present the average results over three random seeds while omitting confidence bands for a clearer comparison.

from Figures 13 to 18 for reference. For each outlier type in Table 2, we evaluate the **superiority rank** of all detectors regarding the F1 score and policy performance, where rank 1 indicates the best performance. A smaller superiority rank implies a more effective detection. Our conclusions are as follows: (1) Conformal MD (MD+C) generally achieves the best detection performance across all considered baselines (except Auto). The superiority of MD+C over MD highlights the crucial role of accurately calibrated thresholds in the online RL detection setting. (2) RMD is less effective than MD and MD+C, performing only on par with TMD and ED. This degradation is due to *the information loss in the dimension reduction of the feature vectors for reducing the computational cost*. Thus, MD+C and MD are preferable to RMD in the computationally demanding online RL setting. (3) The average superiority ranks of all considered detectors are similar in terms of performance and F1 score, verifying the consistency of our results.

**Ablation Study on Double Self-Supervised Detectors.** We conduct an ablation study of double self-supervised detectors on Breakout with random and OOD outliers. Results in Figure 19 of Appendix C.2 show that double self-supervised detectors reduce detection errors and improve detection accuracy.

**Ablation Study on Outlier Proportions.** We also demonstrate the robust detection performance across different proportions of outliers encountered by the agent during training. We conduct experiments on Breakout, and the results are provided in Figure 20 of Appendix C.3.

Superiority Rank	Outlier Type	Random	ED	TMD	MD	RMD	MD+C
Performance	Random	5.7	2.9	4.1	2.5	3.9	<b>1.9</b>
	OOD	5.8	4.5	2.0	3.2	4.1	<b>1.4</b>
	Adversarial	5.7	3.3	4.5	3.0	3.2	<b>1.3</b>
Average	All	5.7	3.6	3.4	2.8	3.8	<b>1.6</b>
F1 Score	Random	6.0	3.1	4.0	3.3	3.4	<b>1.2</b>
	OOD	6.0	4.8	2.0	2.5	4.1	<b>1.6</b>
	Adversarial	6.0	4.2	3.3	2.8	3.7	<b>1.0</b>
Average	All	6.0	3.9	3.1	2.9	3.7	<b>1.3</b>

Table 2: The average superiority rank (1 is best) of anomaly detection methods across all types of outliers in all six environments. Numbers in bold represent the best method.



Figure 5: The clean and noisy state observations in autonomous driving experiments.

### 6.3 Autonomous Driving Environment

To verify the broader applicability of our method, we perform experiments on autonomous driving environments and introduce practical scenarios in which all three types of anomalies commonly occur.

**Random Noise.** Malfunctioning sensors or cameras can introduce random noise into signal observations. For instance, a faulty camera lens may produce distorted images, while a malfunctioning LiDAR sensor might generate erroneous depth measurements. Such random noise can impair the reliability of perception systems in autonomous vehicles.

**Adversarial Attacks.** Adversarial attacks involve intentionally manipulating input signals to disrupt the functioning of RL systems. In the context of autonomous driving, an attacker might tamper with sensor data or traffic signs, resulting in misleading observations and potentially hazardous driving behavior. Adversarial states thus pose a significant threat to the robustness and safety of autonomous driving systems.

**Out-of-Distribution (OOD) States.** Consider a scenario where an RL policy is trained exclusively under sunny weather. Encountering rainy weather poses a challenge, as the observations captured under these conditions deviate from the training data distribution. Such observations are therefore considered Out-of-Distribution (OOD) states.

**Experimental Setup.** We conduct experiments using the CARLA environment (Dosovitskiy et al., 2017). CARLA is an open-source simulator for autonomous driving research known for its high-quality rendering and realistic physics. The environment includes 3D models of static objects, such as buildings, vegetation, traffic signs, and infrastructure, as well as dynamic objects, such as vehicles and pedestrians. The task is to drive safely through the town. In each episode, the vehicle must reach a given goal without collision. The episode ends when the vehicle reaches the goal, collides with an obstacle, or exceeds the time limit.

**Noisy State Observations.** Following the approach used in Atari game settings, we introduce Gaussian noise to simulate random outliers and generate adversarial outliers using adversarial perturbations. For OOD outliers, we leverage CycleGAN-Turbo (Zhu et al., 2017; Parmar et al., 2024), a technique designed for adapting a single-step diffusion model (Ho et al., 2020) to new tasks and domains through adversarial learning (Goodfellow et al., 2014a). This method can perform various image-to-image translation tasks and outperforms existing GAN-based and diffusion-based methods for various scene translation tasks, such as day-to-night conversion and adding/removing weather effects like fog, snow, and rain (Parmar et al., 2024). Specifically, we use CycleGAN-Turbo to create **rainy** and **nighttime** outliers. Examples of different anomaly states are presented in Figure 5.

**Main Results.** Given a fixed dataset and a pre-trained policy, we assess our detection methods across the three types of outliers. Table 3 shows the average accuracy, with MD+C achieving the highest performance

Detection Accuracy (%)	ED	TMD	MD	RMD	MD+C
Random (std $\in [0.005, 0.06]$ )	50.0	60.8	69.8	<b>72.1</b>	61.7
Random (std $\in (0.06, 0.3]$ )	50.0	95.0	95.2	73.6	<b>96.0</b>
Adversarial	50.0	96.7	95.3	73.8	<b>97.5</b>
OOD (Rain)	50.0	96.5	95.5	74.4	<b>97.5</b>
OOD (Night)	50.0	96.5	95.5	74.3	<b>97.5</b>

Table 3: Detection accuracy on the CARLA *town* environment over three types of outliers.

in most scenarios, while RMD performs best in the presence of small random noises. These results suggest that our proposed method effectively detects outliers for realistic problems, such as autonomous driving.

## 7 Discussions and Conclusion

In this paper, we present the first detailed study of a distance-based anomaly detection framework in deep RL, considering random, adversarial, and OOD state outliers in both offline and online settings. The primary detection backbone is based on Mahalanobis distance, and we extend it to robust and distribution-free versions by leveraging robust estimation and conformal prediction techniques. Experiments on Atari games and the autonomous driving environment demonstrate the effectiveness of our proposed methods in detecting the three types of outliers. The conformal MD method achieves the best detection performance in most scenarios, especially in the online setting. Our research contributes to developing safe and trustworthy RL systems in real-world applications.

**Limitations and Future Work.** In the online setting, especially with a high proportion of outliers, it may be preferable to denoise the detected state outliers via some neighboring smoothing techniques, e.g., *mixup* (Wang et al., 2020; Zhang et al., 2018), rather than deleting them directly as performed in this paper. To relax the Gaussian assumption in the hypothesis test of our detection, we can consider other non-parametric methods, such as one-class support vector machines (Choi, 2009) or isolation forests (Liu et al., 2008). A substantial challenge that remains for future work is to devise a more informed detector to distinguish between real “bad” outliers that can cause truly misleading actions and “good” *new samples* collected through exploration, which can potentially benefit the policy learning, especially for image inputs (Zhang & Ranganath, 2023).

## References

- Anastasios Nikolas Angelopoulos, Stephen Bates, Michael Jordan, and Jitendra Malik. Uncertainty sets for image classifiers using conformal prediction. In *International Conference on Learning Representations*, 2021.
- Rina Foygel Barber, Emmanuel J Candes, Aaditya Ramdas, and Ryan J Tibshirani. Conformal prediction beyond exchangeability. *The Annals of Statistics*, 51(2):816–845, 2023.
- M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The Arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, June 2013.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI Gym, 2016. URL <https://arxiv.org/abs/1606.01540>.
- RW Butler, PL Davies, and M Jhun. Asymptotics for the minimum covariance determinant estimator. *The Annals of Statistics*, pp. 1385–1400, 1993.
- Yuanjiang Cao, Xiaocong Chen, Lina Yao, Xianzhi Wang, and Wei Emma Zhang. Adversarial attacks and detection on reinforcement learning-based interactive recommender systems. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 1669–1672, 2020.
- Stephanie C.Y. Chan, Samuel Fishman, Anoop Korattikara, John Canny, and Sergio Guadarrama. Measuring the reliability of reinforcement learning algorithms. In *International Conference on Learning Representations*, 2020.
- Shang-Tse Chen, Cory Cornelius, Jason Martin, and Duen Horng Chau. ShapeShifter: Robust physical adversarial attack on faster R-CNN object detector. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2018, Dublin, Ireland, September 10–14, 2018, Proceedings, Part I 18*, pp. 52–68. Springer, 2019.
- Young-Sik Choi. Least squares one-class support vector machine. *Pattern Recognition Letters*, 30(13):1236–1240, 2009.

- Mohamad H Danesh and Alan Fern. Out-of-distribution dynamics detection: RL-relevant benchmarks and results. *ICML 2021 Workshop on Uncertainty and Robustness in Deep Learning*, 2021.
- Roy De Maesschalck, Delphine Jouan-Rimbaud, and Désiré L Massart. The Mahalanobis distance. *Chemo-metrics and Intelligent Laboratory Systems*, 50(1):1–18, 2000.
- Hao Dong, Zihan Ding, and Shanghang Zhang. *Deep Reinforcement Learning: Fundamentals, Research and Applications*. Springer Nature, 2020.
- Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *Conference on Robot Learning*, pp. 1–16. PMLR, 2017.
- Nebrase Elmrabit, Feixiang Zhou, Fengyin Li, and Huiyu Zhou. Evaluation of machine learning algorithms for anomaly detection. In *2020 international conference on cyber security and protection of digital services (cyber security)*, pp. 1–8. IEEE, 2020.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014a.
- Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *International Conference on Learning Representations (ICLR)*, 2014b.
- R Grübel. A minimal characterization of the covariance matrix. *Metrika*, 35(1):49–52, 1988.
- Shangding Gu, Long Yang, Yali Du, Guang Chen, Florian Walter, Jun Wang, and Alois Knoll. A review of safe reinforcement learning: Methods, theory and applications, 2024. URL <https://arxiv.org/abs/2205.10330>.
- Tom Haider, Karsten Roscher, Felipe Schmoeller da Roza, and Stephan Günnemann. Out-of-distribution detection for reinforcement learning agents with probabilistic dynamics models. In *Proceedings of the 2023 International Conference on Autonomous Agents and Multiagent Systems*, pp. 851–859, 2023.
- Johanna Hardin and David M Rocke. The distribution of robust distances. *Journal of Computational and Graphical Statistics*, 14(4):928–946, 2005.
- Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- Waleed Hilal, S Andrew Gadsden, and John Yawney. Financial fraud: a review of anomaly detection techniques and recent advances. *Expert systems With applications*, 193:116429, 2022.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Liyuan Hu, Mengbing Li, Chengchun Shi, Zhenke Wu, and Piotr Fryzlewicz. Doubly inhomogeneous reinforcement learning, 2022. URL <https://arxiv.org/abs/2211.03983>.
- Sandy Huang, Nicolas Papernot, Ian Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial attacks on neural network policies. *International Conference on Learning Representations (ICLR) workshop*, 2017.
- Peter J Huber. *Robust statistics*, volume 523. John Wiley & Sons, 2004.
- Mia Hubert and Michiel Debruyne. Minimum covariance determinant. *Wiley interdisciplinary reviews: Computational statistics*, 2(1):36–43, 2010.
- Vladislav Ishimtsev, Alexander Bernstein, Evgeny Burnaev, and Ivan Nazarov. Conformal  $k$ -nn anomaly detector for univariate data streams. In *Conformal and Probabilistic Prediction and Applications*, pp. 213–227. PMLR, 2017.

- Scott Jordan, Yash Chandak, Daniel Cohen, Mengxue Zhang, and Philip Thomas. Evaluating the performance of reinforcement learning algorithms. In *International Conference on Machine Learning*, pp. 4962–4973. PMLR, 2020.
- Ryo Kamoi and Kei Kobayashi. Why is the mahalanobis distance effective for anomaly detection?, 2020. URL <https://arxiv.org/abs/2003.00402>.
- Kiarash Kazari, Ezzeldin Shereen, and György Dán. Decentralized anomaly detection in cooperative multi-agent reinforcement learning. In *IJCAI*, pp. 162–170, 2023.
- B Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A Al Sallab, Senthil Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(6):4909–4926, 2021.
- William R Klecka, Gudmund R Iversen, and William R Klecka. *Discriminant analysis*, volume 19. Sage, 1980.
- Rikard Laxhammar and Göran Falkman. Sequential conformal anomaly detection in trajectories based on Hausdorff distance. In *14th international conference on information fusion*, pp. 1–8. IEEE, 2011.
- Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in Neural Information Processing Systems*, 31, 2018.
- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems, 2020. URL <https://arxiv.org/abs/2005.01643>.
- Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *International Conference on Learning Representations (ICLR)*, 2016.
- Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *2008 Eighth IEEE International Conference on Data Mining*, pp. 413–422. IEEE, 2008.
- Ricardo A Maronna and Víctor J Yohai. Robust estimation of multivariate location and scatter. *Wiley StatsRef: Statistics Reference Online*, pp. 1–12, 2014.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928–1937. PMLR, 2016.
- Robert Müller, Steffen Illium, Thomy Phan, Tom Haider, and Claudia Linnhoff-Popien. Towards anomaly detection in reinforcement learning. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*, pp. 1799–1803, 2022.
- Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. Deep learning for anomaly detection: A review. *ACM Computing Surveys (CSUR)*, 54(2):1–38, 2021.
- Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alex Gammerman. Inductive confidence machines for regression. In *Machine learning: ECML 2002: 13th European conference on machine learning Helsinki, Finland, August 19–23, 2002 proceedings 13*, pp. 345–356. Springer, 2002.
- Gaurav Parmar, Taesung Park, Srinivasa Narasimhan, and Jun-Yan Zhu. One-step image translation with text-to-image models, 2024. URL <https://arxiv.org/abs/2403.12036>.



- Anay Pattanaik, Zhenyi Tang, Shuijing Liu, Gautham Bommanan, and Girish Chowdhary. Robust deep reinforcement learning with adversarial attacks. *Advances in Neural Information Processing Systems*, 2017.
- Andrew Patterson, Samuel Neumann, Martha White, and Adam White. Draft: Empirical design in reinforcement learning. *Journal of Artificial Intelligence Research*, 1, 2020.
- Jie Ren, Stanislav Fort, Jeremiah Liu, Abhijit Guha Roy, Shreyas Padhy, and Balaji Lakshminarayanan. A simple fix to Mahalanobis distance for improving near-OOD detection. *International Conference on Machine Learning (ICML) workshop*, 2021.
- David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Peter J Rousseeuw. Least median of squares regression. *Journal of the American Statistical Association*, 79(388):871–880, 1984.
- Peter J Rousseeuw and Bert C Van Zomeren. Unmasking multivariate outliers and leverage points. *Journal of the American Statistical Association*, 85(411):633–639, 1990.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *Proceedings of The 32nd International Conference on Machine Learning*, pp. 1889–1897. PMLR, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *International Conference on Learning Representations*, 2018.
- Andreas Sedlmeier, Thomas Gabor, Thomy Phan, Lenz Belzner, and Claudia Linnhoff-Popien. Uncertainty-based out-of-distribution classification in deep reinforcement learning. In *Proceedings of the 12th International Conference on Agents and Artificial Intelligence*. SCITEPRESS - Science and Technology Publications, 2020. doi: 10.5220/0008949905220529.
- Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(3), 2008.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *International Conference on Learning Representations (ICLR)*, 2013.
- Jiaye Teng, Chuan Wen, Dinghuai Zhang, Yoshua Bengio, Yang Gao, and Yang Yuan. Predictive inference with feature conformal prediction. In *The Eleventh International Conference on Learning Representations*, 2023.
- Ryan J Tibshirani, Rina Foygel Barber, Emmanuel Candes, and Aaditya Ramdas. Conformal prediction under covariate shift. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033, 2012. doi: 10.1109/IROS.2012.6386109.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.

- Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. *Algorithmic learning in a random world*, volume 29. Springer, 2005.
- Chen Wang, Sarah Erfani, Tansu Alpcan, and Christopher Leckie. Oil-ad: An anomaly detection framework for sequential decision sequences, 2024. URL <https://arxiv.org/abs/2402.04567>.
- Kaixin Wang, Bingyi Kang, Jie Shao, and Jiashi Feng. Improving generalization in reinforcement learning with mixture regularization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 7968–7978. Curran Associates, Inc., 2020.
- Edwin B Wilson and Margaret M Hilferty. The distribution of chi-square. *Proceedings of the National Academy of Sciences of the United States of America*, 17(12):684, 1931.
- Wei Xiao, Xiaolin Huang, Fan He, Jorge Silva, Saba Emrani, and Arin Chaudhuri. Online robust principal component analysis with change point detection. *IEEE Transactions on Multimedia*, 22(1):59–68, 2019.
- Hongming Zhang and Tianyang Yu. Taxonomy of reinforcement learning algorithms. *Deep reinforcement learning: Fundamentals, research and applications*, pp. 125–133, 2020.
- Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.
- Lily H Zhang and Rajesh Ranganath. Robustness to spurious correlations improves semantic out-of-distribution detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 15305–15312, 2023.
- Yingying Zhang, Chengchun Shi, and Shikai Luo. Conformal off-policy prediction. In *International Conference on Artificial Intelligence and Statistics*, pp. 2751–2768. PMLR, 2023.
- Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pp. 2223–2232, 2017.

## A Proof of Proposition 1

*Proof.* We show that for each action class  $c$ , the square of Mahalanobis distance  $d$  is identically independent Chi-squared distributed under the Gaussian assumption. Without loss of generality, we denote  $\mu$  and  $\Sigma$  as the mean and variance matrix of the closest class-conditional Gaussian distribution. We need to show  $d = (f(\mathbf{s}) - \mu)^\top \Sigma^{-1} (f(\mathbf{s}) - \mu)$  is Chi-squared distributed. Firstly, by eigenvalue decomposition, we have

$$\Sigma^{-1} = \sum_{k=1}^p \lambda_k^{-1} u_k u_k^\top, \quad (9)$$

where  $\lambda_k$  and  $u_k$  are the  $k$ -th eigenvalue and eigenvector of  $\Sigma$ . Plugging it into the form of  $d$ , we immediately obtain

$$\begin{aligned} d &= (f(\mathbf{s}) - \mu)^\top \Sigma^{-1} (f(\mathbf{s}) - \mu) \\ &= (f(\mathbf{s}) - \mu)^\top \left( \sum_{k=1}^p \lambda_k^{-1} u_k u_k^\top \right) (f(\mathbf{s}) - \mu) \\ &= \sum_{k=1}^p \lambda_k^{-1} (f(\mathbf{s}) - \mu)^\top u_k u_k^\top (f(\mathbf{s}) - \mu) \\ &= \sum_{k=1}^p \left[ \lambda_k^{-\frac{1}{2}} u_k^\top (f(\mathbf{s}) - \mu) \right]^2 \\ &= \sum_{k=1}^p \mathbf{X}_k^2, \end{aligned} \quad (10)$$

where  $\mathbf{X}_k^2$  is a new Gaussian variable that results from the linear transform of a Gaussian distribution  $f(\mathbf{s})$  where  $f(\mathbf{s}) \sim \mathcal{N}(\mu, \Sigma)$ . Therefore, the resulting variance  $\sigma_k^2$  can be derived as

$$\sigma_k^2 = \lambda_k^{-\frac{1}{2}} u_k^\top \Sigma \lambda_k^{-\frac{1}{2}} u_k = \lambda_k^{-1} u_k^\top \left( \sum_{j=1}^p \lambda_j u_j u_j^\top \right) u_k = \sum_{j=1}^p \lambda_k^{-1} \lambda_j u_k^\top u_j u_j^\top u_k \quad (11)$$

As the  $\mu_j$  and  $\mu_k$  are orthogonal if  $j \neq k$ , the variance  $\sigma_k^2$  can be further reduced to

$$\sigma_k^2 = \lambda_k^{-1} \lambda_k u_k^\top u_k u_k^\top u_k = \|u_k\|^2 \|u_k\|^2 = 1. \quad (12)$$

Each  $\mathbf{X}_k$  is a standard Gaussian distribution. Then we have  $d$ , the square of Mahalanobis distance, Chi-squared distributed, i.e.,  $d \sim \chi^2(p)$ , independent of the action class  $c$ . Without loss of generality, the smallest  $d$  among all action classes, i.e.,  $M(\mathbf{s})$ , is also a Chi-squared distribution. That is to say,  $M(\mathbf{s}) \sim \chi^2(p)$ .  $\square$

## B Results in Offline Setting

### B.1 Results across Different Noise Strengths

We provide detailed detection accuracy of various detection methods across different noise strengths in Table 4.

### B.2 Visualization of Outlier States on Six Games

We plot the outlier states on Breakout, Asterix, and SpaceInvaders games in Figure 6 and outliers states on Enduro, FishingDerby, and Tutankham in Figure 7.

Table 4: Detection accuracy (%) of our MD, Robust MD, and conformal MD strategies compared with other baseline methods on six Atari games with  $\alpha = 0.05$ .

Games	Outliers	Perturbation	ED	TMD	MD	RMD	MD+C	Games	Outliers	Perturbation	ED	TMD	MD	RMD	MD+C
Breakout	Random	std=0.02	50.0	52.2	56.5	<b>66.1</b>	55.5	Enduro	Random	std=0.1	49.4	44.7	48.2	<b>76.6</b>	48.6
		std=0.04	56.4	67.8	71.4	<b>76.3</b>	70.0			std=0.2	48.6	58.4	72.2	<b>80.3</b>	54.5
	Adversarial	$\epsilon=0.001$	80.4	87.4	89.4	80.0	<b>90.0</b>		Adversarial	$\epsilon=0.001$	93.1	90.8	95.2	80.5	<b>97.4</b>
		$\epsilon=0.01$	87.2	90.7	92.5	80.7	<b>94.5</b>			$\epsilon=0.01$	94.6	90.8	95.2	80.3	<b>97.5</b>
	OOD	Asterix	53.5	47.7	51.2	<b>81.3</b>	51.8		OOD	FishingDerby	63.7	72.9	78.5	<b>80.3</b>	53.9
		SpaceInvaders	46.5	47.4	49.8	<b>76.1</b>	51.2			Tutankham	50.2	52.6	61.0	<b>80.2</b>	52.5
Asterix	Random	std=0.1	42.8	42.9	48.0	<b>66.2</b>	49.1	FishingDerby	Random	std=0.2	49.0	48.7	51.9	<b>83.5</b>	51.0
		std=0.2	45.8	49.1	71.1	<b>76.1</b>	60.4			std=0.3	49.1	83.8	86.4	<b>87.6</b>	79.5
	Adversarial	$\epsilon=0.001$	83.9	85.2	91.1	75.7	<b>93.3</b>		Adversarial	$\epsilon=0.001$	82.4	92.9	<b>97.5</b>	87.3	97.3
		$\epsilon=0.01$	84.5	85.8	91.5	75.9	<b>94.0</b>			$\epsilon=0.01$	88.2	92.9	<b>97.5</b>	87.5	97.4
	OOD	Breakout	42.3	43.0	48.1	<b>76.1</b>	51.9		OOD	Enduro	49.0	56.7	60.6	<b>75.6</b>	59.9
		SpaceInvaders	37.9	38.6	43.7	<b>54.2</b>	47.5			Tutankham	53.1	55.1	57.8	<b>75.8</b>	55.8
SpaceInvaders	Random	std=0.02	50.1	49.7	55.1	<b>74.7</b>	53.8	Tutankham	Random	std=0.04	50.0	48.5	49.3	<b>75.7</b>	51.6
		std=0.04	54.0	82.6	<b>89.5</b>	83.6	86.9			std=0.06	50.0	47.2	49.0	<b>78.3</b>	52.7
	Adversarial	$\epsilon=0.001$	68.6	90.5	95.5	83.2	<b>95.9</b>		Adversarial	$\epsilon=0.001$	60.0	88.3	<b>93.7</b>	78.8	81.2
		$\epsilon=0.01$	76.2	91.8	96.3	83.6	<b>96.8</b>			$\epsilon=0.01$	62.1	89.5	<b>94.7</b>	78.6	93.1
	OOD	Breakout	45.7	52.7	52.6	<b>82.9</b>	50.6		OOD	Enduro	60.0	89.6	<b>94.7</b>	78.9	90.9
		Asterix	44.7	60.4	50.2	<b>83.4</b>	49.7			FishingDerby	50.0	83.5	<b>89.2</b>	78.4	64.7

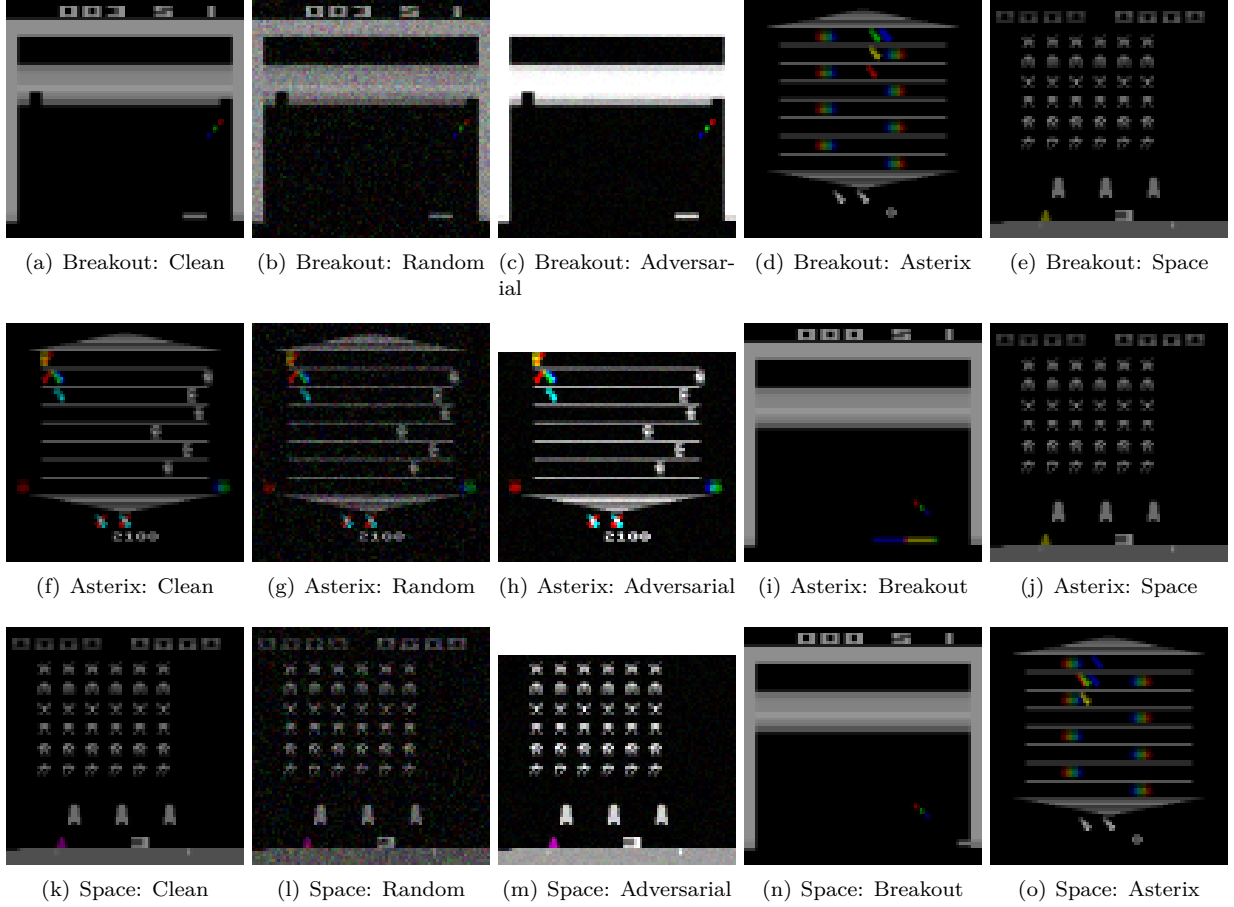


Figure 6: Visualization of various state outliers on Breakout, Asterix, and SpaceInvaders games.

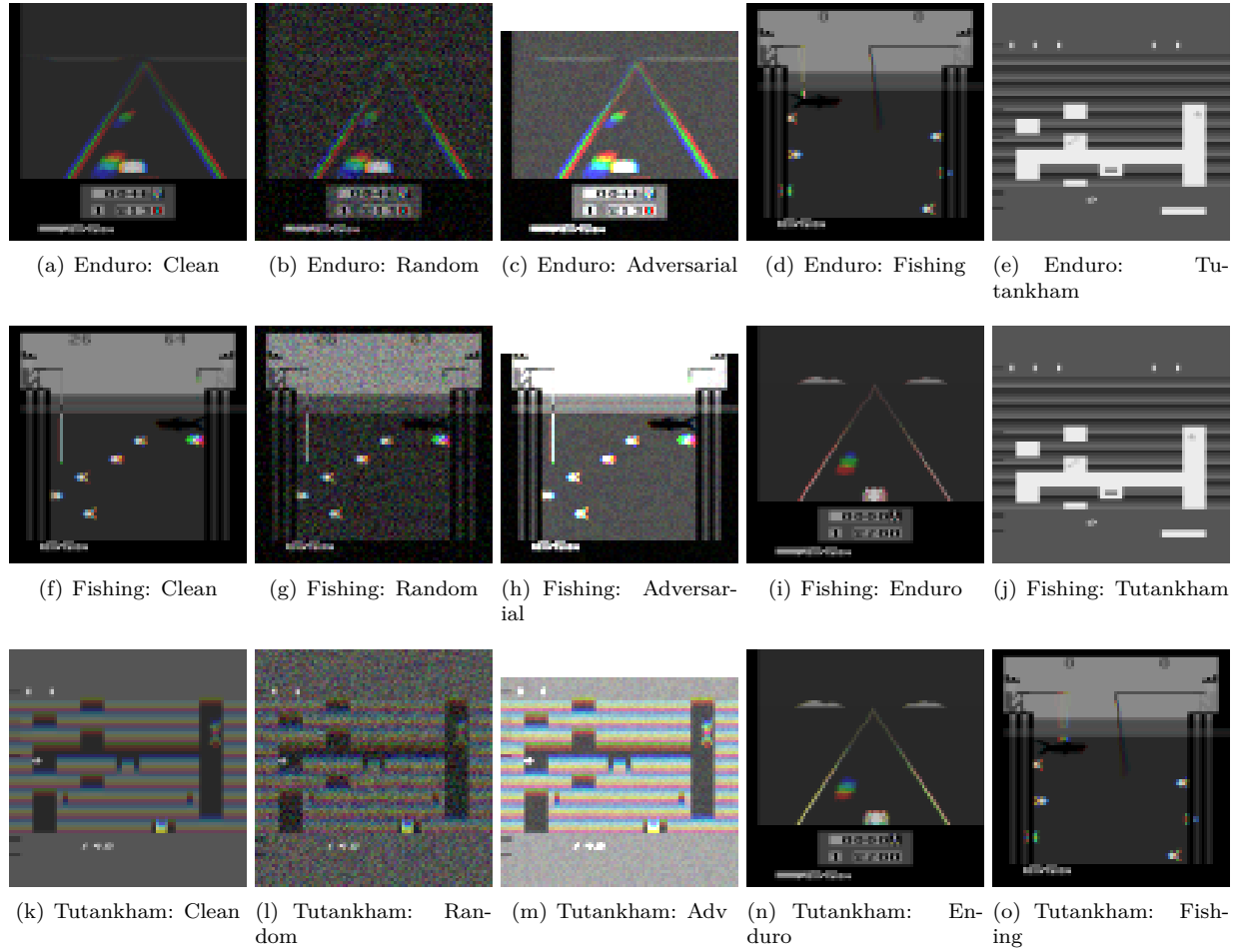


Figure 7: Visualization of various state outliers on Enduro, FishingDerby, and Tutankham games.

### B.3 Effectiveness of Robust MD

We take the cubic root of the Mahalanobis distances, yielding approximately normal distributions (Wilson & Hilferty, 1931). In this experiment, 250 clean states are drawn from the replay buffer, and 50 abnormal states are drawn from each of the three types of outliers. We reduce the state feature dimension to 2 via t-SNE and compute Mahalanobis distances of these two kinds of states to their centrality within each action class under the estimation based on MD or Robust MD, respectively. Figure 8 suggests that Robust MD separates inliers and outliers better than MD on Breakout within a random action class, indicating its effectiveness in detecting RL evaluation. Similar results are also given in other games.

We plot the distributions of inliers and three types of outliers on SpaceInvaders and Asterix games in Figure 9 and 10, respectively. It is worth noting that Robust MD is also capable of enlarging the separation of distributions between inliers and both random and adversarial outliers on SpaceInvaders game, while its benefit seems to be negligible on OOD outliers (Breakout) on SpaceInvaders games as well as in Asterix game. We speculate that it is determined by the game’s difficulty. Specifically, the PPO algorithm can

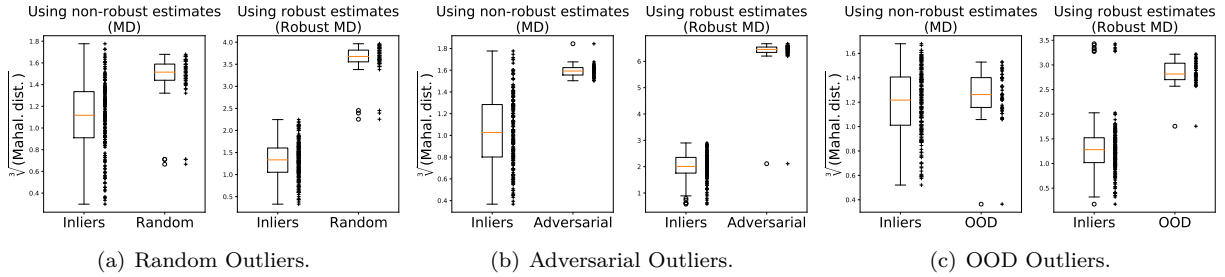


Figure 8: Boxplot of distributions between inliers and three types of outliers in an action class on Breakout game.

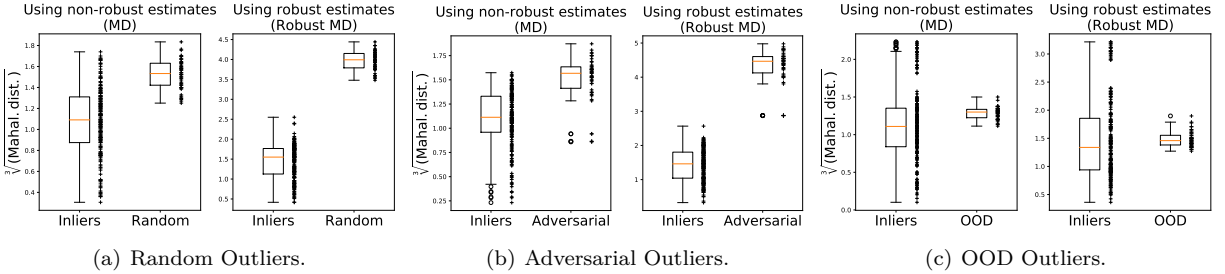


Figure 9: Boxplot of distributions between inliers and three types of outliers in an action class on SpaceInvaders game.

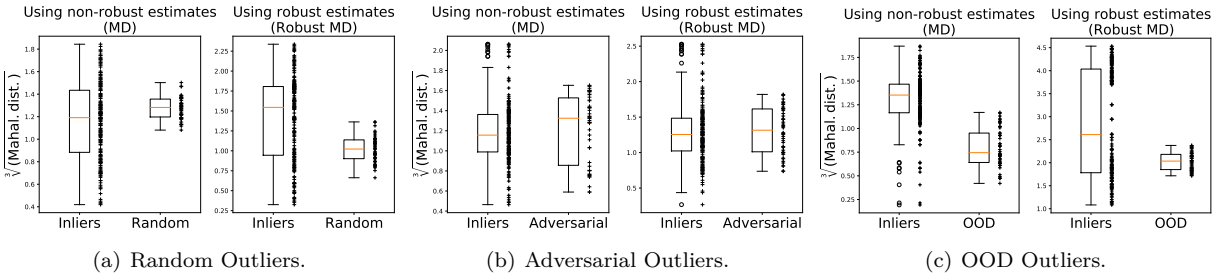


Figure 10: Boxplot of distributions between inliers and three types of outliers in an action class on Asterix game.

achieve desirable performance on the simple Breakout game, thus yielding informative feature space vectors. By contrast, there is room for the generalization of PPO on both SpaceInvaders and Asterix games, such that Robust MD might not help when handling the less meaningful state feature vectors in these two games.

#### B.4 Sensitivity Analysis

We provide the sensitivity analysis of Robust MD in terms of the PCA dimension in Figure 11. The impact of the number of principal components on the detection performance for robust MD detection is shown in Figure 11. The detection accuracy over all considered outliers improves as the number of principal components increases, except for a slight decline for random and adversarial outliers (red and blue lines) on the Breakout game. The increase implies that the subspace spanned by principal components with small explained variance also contains valuable information for detecting anomalous states from in-distribution states, which coincides with the conclusion in (Kamoi & Kobayashi, 2020).

The result of MD estimation manifests in Figure 12. It suggests that there is still an ascending tendency of detection accuracy as the number of principal components increases.

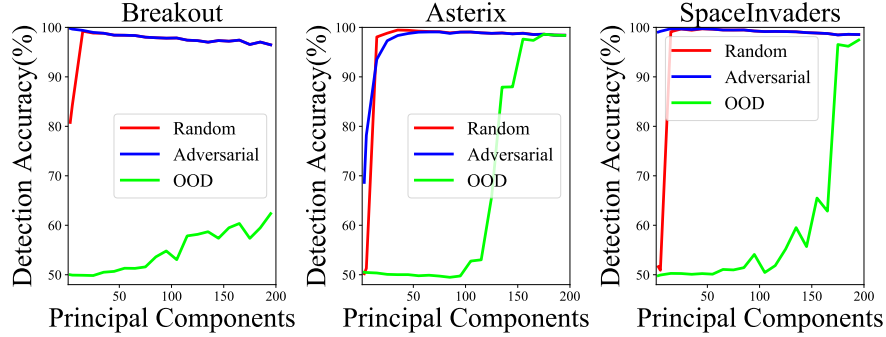


Figure 11: Detection performance under **Robust MD** as the number of principal components increases.

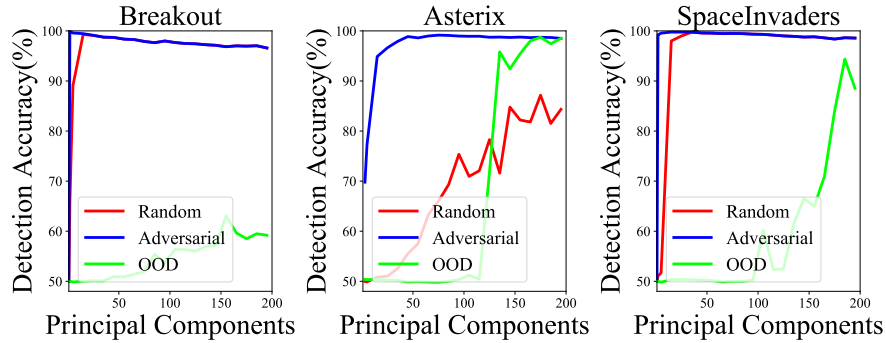


Figure 12: Detection performance under **MD** as the number of principal components increases.

## C Results in Online Setting

### C.1 Setup and Full Main Results

As a supplement to the results on the main pages, we provide the whole results on all six Atari games from Figure 13 to Figure 18. The "Mean Score" in the first row indicates the accumulated rewards of PPO, and the "F1 Score" in the second row shows the detection performance during RL training. The F1 score is computed based on precision and recall. We also find that the cumulative reward is not strongly correlated with detection ability in some games. A high detection accuracy may only improve the cumulative reward to a small degree. This suggests that we need more metrics to measure the effect of our detection performance more effectively. Hyperparameters in our methods are shown in Table 5.

Hyperparameter	Value
Confidence level ( $1-\alpha$ )	1-0.05
Moving window size ( $m$ )	5120
Sample size ( $N_c$ )	2560
Iteration ( $K$ )	$\approx 10000$ (1e7 steps in total)
Environment number ( $N$ )	8
Horizon ( $T$ )	128

Table 5: Hyper-parameters in the training phase. RL-related parameters are the same as those of the PPO algorithm.

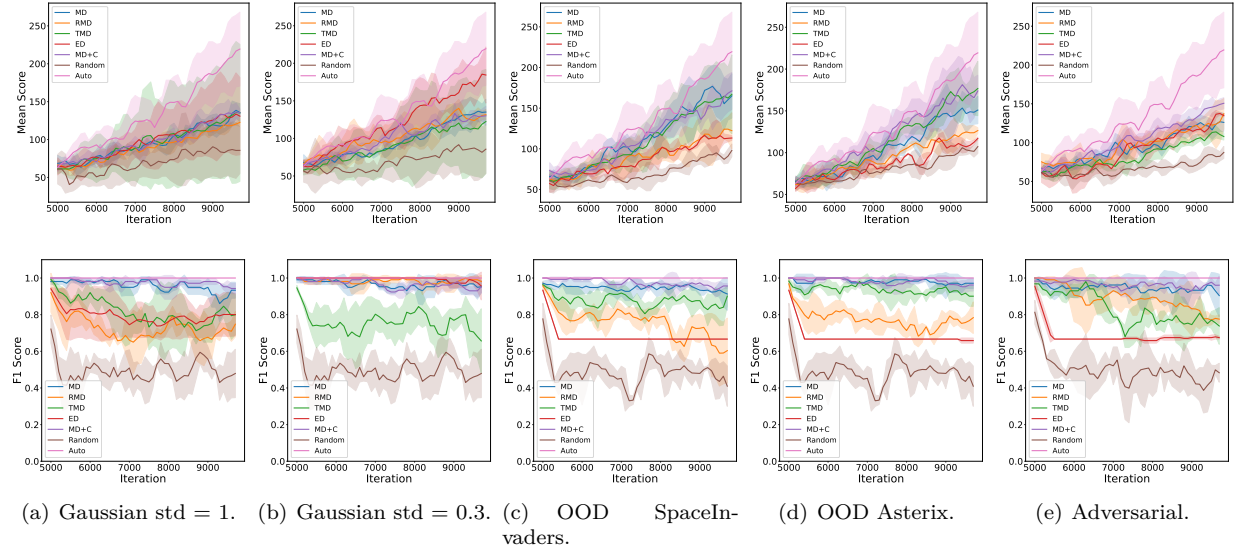


Figure 13: Detection performance across various state outliers in the online training on Breakout. "Mean Score" in the first row indicates the accumulated rewards, "accuracy" and "F1 Score" evaluate the detection performance during training.



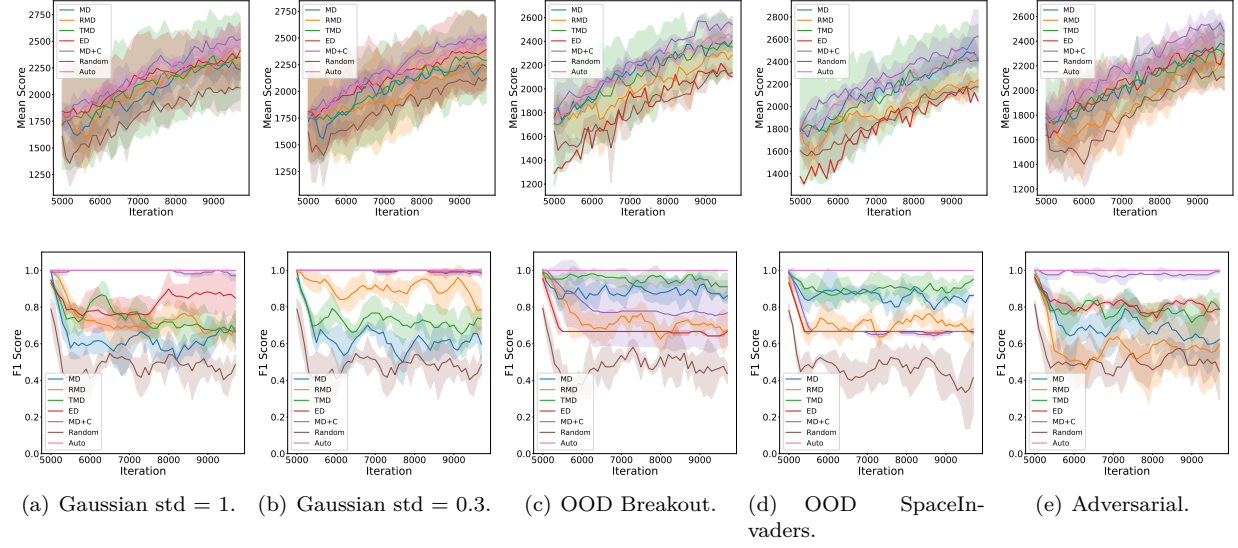


Figure 14: Detection performance across various state outliers in the online training on Asterix.

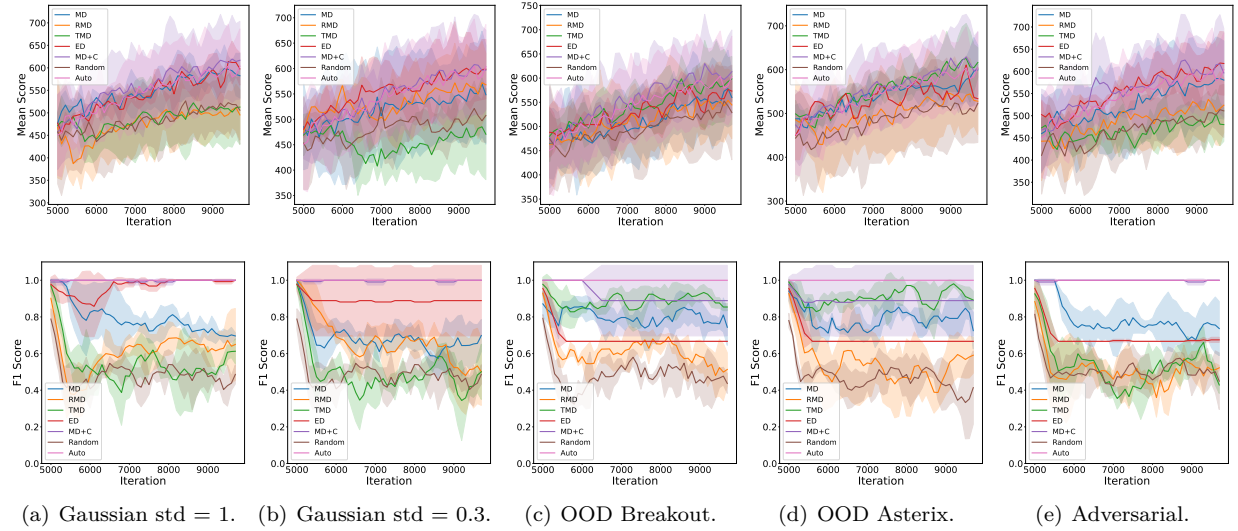


Figure 15: Detection performance across various state outliers in the online training on SpaceInvaders.

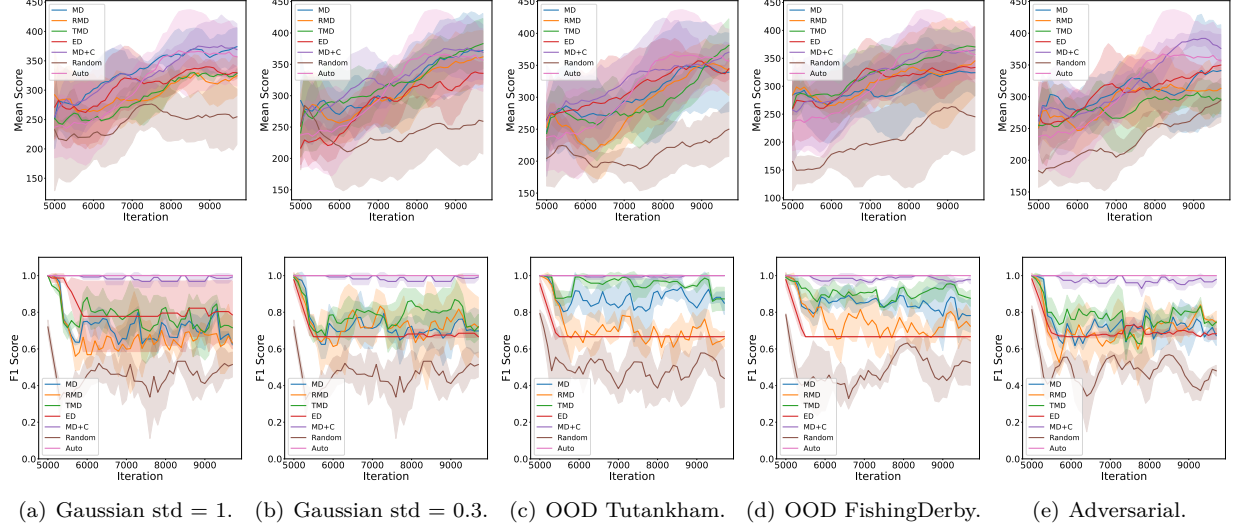


Figure 16: Detection performance across various state outliers in the training phase on Enduro.

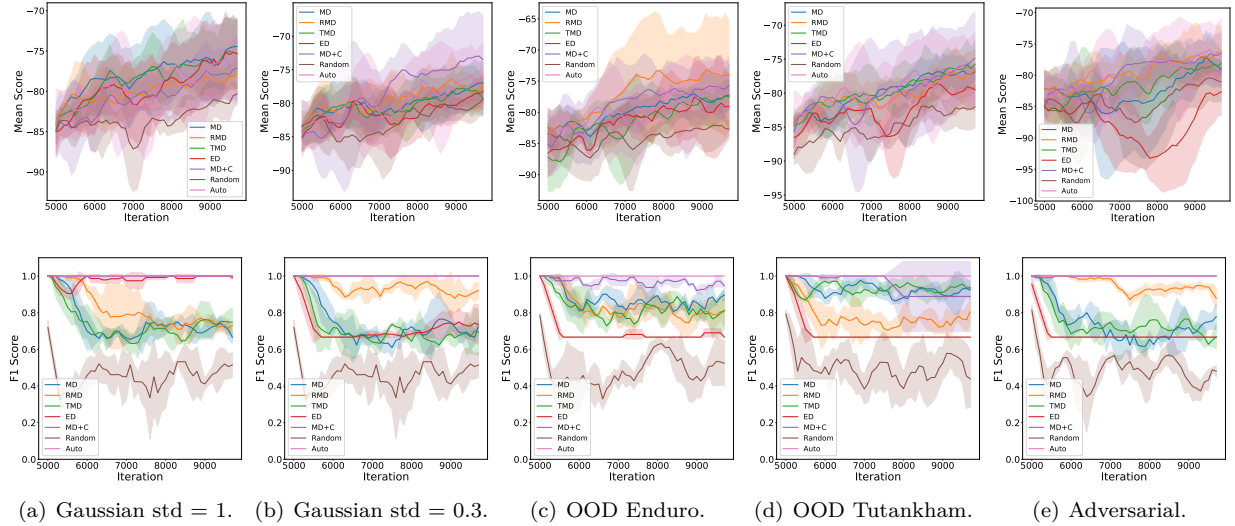


Figure 17: Detection performance across various state outliers in the online training on FishingDerby.

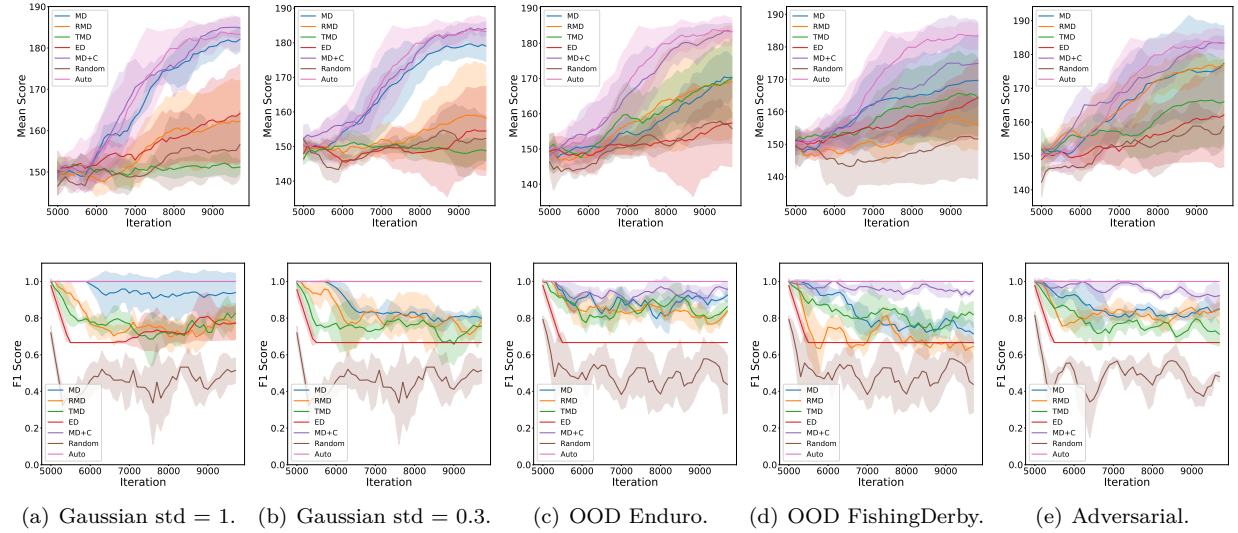


Figure 18: Detection performance across various state outliers in the online training on Tutankham.

## C.2 Ablation Study on Double Anomaly Detectors

Figure 19 reveals that double self-supervised detectors can help adjust the detection errors and improve the detection accuracy compared with the single detector. MD with double detectors outperforms MD with a single detector significantly, although RMD with double detectors is comparable to RMD with a single detector.

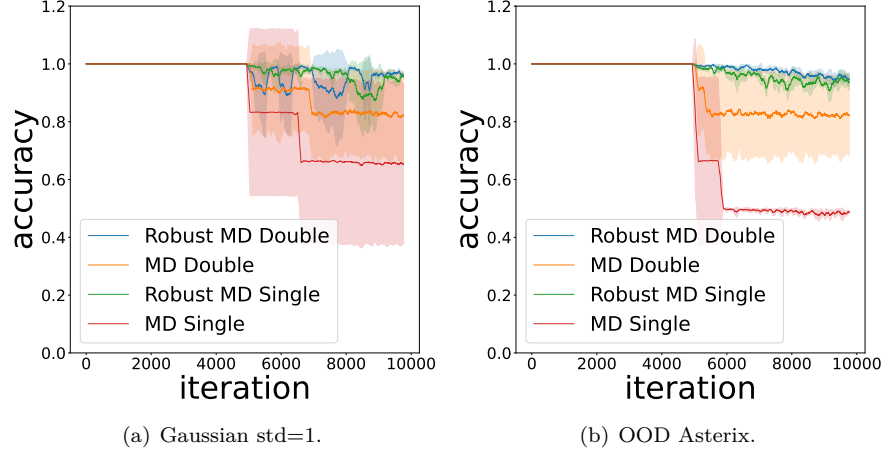


Figure 19: The detection accuracy with and without double self-supervised detectors on Breakout with random and OOD outliers on Breakout.

## C.3 Ablation Study on Number of Noisy Environments

We train PPO in two, four, or six noisy environments with random and OOD outliers among all eight parallel environments. We use PCA to reduce the feature vectors to 50 dimensions and estimate the detector using Robust MD. Figure 20 illustrates that compared with the **Auto** baseline, our RMD method is robust when encountering different ratios of outliers, especially with a higher contamination ratio. The dashed lines in different colors represent **Auto** baselines that correspond to the different number of noisy environments. The training performance with our detection method gradually approaches the ideal baselines, i.e., **Auto**.

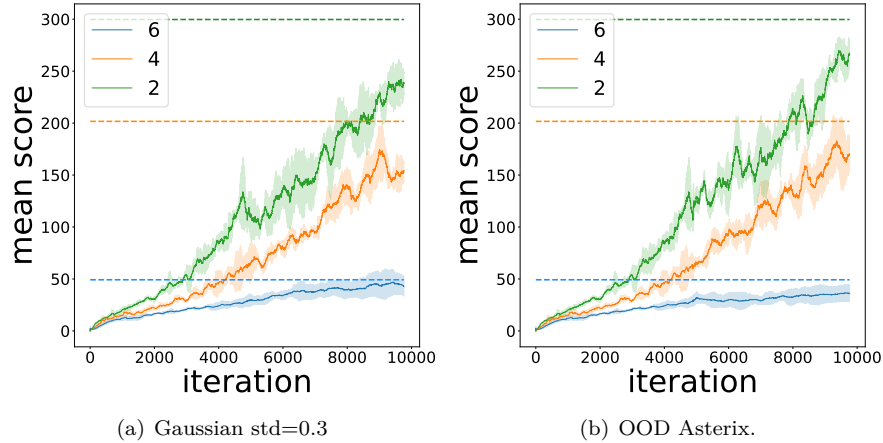


Figure 20: Training performance under Robust MD detection under different proportions of outlier exposure on Breakout (2, 4, 6 out of 8 environments).