
Human-in-the-Loop Out-of-Distribution Detection with False Positive Rate Control

Harit Vishwakarma^{*1} Heguang Lin^{*1} Ramya Korlakai Vinayak²

Abstract

Robustness to Out-of-Distribution (OOD) samples is essential for the successful deployment of machine learning models in the open world. Since it is not possible to have a priori access to a variety of OOD data before deployment, several recent works have focused on designing scoring functions to quantify OOD uncertainty. These methods often find a threshold that achieves 95% true positive rate (TPR) on the In-Distribution (ID) data used for training and uses this threshold for detecting OOD samples. However, this can lead to very high FPR as seen in a comprehensive evaluation in the Open-OOD benchmark, the FPR can range between 60 to 96% on several ID and OOD dataset combinations. In contrast, practical systems deal with a variety of OOD samples on the fly and critical applications, e.g., medical diagnosis, demanding guaranteed control of the false positive rate (FPR). To meet these challenges, we propose a mathematically grounded framework for human-in-the-loop OOD detection, wherein expert feedback is used to update the threshold. This allows the system to adapt to variations in the OOD data while adhering to the quality constraints. We propose an algorithm that uses any time-valid confidence intervals based on the Law of Iterated Logarithm (LIL). Our theoretical results show that the system meets FPR constraints while minimizing the human feedback for points that are in-distribution. Another key feature of the system is that it can work with any existing post-hoc OOD uncertainty-quantification methods. We evaluate our system empirically on a mixture of benchmark OOD datasets in image classification tasks on CIFAR-10 and CIFAR-100 as in distribu-

tion datasets and show that our method can maintain FPR at most 5% while maximizing TPR.

1. Introduction

Deploying machine learning (ML) models in the open world makes it subject to out-of-distributions (OOD) inputs — in the classification setup OOD data points are those that do not belong to any of the classes in the training data. The modern ML models, in particular deep neural networks, can fail silently with high confidence on OOD points (Nguyen et al., 2015; Amodei et al., 2016) rather than flagging them as OOD and asking for human intervention as they are not designed to do so. Such failures can have serious consequences in high-risk applications e.g. medical diagnosis, autonomous driving, etc. For a successful deployment of an ML model in the open world, we need mechanisms to ensure robustness to the OOD inputs.

Many recent works have addressed this problem by proposing post-hoc methods for OOD detection (Liang et al., 2017; Lee et al., 2018; Liu et al., 2020; Ming et al., 2022). Broadly, these works propose methods to quantify a *score* that can be used to decide OOD vs ID label for a given point. Many of these methods are based on the distance between data points or a model’s confidence score in prediction. For a detailed survey of literature in the area of generalized OOD detection, see (Yang et al., 2021b). However, many of these works are largely limited to static settings where the ID data which is available in plenty for training and validating the ML system is used to set a threshold on the scores used for OOD detection (Liang et al., 2017; Liu et al., 2020; Ming et al., 2022). This is usually done by setting a threshold that achieves a certain level of true positive rate (TPR), e.g., TPR of 95%. However, this can lead to a very high false positive rate (FPR) e.g., ranging between 60% to 90% on several benchmarked ID and OOD dataset combinations (Yang et al., 2022). Furthermore, even if the ID data distribution remains the same after deployment, the OOD data could vary, resulting in highly fluctuating FPR. Thus, having a small and fixed amount of OOD data collected a priori to validate the FPR at a given threshold would not help in guaranteeing FPR.

In critical applications, the consequences of classifying an

^{*}Equal contribution ¹Dept. of Computer Sciences, University of Wisconsin-Madison, WI, USA ²Dept. of Electrical and Computer Engineering, University of Wisconsin-Madison, WI, USA. Correspondence to: Harit Vishwakarma <hvishwakarma@cs.wisc.edu>.

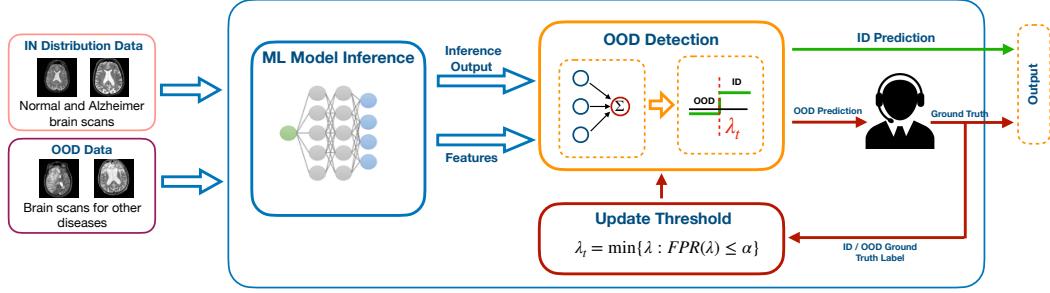


Figure 1. Illustration of OOD detection with human-in-the-loop with FPR control. In this example, the ID data is of brain scans of normal people and those with Alzheimer’s disease. The OOD data could be anything else, e.g. brain scans of patients with some other diseases.

OOD point as ID (false positive) could be more catastrophic than classifying an ID point as OOD (false negative), e.g. in medical diagnosis, when in doubt it is better to classify a brain scan as OOD and defer the decision to human experts rather than for the ML model to give it a disease label or classify as normal assuming it to be ID. Therefore, it is crucial to guarantee that the false positive rate (FPR) is below a certain acceptable rate, e.g., FPR below 5%. Since the availability of exact type of OOD data that the system can encounter during deployment is rare, it is crucial to design systems that can adapt to the OOD data while controlling the FPR during deployment.

Goal: Develop human-in-the-loop out-of-distribution detection system that has guaranteed false positive rate control while minimizing the amount of human intervention needed.

In this work we present a framework for human-in-the-loop Out-of-Distribution (OOD) detection, ensuring strict control over the false positive rate (FPR) while adapting to diverse OOD data.

Our Contributions: Toward this goal, we make the following contributions:

1. **Human-in-the-loop OOD detection framework:** We propose a novel mathematically grounded framework that incorporates expert human feedback to adaptively update the OOD detection threshold, ensuring robustness to variations in OOD data encountered after deployment. Our framework can be used with any scoring function.
2. **Guaranteed FPR control:** Our approach leverages mathematically grounded confidence intervals based on the Law of Iterated Logarithm to meet false positive rate (FPR) constraints while minimizing the need for human feedback on in-distribution points. For stationary settings, we provide theoretical guarantees for our proposed framework and algorithm on controlling FPR at the desired level at all times and also provide a bound on the time taken to reach a given level of optimality. Using the insight from this analysis, we also propose an ap-

proach for settings with change points that reduces the duration of violation of FPR control before adapting to the change.

3. Empirical validation on benchmark datasets: We evaluate our framework through extensive simulations both in stationary and distribution shift settings. Through experiments on benchmark OOD datasets in image classification tasks, we demonstrate the practical effectiveness of our proposed approaches.

The paper is structured as follows: Section 2 presents the framework in detail, while Section 3 provides theoretical guarantees on False Positive Rate (FPR) control. In Section 4, we conduct a comprehensive empirical evaluation of the proposed system. The proof details and extensive experimental results are in Appendix.

2. Human-in-the-Loop OOD Detection

In this section, we discuss our proposed system in detail. Recall that we are motivated by two facts. First, the type of OOD samples the system will encounter after deployment are often not available during development, hence we need to build OOD detection systems that can adapt to various kinds of OOD data that it encounters on-the-fly after deployment. Second, in many critical applications, the cost of false positives i.e. misclassifying an OOD point as ID can have more severe consequences than misclassifying an ID point as OOD. E.g., in medical diagnosis, when in doubt it is preferable to classify a brain scan as OOD and seek the input of a radiologist, rather than labeling it with a disease or as normal using the machine-learned classifier. We propose a human-in-the-loop OOD detection system (Figure 1) that can work with any ML inference model and scoring function for OOD detection. We begin by describing the problem setting and then discuss each component of our system in detail. See algorithm 1 for step-by-step details.

Data Stream: Let $\mathcal{X} \subseteq \mathbb{R}^d$ denote the feature space of the data points. The OOD detection system is expected to classify an incoming data point as either “1” i.e. ID (In Distribution) or “−1” stands i.e. OOD (Out Of Distribution), i.e. the label space is $\mathcal{Y} = \{-1, 1\}$. Let the distribution of

ID and OOD data be denoted by \mathcal{D}_{id} and \mathcal{D}_{ood} respectively. Let $x_t \in \mathcal{X}$ be the sample received at the time t . Let $y_t \in \{-1, 1\}$ be the true label for x_t with respect to ID or OOD classification. Assume x_t are independent and drawn according to the following mixture model, $x_t \sim (1-\gamma)\mathcal{D}_{id} + \gamma\mathcal{D}_{ood}$, where $\gamma \in [0, 1]$ is the fraction of OOD in the mixed stream. Note, \mathcal{D}_{id} , \mathcal{D}_{ood} and γ are *unknown*.

Scoring function for OOD detection: After receiving data point x_t , the system uses a given scoring function, $g : \mathcal{X} \mapsto \mathcal{S} \subseteq \mathbb{R}$, to compute a real-valued score quantifying the uncertainty of the point being ID or OOD. There are several scoring functions developed recently e.g. energy-based scores (Liu et al., 2020), Mahalanobis distance-based scores (Lee et al., 2018) etc. Our system can use any of these post-hoc OOD uncertainty quantification functions. We emphasize that our aim is not to design a new OOD uncertainty quantification method (scoring function), instead, we propose a system in which any such g can be plugged in and it can control the FPR.

Algorithm 1 Human in the Loop OOD Detection

Input: FPR threshold α , window size N_w ,

```

1: sampling probability  $p \in (0, 1)$  Scoring function  $g : \mathcal{X} \mapsto \mathbb{R}$ ,
2:  $S_0 = \emptyset, \hat{\lambda}_0 = \infty$ 
3: for  $t = 1, 2, \dots$  do
4:   Receive data point  $x_t$ 
5:    $s_t = g(x_t)$ 
6:   if  $s_t \leq \hat{\lambda}_{t-1}$  then
7:      $l_t = 1$ 
8:   else
9:      $l_t \sim \text{Bernoulli}(p)$ 
10:  end if
11:  if  $l_t = 1$  then
12:     $y_t = \text{GetExpertLabel}(x_t)$ 
13:     $S_t = S_{t-1} \cup \{(s_t, y_t)\}$ 
14:  end if
15:   $\hat{\lambda}_t = \text{SolveOptForLambda}(S_t, N_w, \alpha)$ 
16:   $\hat{y}_t = \text{sign}(s_t - \hat{\lambda}_t)$ 
17:  if  $l_t = 1$  then
18:    Output  $y_t$ 
19:  else
20:    Output  $\hat{y}_t$ 
21:  end if
22: end for
```

Algorithm 2 SolveOptForLambda

Input: FPR threshold α , S_t

```

1:  $\hat{\lambda}_t := \arg \min_{\lambda \in \Lambda} \lambda \text{ s.t. } \widehat{\text{FPR}}(\lambda, t) + \psi(t, \delta) \leq \alpha$ 
2: Output  $\hat{\lambda}_t$ 
```

Denote the score computed for point x_t as $s_t = g(x_t)$.

To be consistent across various scoring functions, let a higher score indicate ID and a lower score indicate OOD points. After computing the uncertainty score s_t the system needs to decide whether x_t is OOD or ID, which is done using a threshold-based classifier parameterized with $\lambda \in \Lambda \subseteq \mathbb{R}$: $h_\lambda(g(x)) = \text{sign}(g(x) - \lambda)$. Here we assume $\Lambda = [\Lambda_{\min}, \Lambda_{\max}]$ is a bounded subset of \mathbb{R} . The threshold-based prediction is common in the OOD detection literature (Liu et al., 2020; Lee et al., 2018). Since OOD data is usually not available during development, a common practice is to find a threshold $\hat{\lambda}$ that correctly classifies at least 95% of the ID data used for training/validation of the ML system as ID, i.e., $\hat{\lambda}$ is chosen for achieving 95% TPR. While simple, a drawback of this approach is that it can result in an exceedingly high FPR, as demonstrated by a thorough examination conducted in the Open-OOD benchmark, where the FPR can range from 60% to 96% on various combinations of ID and OOD datasets. In contrast, real-world systems must handle a diverse range of OOD samples in real time, and for critical applications such as medical diagnosis, it is imperative to ensure control over the FPR. The population level FPR and TPR for any $\lambda \in \Lambda$ are defined as follows,

$$\text{FPR}(\lambda) = \mathbb{E}_{s \sim \mathcal{D}_{ood}}[\mathbf{1}\{s > \lambda\}] \quad \text{and.} \quad (1)$$

$$\text{TPR}(\lambda) = \mathbb{E}_{s \sim \mathcal{D}_{id}}[\mathbf{1}\{s > \lambda\}] \quad (2)$$

Note that the cumulative distribution function (CDF) of \mathcal{D}_{ood} , $\text{CDF}_{\mathcal{D}_{ood}}(\lambda) = \mathbb{E}_{s \sim \mathcal{D}_{ood}}[\mathbf{1}\{s \leq \lambda\}]$. Therefore, $\text{FPR}(\lambda) = 1 - \text{CDF}_{\mathcal{D}_{ood}}(\lambda)$. Similarly, $\text{TPR}(\lambda) = 1 - \text{CDF}_{\mathcal{D}_{id}}(\lambda)$. Since the CDF of any distribution is a monotonic function, both FPR and TPR are monotonic in λ .

Expert feedback and importance sampling: In our proposed system, we choose λ adaptively using human feedback so that the FPR is maintained below the user-specified rate of α . One can of course achieve this trivially by setting $\lambda_t = \Lambda_{\max}$, i.e., always getting human feedback. Therefore, in addition to controlling the FPR, we want to minimize the human feedback solicited by the system. This is equivalent to maximizing the true positive rate. That is, $\lambda_t := \arg \max_{\lambda} \text{TPR}(\lambda)$ subject to $\text{FPR}(\lambda) \leq \alpha$. Since the TPR is monotonic in λ , this can be re-written as,

$$\lambda_t^* := \arg \underset{\lambda \in \Lambda}{\text{minimize}} \lambda, \quad \text{subject to} \quad \text{FPR}(\lambda) \leq \alpha. \quad (\text{P1})$$

Optimal threshold, denoted by λ^* , is the smallest λ such that $\text{FPR}(\lambda^*) = \alpha = 1 - \text{CDF}_{\mathcal{D}_{out}}(\lambda^*)$

(see Figure 2). When the distribution of the OOD points, \mathcal{D}_{ood} , is not changing, $\lambda_t^* =: \lambda^*$. Note that, γ , the mixture ratio, or the distribution of the ID points \mathcal{D}_{id} changing does not affect the value of the optimal threshold. As we do not have access to the true FPR and TPR values, we cannot solve the optimization problem (P1). Instead, we have to estimate λ_t^* using the observations made till time t . Thus,

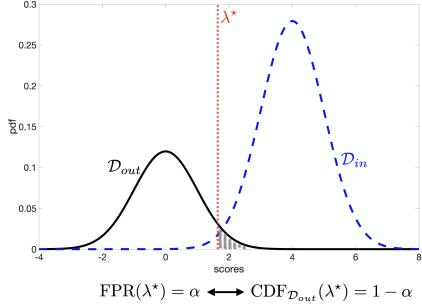


Figure 2. Optimal λ^* for the optimization problem (P1) with $\alpha = 0.05$ and $x_t \stackrel{i.i.d.}{\sim} 0.7 \mathcal{D}_{in} + 0.3 \mathcal{D}_{out}$, where the scores of \mathcal{D}_{in} and \mathcal{D}_{out} are distributed as $\mathcal{N}(4, 1)$ and $\mathcal{N}(0, 1)$ respectively.

at each time point our goal is to find $\hat{\lambda}_t \in \Lambda$ such that the FPR when using $\hat{\lambda}_t$ as the threshold in h_λ , denoted by $FPR(\hat{\lambda}_t) \leq \alpha$.

Ideally, we want to avoid human feedback for points with a score greater than $\hat{\lambda}_t$, i.e., those points that are determined as ID by the system. However, in order to have an unbiased estimate of the FPR and also to allow for potential change in the distribution of OOD samples and therefore change in true FPR, we allow for human feedback with a small probability p for points predicted as in-distribution by the system to be able to detect a change.

FPR estimation and adapting the threshold: At each time t , we observe $x_t \stackrel{i.i.d.}{\sim} (1-\gamma)\mathcal{D}_{id} + \gamma\mathcal{D}_{ood}$, and $s_t = g(x_t)$ is the corresponding score. If $s_t \leq \hat{\lambda}_{t-1}$, then it is considered an OOD point and hence gets a human label for it and we get to know whether it is in fact OOD or ID. If $s_t > \hat{\lambda}_{t-1}$, then it is considered an ID point and hence gets a human label only with probability p . So, we get to know whether it is truly ID or not with probability p . Now we have to update the threshold, $\hat{\lambda}_t$, such that the $FPR(\hat{\lambda}_t) \leq \alpha$ for all t , while trying to maximize TPR($\hat{\lambda}_t$). Our approach is based on constructing an unbiased estimator of $FPR(\lambda)$ using the OOD samples received till time t and in conjunction with confidence intervals for $FPR(\lambda)$ for all thresholds $\lambda \in \Lambda$ that is valid for all times simultaneously. Together, at each time t , these give us a reliable upper bound on the true $FPR(\lambda)$ for all λ enabling us to find the smallest λ such that the upper bound on $FPR(\lambda)$ is at most α . Let $S_t^{(o)} = \{s_1^{(o)}, \dots, s_{N_t^{(o)}}^{(o)}\}$ denote the scores of the points that have been truly identified as OOD points from human feedback. We estimate the FPR as follows,

$$\widehat{FPR}(\lambda, t) = \frac{1}{N_t^{(o)}} \sum_{u=1}^{N_t^{(o)}} Z_u(\lambda), \quad \text{where,} \quad (3)$$

$$Z_u(\lambda) := \begin{cases} 1(s_u^{(o)} > \lambda), & \text{if } s_u^{(o)} \leq \hat{\lambda}_{u-1} \\ \frac{1}{p} 1(s_u^{(o)} > \lambda), & \text{w.p. } p \text{ if } s_u^{(o)} > \hat{\lambda}_{u-1} \\ 0, & \text{w.p. } 1-p \text{ if } s_u^{(o)} > \hat{\lambda}_{u-1} \end{cases}. \quad (4)$$

Lemma 1. $\widehat{FPR}(\lambda, t)$ as defined in eq. (6) is an unbiased estimate of the true FPR(λ), i.e., $\mathbb{E}[\widehat{FPR}(\lambda, t)] = FPR(\lambda)$.

Finding threshold using a UCB on FPR: We use this estimated FPR with an upper confidence bound (UCB) to replace the unknown true FPR in the optimization problem (P1) to obtain the following optimization problem (P2),

$$\hat{\lambda}_t := \arg \min_{\lambda \in \Lambda} \lambda \text{ s.t. } \widehat{FPR}(\lambda, t) + \psi(t, \delta) \leq \alpha, \quad (P2)$$

where the term $\psi(t, \delta)$ is a time-varying upper confidence which is simultaneously valid for all λ for all time with probability at least $1 - \delta$ for any given $\delta \in (0, 1)$. The minimization problem can be solved in many ways. We use a binary search procedure where we search over a grid on $[\Lambda_{\min}, \Lambda_{\max}]$ with grid-size ν . The procedure 2 searches for a smallest λ such that $\widehat{FPR}(\lambda, t) + \psi(t, \delta) \leq \alpha$. It uses eq. (6) to compute the empirical FPR at various thresholds and the confidence interval $\psi(t, \delta)$ given in eq. (5). Details of the binary search procedure are in the Appendix.

Upper confidence bound (UCB): At each time point, the algorithm estimates the FPR using a finite number of samples at all thresholds. We need confidence intervals that are valid for all thresholds at all time points to ensure the algorithm has reliable upper bounds on the FPR. In particular, we use the Law of iterated logarithm(LIL) (Khinchine, 1924) based confidence bounds that are known to be tight. In our setting, due to the importance sampling, the samples become conditionally dependent. This dependence prevents direct application of known results like (Howard & Ramdas, 2022). In section 3 we build upon the LIL bounds for martingales (Balsubramani, 2015) and derive a confidence interval bound that is valid in our setting (see equation (5)),

$$\psi(t, \delta) := \sqrt{\frac{3c_t}{N_t^{(o)}} \left[2 \log \log \left(\frac{3c_t N_t^{(o)}}{2} \right) + \log \left(\frac{2|\Lambda|}{\delta} \right) \right]}, \quad (5)$$

where $c_t = 1 - \beta_t + \frac{\beta_t}{p^2}$, $\beta_t = \frac{N_t^{(o,p)}}{N_t^{(o)}}$ and $N_t^{(o,p)}$ is the number of points sampled using importance sampling until time t and $\nu \in (0, 1)$ is a given discretization parameter.

Handling distribution shift: One of the motivations for the system is to be able to adapt to the variations of the OOD data. As long as \mathcal{D}_{ood} does not change, changes in the \mathcal{D}_{id} or the mixing ratio γ do not affect the true FPR. However,

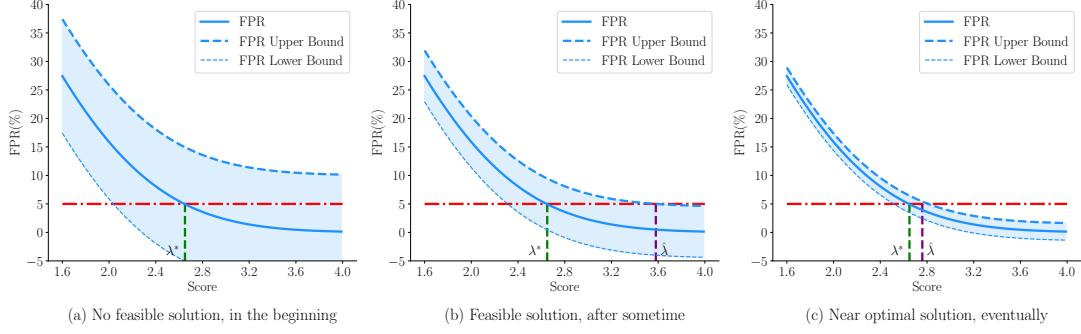


Figure 3. Illustration of the confidence intervals on FPR and their effect on threshold estimation. We expect as the system receives more OOD samples the confidence intervals will shrink and lead to a better threshold.

the true FPR does get affected when \mathcal{D}_{ood} changes. When there is a change in \mathcal{D}_{ood} , estimating the FPR using all the acquired samples so far can lead to inaccurate estimates as the current estimate is highly influenced by scores that are far behind in time from the previous \mathcal{D}_{ood} . This can lead to inaccurate estimation of λ and erroneous predictions for the current time. To overcome this challenge, we propose a sliding window-based approach where the user can set a window size $N_w > 0$ and the system will only estimate the FPR and the confidence intervals using the most recent N_w OOD samples. This will allow the system to adapt the threshold that is well aligned with the new distribution(s) of OOD samples. Next, we provide theoretical guarantees for the proposed algorithm in the setting when \mathcal{D}_{ood} does not change over time.

3. Theoretical Guarantees

In this section, we provide theoretical analysis and guarantees for Algorithm 1 when the distributions are fixed. Also, assume that the scores $g(x)$ have sub-Gaussian tails. Here we assume that \mathcal{D}_{ood} is not changing. We provide anytime valid confidence intervals on the FPR at all thresholds which are used in the optimization problem (P2), using which we can guarantee that the FPR is always controlled. We also provide a bound on the time taken to reach feasibility, i.e., for the constraint in Equation (P2) to be feasible. Furthermore, we also provide the bound on time taken to reach η -Optimality which is defined as follows,

Definition 1. (η -Optimality) The system is said to be operating in the η -Optimal regime after some time point T_η , if $FPR(\lambda^*) - FPR(\hat{\lambda}_t) \leq \eta$.

Theorem 1. Let $\alpha \in (0, 1)$ and $\delta \in (0, 1)$. Let $x_t \stackrel{i.i.d}{\sim} (1 - \gamma)\mathcal{D}_{id} + \gamma\mathcal{D}_{ood}$. If Algorithm 1 uses the optimization problem (P2) to find the thresholds with the upper confidence term $\psi(N_t^{(o)}, \delta)$ given by equation (5), then, with probability at least $1 - \delta$,

1. (Controlled FPR) For all $t \geq 1$, $FPR(\hat{\lambda}_t) \leq \alpha$.
2. (Time to reach feasibility) Let T_f be such that $N_{T_f}^{(o)} \geq$

$\frac{C_1}{\alpha^2} \log\left(\frac{C_2}{\delta} \log\left(\frac{C_3}{\alpha}\right)\right)$ then for any $t \geq T_f$ the algorithm will find a feasible threshold, $\hat{\lambda}_t$ such that $\widehat{FPR}(\hat{\lambda}_t) + \psi(N_t^{(o)}) \leq \alpha$.

3. (Time to reach η -Optimality) Let T_{opt} be such that $N_{T_{opt}}^{(o)} \geq \frac{4C_1}{\eta^2} \log\left(\frac{C_2}{\delta} \log\left(\frac{2C_3}{\eta}\right)\right)$ and $\widehat{FPR}(\hat{\lambda}_{T_{opt}}) \geq [FPR(\lambda^*) - \eta/2, \alpha]$, then then for any $t \geq T_{opt}$ if the $\hat{\lambda}_t$ satisfy the η -Optimality condition in definition 1.

The above theorem establishes key properties of Algorithm 1 and provides insights into its behavior and performance guarantees. We now discuss each property in detail, along with their implications.

Controlled false positive rate: The first property of Theorem 2 ensures that the Algorithm 1 effectively controls the False Positive Rate (FPR) throughout its operation. Specifically, it guarantees that for all time steps $t \geq 1$, the FPR of the estimated threshold obtained by the algorithm will be less than or equal to a predetermined threshold α . This property is crucial in applications where accurately controlling the rate of false positives is essential. By limiting the FPR to a predefined threshold, Algorithm 1 provides a reliable mechanism for distinguishing between in-distribution and out-of-distribution samples, reducing the likelihood of erroneous classifications.

Time to reach feasibility: The second property of Theorem 2 concerns the time it takes for Algorithm 1 to find a feasible threshold. It provides conditions under which the algorithm is guaranteed to discover a suitable threshold, $\hat{\lambda}_t$, that has FPR at most α . It is contingent upon the feasible time T_f , the time step at which a sufficient number of observations $N_{T_f}^{(o)}$ is obtained so that the confidence interval $\psi(t, \delta) \leq \alpha$.

Time to reach η -Optimality: The third property provides a bound on the time T_{opt} after which the Algorithm 1 achieves the η -Optimal regime. This regime implies that the algorithm operates in a state where the difference between the FPR of the true optimal threshold, $FPR(\lambda^*)$, and the FPR of the estimated threshold $FPR(\hat{\lambda}_t)$, is within an accept-

able range of η . The theorem says that, if the estimated FPR at time step T_{opt} , denoted as $\bar{\text{FPR}}(\hat{\lambda}_{T_{opt}})$, is within the range $[\text{FPR}(\lambda^*) - \eta/2, \alpha]$ and the confidence interval $\psi(T_{opt}, \delta) \leq \eta/2$. Then for all time points after T_{opt} the algorithm will find a $\hat{\lambda}_t$ that satisfies the η -Optimality condition. T_{opt} is defined in terms of the time point when the number of acquired OOD samples $N_{T_{opt}}^{(o)}$ becomes at least $\frac{4C_1}{\eta^2} \log \left(\frac{C_2}{\delta} \log \left(\frac{2C_3}{\eta} \right) \right)$. We require these many samples to guarantee the confidence intervals $\psi(t, \delta)$ are sufficiently small (of the order of η in this case) so that when the empirical estimate of FPR is very close to α we know that the algorithm will return a threshold satisfying η -Optimality.

If γ is not changing, it is very easy to bound T_f and T_{opt} . When t is large enough, with high probability γ fraction of what is observed is going to be OOD. So, for T_f , N_f will concentrate around γT_f . And similarly for T_{opt} , N_{opt} while also accounting for importance sampling. The details of the proof of the statements in the main theorem are provided in the appendix. Here we provide the key results and outline the key ideas of the proof.

Proof outline and discussion: The main technical challenge is to obtain accurate confidence intervals $\psi(t, \delta)$ that are simultaneously valid with high probability for the FPR estimates at all time points and all thresholds. Fortunately, there is a rich line of work that provide tight confidence intervals valid for all times based on the Law of Iterated Logarithm (LIL) (Khinchine, 1924; Kolmogorov, 1929; Smirnov, 1944). Non-asymptotic versions of LIL have been proved in various settings e.g. multi-armed bandits (Jamieson et al., 2013), quantile estimation (Howard & Ramdas, 2022). Roughly speaking, these bounds provide confidence intervals that are $\mathcal{O}(\sqrt{\log \log(t)/t})$ and are known to be tight. However, most of these works assume i.i.d samples. In our setting, once we reach a feasible regime, our treatment of observing the human feedback is dependent on whether the score is above or below $\hat{\lambda}_{t-1}$ which itself is estimated using all the past data which creates dependence. We handle this by first showing that there is a martingale structure that we then exploit by using LIL results for martingales (Balsubramani, 2015). A limitation of (Balsubramani, 2015) is that it can only provide us confidence intervals valid for FPR estimate for a given threshold λ . However, we need intervals that are simultaneously valid for all λ as well. Building upon the work in (Balsubramani, 2015) we derive confidence intervals that are simultaneously valid for all t and finitely many thresholds. Equation (5) shows the $\psi(t, \delta)$ we obtain. Please see the Appendix for detailed proofs.

4. Empirical Evaluation

Methods: We evaluate our method on synthetic and real-world image classification datasets. We compare our (a) LIL confidence interval based method against (b) No-UCB:

which does not use any confidence intervals (c) Hoeffding: which uses the confidence intervals from Hoeffding's inequality. We consider two variations of each method, one without using a window and the other using a window size. We expect that No-UCB will violate the FPR constraint since it is not accounting for the uncertainty in the estimates. While the methods that accurately account for the uncertainty using confidence intervals like LIL, and Hoeffding are expected to adhere to the FPR constraints. We note, that the confidence intervals from Hoeffding inequality are not theoretically valid for these settings but are a reasonable choice for a practitioner, and in our evaluation, we do not observe significant differences between Hoeffding and LIL-based bounds end results. We use $\alpha = 0.05$, $\delta = 0.2$, and importance sampling probability $p = 0.2$ through all the empirical evaluations.

The theoretical LIL bound in 5 has constants that can be pessimistic in practice. We get around this by using a LIL-Heuristic bound which has the same form as in equation (5) but with different constants in particular we consider the form in equation LIL-Heuristic. We find the constants C_1, C_2, C_3 using a simulation on estimating the bias of a coin with different constants and picking the ones so that the observed failure probability is below 5%. More details on this are in the Appendix.

$$C_1 \sqrt{\frac{c_t}{N_t^{(o)}} \left(\log \log(C_2 c_t N_t^{(o)}) + \log \left(\frac{C_3}{\delta} \right) \right)}. \quad (\text{LIL-Heuristic})$$

Synthetic data setup: We simulate the OOD and ID scores using a mixture of two Gaussians $\mathcal{N}_{id}(\mu = 5.5, \sigma = 4)$ and $\mathcal{N}_{ood}(\mu = -6, \sigma = 4)$. To simulate distribution change we change the OOD distribution to $\mathcal{N}_{ood}(\mu = -5, \sigma = 4)$ at $t = 55k$. We shuffle 100K OOD and 10K ID samples and run the methods. We show the results in figure 4.

Real data setup: We evaluate our proposed system empirically on two sets of benchmarks from OpenOOD Benchmark (Yang et al., 2022). Here we show the results on CIFAR-10 ID dataset and show the results on CIFAR-100 ID dataset in the appendix. CIFAR-10 (Krizhevsky et al., 2009). CIFAR-10 is a 10-class dataset for general object classification. We use the official testing datasets as the ID dataset. The OOD datasets for CIFAR-10 consists of CIFAR100, SVHN (Netzer et al., 2011), TinyImageNet (Krizhevsky et al., 2017) (1,207 images are removed from TinyImageNet since they belong to CIFAR-10 (Yang et al., 2021a)), MNIST (Deng, 2012), Texture (Kylberg, 2011), Places365 (Zhou et al., 2017) (1,305 images are removed due to semantic overlaps). We use a pre-train ResNet-18 with 94.3% accuracy for all the experiments on CIFAR-10.

Data Stream: For the non-distribution shifted setting, we combine all six OOD datasets for a joint OOD distribu-

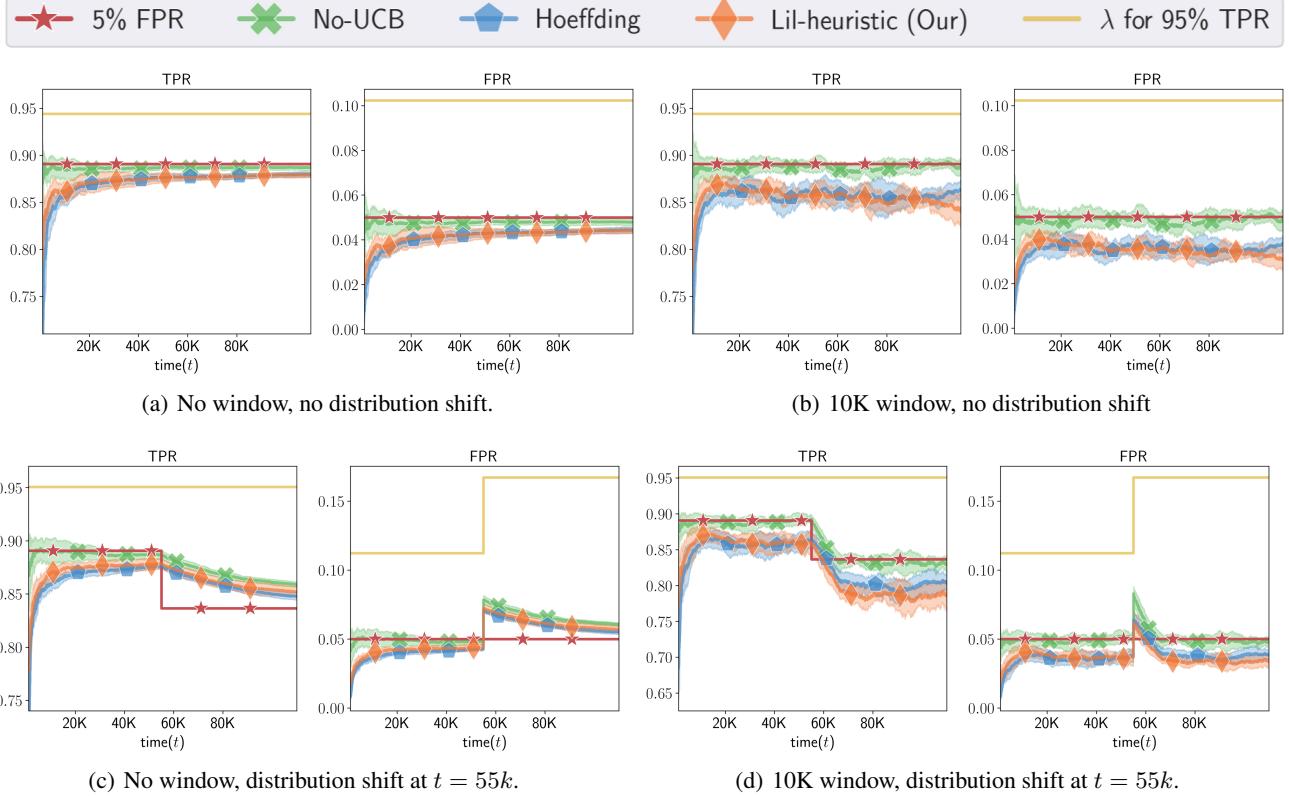


Figure 4. Experiments on Synthetic Data. Each method is repeated 20 times. The mean and standard deviation is shown. The distribution shift starts at $t = 55k$ for figure 4(c), 4(d).

tion. For shifted distribution setting, we sample a portion of OOD samples from three OOD datasets and then sample the shifted distribution samples from the rest of the OOD datasets. We randomly sample 90k OOD and 9k ID samples.

Computing OOD Scores: Accurately detecting OOD points in the online setting needs a good scoring function that separates the ID and OOD points at some threshold score λ . We leverage existing works on the construction of the scoring function. We use ODIN (Liang et al., 2017), Mahalanobis Distance (Lee et al., 2018), Energy Score (Liu et al., 2020), SSD (Sehwag et al., 2021), VIM (Wang et al., 2022), and KNN (Sun et al., 2022) scores for the evaluation. We use an open-source codebase, OpenOOD (Yang et al., 2022), to implement all the methods. Due to space limitation, we present results for KNN (Sun et al., 2022) score here. For more details on these scores and results on the rest of the scores please see the Appendix.

Discussion: As expected, in the fixed distributions setting in both synthetic and real data settings (figures 4(a), and 9(a), respectively), we see that not using a UCB leads to violation of FPR constraints and the methods with LIL-Heuristic, Hoeffding based intervals are able to maintain the FPR below the user given threshold 5%. Moreover, all the methods improve as they acquire more samples with

time and eventually reach very close to the optimal solution. When we run these methods with a window size of 10k in the fixed distribution setting (figures 7(b), 9(b)) we observe similar behavior except with a bit more variance since with a fixed window the confidence intervals are not shrinking with time. Though the windowed setting is more useful when the distributions change and not so much of use in the fixed distribution case, nevertheless we show this experiment to validate our understanding of the fixed distribution setting.

Moving forward, we investigate the case where the distributions change at a specific time point. In such scenarios, we find that the windowed approach adapts more rapidly compared to the method without a window (see figures 4(c), 4(d), 9(c), 9(d)). The use of a fixed window allows the algorithm to quickly adjust to the change, whereas without a window, the adaptation process is significantly delayed.

In summary, our findings demonstrate that the choice of using a windowed approach or not depends on the nature of the data and the presence of distribution changes. The windowed approach proves advantageous in scenarios where rapid adaptation is crucial, while the non-windowed approach can still be effective, albeit potentially with longer adaptation times. The consistency between our observations in the synthetic and the real-data evaluations provides strong

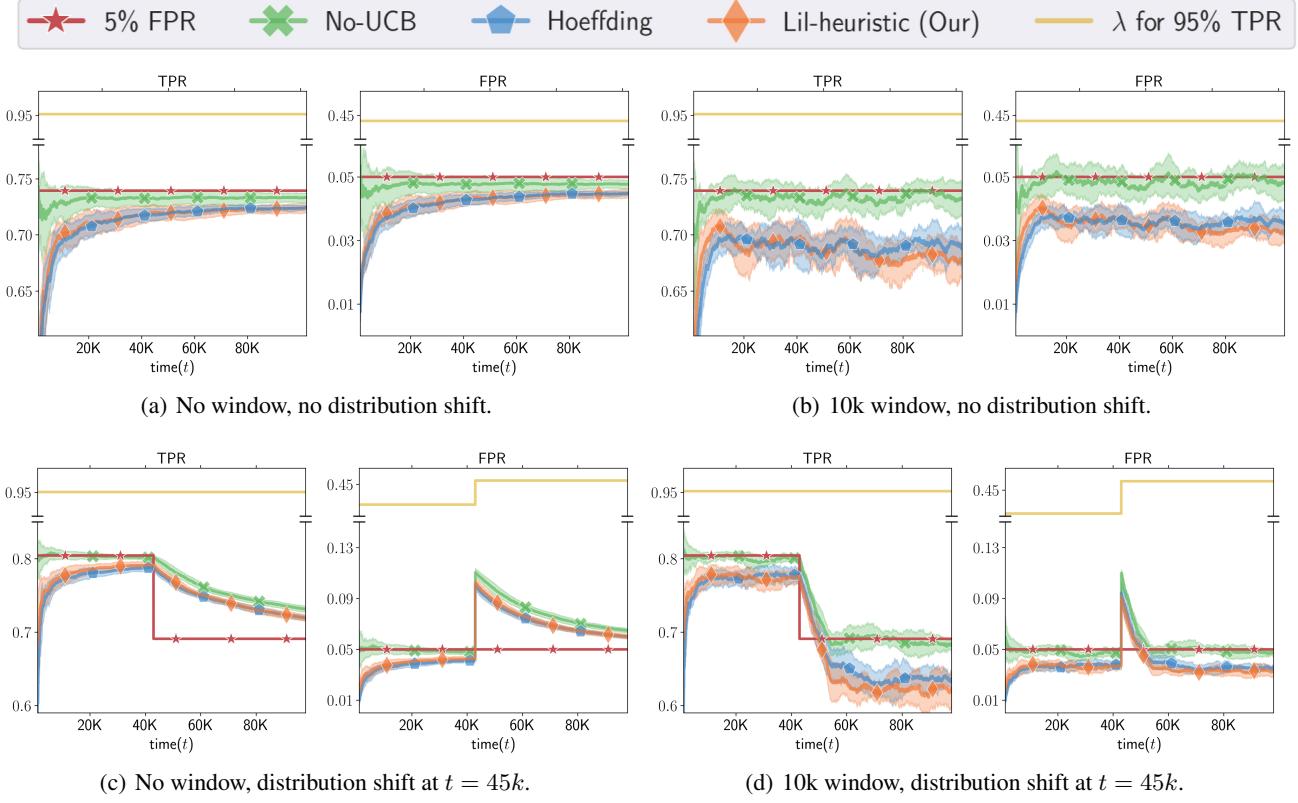


Figure 5. Experiments on the KNN method with cifar10 as the ID dataset. In the distribution shift setting, the OOD datasets are initially a mixture of MNIST, Places365, and Cifar100. After the distribution shift, OOD datasets consist of TinyImageNet, SVHN, and SVHN.

evidence of the reliability and effectiveness of our method.

5. Related Works

Out-Of-Distribution detection: The problem of OOD detection has been addressed in many recent works where the main contributions have been methods to quantify a score (uncertainty) which gives a better separation of OOD and ID data points. Liang et al. (Liang et al., 2017) proposed ODIN, which uses temperature scaling to separate the softmax score distributions between ID and OOD images. Liu et al. (Liu et al., 2020) proposed using energy score to perform OOD detection on pre-trained neural classifiers. Lee et al. (Lee et al., 2018), Sehwag et al. (Sehwag et al., 2021), and Ming et al. (Ming et al., 2022) proposed mahalanobis distance-based scores to detect OOD samples. While these methods perform well, the evaluation setup is rather static and does not reflect the real-world deployment scenario, wherein the system has to adapt to new and evolving OOD data. In our work we are proposing a simple and extensible system for online OOD detection. Moreover, the system can also adapt with labels from humans on selected points.

Online anomaly detection: There is rich literature on anomaly (or outlier) detection in offline settings (Chandola et al., 2009; Campos et al., 2016; Chalapathy & Chawla,

2019). However, our setting is akin to the online anomaly (outlier) detection – wherein the system receives samples one at a time and it has to figure out the outliers or anomalous behavior within a given window of time. Some of the notable works along this line are (Subramanian et al., 2006; Angiulli & Fassetti, 2007; Zhang et al., 2013). The methods proposed are unsupervised and perform density or distance-based detection.

Outlier detection with human in the Loop: The notion of outlier may not always be based on statistical rarity and might need input from humans to learn the notion of outlier in the application of interest. Some of the recent works (Chai et al., 2020; Islam et al., 2018) have given methods for outlier detection in offline setting leveraging human inputs. The focus has been on minimizing human effort by figuring out some candidate outliers and designing good questions and context for getting human inputs.

While there are a number of works on outlier or OOD detection, the main focus has been on designing methods (scoring functions) to distinguish inliers vs outliers mostly in the offline setting. Our work is rather complementary – we consider a deployed OOD system that can work with any scoring function and propose ways for online adaptation of this system based on human inputs.

6. Conclusion

We presented a mathematically grounded framework for human-in-the-loop Out-of-Distribution (OOD) detection. By incorporating expert feedback and utilizing confidence intervals based on the Law of Iterated Logarithm (LIL), our approach maintains control over false positive rates (FPR) while maximizing true positive rates (TPR). The empirical evaluations on synthetic data and image classification tasks demonstrate the effectiveness of our method in maintaining FPR at or below 5% while achieving high TPR. Our work gives a promising solution for addressing the challenge of robustness to OOD samples in real-world applications.

References

- Amodei, D., Olah, C., Steinhardt, J., Christiano, P. F., Schulman, J., and Mané, D. Concrete problems in AI safety. *CoRR*, abs/1606.06565, 2016.
- Angiulli, F. and Fassetti, F. Detecting distance-based outliers in streams of data. In *Proceedings of the Sixteenth ACM Conference on Conference on Information and Knowledge Management*, CIKM '07, pp. 811–820, 2007. ISBN 9781595938039.
- Balsubramani, A. Sharp finite-time iterated-logarithm martingale concentration, 2015.
- Campos, G. O., Zimek, A., Sander, J., Campello, R. J., Micenková, B., Schubert, E., Assent, I., and Houle, M. E. On the evaluation of unsupervised outlier detection: Measures, datasets, and an empirical study. *Data Min. Knowl. Discov.*, 30(4):891–927, 2016. ISSN 1384-5810. doi: 10.1007/s10618-015-0444-8. URL <https://doi.org/10.1007/s10618-015-0444-8>.
- Chai, C., Cao, L., Li, G., Li, J., Luo, Y., and Madden, S. Human-in-the-loop outlier detection. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, SIGMOD '20, pp. 19–33, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450367356. doi: 10.1145/3318464.3389772. URL <https://doi.org/10.1145/3318464.3389772>.
- Chalapathy, R. and Chawla, S. Deep learning for anomaly detection: A survey. *arXiv preprint arXiv:1901.03407*, 2019.
- Chandola, V., Banerjee, A., and Kumar, V. Anomaly detection: A survey. *ACM Comput. Surv.*, 41(3), jul 2009. ISSN 0360-0300. doi: 10.1145/1541880.1541882. URL <https://doi.org/10.1145/1541880.1541882>.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Deng, L. The mnist database of handwritten digit images for machine learning research [best of the web]. *IEEE signal processing magazine*, 29(6):141–142, 2012.
- Howard, S. R. and Ramdas, A. Sequential estimation of quantiles with applications to A/B testing and best-arm identification. *Bernoulli*, 28(3):1704 – 1728, 2022. doi: 10.3150/21-BEJ1388. URL <https://doi.org/10.3150/21-BEJ1388>.
- Islam, M. R., Das, S., Doppa, J. R., and Natarajan, S. Glad: Glocalized anomaly detection via human-in-the-loop learning. *arXiv preprint arXiv:1810.01403*, 2018.
- Jamieson, K., Malloy, M., Nowak, R., and Bubeck, S. lil' ucb : An optimal exploration algorithm for multi-armed bandits, 2013.
- Khinchine, A. Über einen satz der wahrscheinlichkeitsrechnung. *Fundamenta Mathematicae*, 6:9–20, 1924.
- Kolmogorov, A. Über das gesetz des iterierten logarithmus. *Mathematische Annalen*, 101:126–135, 1929.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- Kylberg, G. *Kylberg texture dataset v. 1.0*. Centre for Image Analysis, Swedish University of Agricultural Sciences and ..., 2011.
- Lee, K., Lee, K., Lee, H., and Shin, J. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems*, 31, 2018.
- Liang, S., Li, Y., and Srikant, R. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017.
- Liu, W., Wang, X., Owens, J., and Li, Y. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems*, 33:21464–21475, 2020.
- Ming, Y., Sun, Y., Dia, O., and Li, Y. Cider: Exploiting hyperspherical embeddings for out-of-distribution detection. *arXiv preprint arXiv:2203.04450*, 2022.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. 2011.

Nguyen, A. M., Yosinski, J., and Clune, J. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *CVPR*, pp. 427–436. IEEE Computer Society, 2015.

Sehwag, V., Chiang, M., and Mittal, P. Ssd: A unified framework for self-supervised outlier detection. *arXiv preprint arXiv:2103.12051*, 2021.

Smirnov, N. Approximate laws of distribution of random variables from empirical data. *Uspekhi Matematicheskikh Nauk*, 10:179–206, 1944.

Subramaniam, S., Palpanas, T., Papadopoulos, D., Kalogeraki, V., and Gunopoulos, D. Online outlier detection in sensor data using non-parametric models. In *Proceedings of the 32nd International Conference on Very Large Data Bases*, VLDB ’06, pp. 187–198. VLDB Endowment, 2006.

Sun, Y., Ming, Y., Zhu, X., and Li, Y. Out-of-distribution detection with deep nearest neighbors. In *International Conference on Machine Learning*, pp. 20827–20840. PMLR, 2022.

Wang, H., Li, Z., Feng, L., and Zhang, W. Vim: Out-of-distribution with virtual-logit matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4921–4930, 2022.

Yang, J., Wang, H., Feng, L., Yan, X., Zheng, H., Zhang, W., and Liu, Z. Semantically coherent out-of-distribution detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8301–8309, 2021a.

Yang, J., Zhou, K., Li, Y., and Liu, Z. Generalized out-of-distribution detection: A survey. *arXiv preprint arXiv:2110.11334*, 2021b.

Yang, J., Wang, P., Zou, D., Zhou, Z., Ding, K., Peng, W., Wang, H., Chen, G., Li, B., Sun, Y., Du, X., Zhou, K., Zhang, W., Hendrycks, D., Li, Y., and Liu, Z. Openood: Benchmarking generalized out-of-distribution detection, 2022.

Zhang, Y., Meratnia, N., and Havinga, P. J. Distributed online outlier detection in wireless sensor networks using ellipsoidal support vector machine. *Ad Hoc Networks*, 11 (3):1062–1074, 2013.

Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., and Torralba, A. Places: A 10 million image database for scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 40(6):1452–1464, 2017.

Appendix

The appendix is organized as follows. We summarize the notation in Table 1. Then we give the proof of the main theorem (Theorem 2) and the proofs of supporting lemmas. Further we provide additional experiments and insights from them.

Glossary

The notation is summarized in Table 1 below.

Symbol	Definition
\mathcal{X}	feature space.
\mathcal{Y}	label space, $\{+1, -1\}$, +1 for ID and -1 for OOD .
$\mathcal{D}_{id}, \mathcal{D}_{ood}$	distributions of ID and OOD points.
γ	mixing ratio of OOD and ID distributions.
λ	threshold for OOD classification.
$FPR(\lambda)$	population level false positive rate with threshold λ .
$TPR(\lambda)$	population level true positive rate with threshold λ .
$\widehat{FPR}(\lambda, t)$	empirical FPR at time t , adjusted to account for importance sampling (see eq. (6)).
λ^*	the optimal threshold for OOD classification s.t. $FPR(\lambda) \leq \alpha$ and $TPR(\lambda)$ is maximized.
$\hat{\lambda}_t$	the estimated threshold at round t .
x_t, y_t	sample and the true label at time t .
g	OOD uncertainty quantification (score) function.
$s_u^{(o)}$	score of u^{th} OOD sample.
$i_u^{(o)}$	indicator variable denoting whether $s_u^{(o)}$ was importance sampled or not.
$N_t^{(o)}$	number of OOD points till time t .
$N_t^{(o,p)}$	number of OOD points sampled using importance sampling until time t .
β_t	it is equal to $N_t^{(o,p)} / N_t^{(o)}$.
p	probability for importance sampling.
δ	failure probability.
α	user given upper bound on FPR that the algorithm needs to maintain.
η	the algorithm is in η -optimality if close $FPR(\lambda^*) - FPR(\hat{\lambda}_t) \leq \eta$.
$\Lambda_{\min}, \Lambda_{\max}$	the minimum and maximum scores(thresholds) considered by the algorithm.
ν	discretization parameter for the interval $[\Lambda_{\min}, \Lambda_{\max}]$ set by the user.
$\psi(t, \delta)$	LIL based confidence interval at time t .

Table 1. Glossary of variables and symbols used in this paper.

7. Proofs

At each time t , we observe $x_t \stackrel{i.i.d.}{\sim} (1 - \gamma)\mathcal{D}_{id} + \gamma\mathcal{D}_{ood}$, and $s_t = g(x_t)$ is the corresponding score. If $s_t \leq \hat{\lambda}_{t-1}$, then it is considered an OOD point and hence gets a human label for it and we get to know whether it is in fact OOD or ID. If $s_t > \hat{\lambda}_{t-1}$, then it is considered an ID point and hence gets a human label only with probability p . So, we get to know whether it is truly ID or not with probability p . Now we have to update the threshold, $\hat{\lambda}_t$, such that the $FPR(\hat{\lambda}_t) \leq \alpha$ for all t , while trying to maximize $TPR(\hat{\lambda}_t)$. Our approach is based on constructing an unbiased estimator of $FPR(\lambda)$ using the OOD samples received till time t and in conjunction with confidence intervals for $FPR(\lambda)$ for all thresholds $\lambda \in \Lambda$ that is valid for all times simultaneously. Together, at each time t , these give us a reliable upper bound on the true $FPR(\lambda)$ for all λ enabling us to find the smallest λ such that the upper bound on $FPR(\lambda)$ is at most α . Let $S_t^{(o)} = \{s_1^{(o)}, \dots, s_{N_t^{(o)}}^{(o)}\}$ denote the scores of the points that have been truly identified as OOD points from human feedback. We estimate the FPR as follows,

$$\widehat{FPR}(\lambda, t) = \frac{1}{N_t^{(o)}} \sum_{u=1}^{N_t^{(o)}} Z_u(\lambda), \text{ where } Z_u(\lambda) := \begin{cases} \mathbf{1}(s_u^{(o)} > \lambda), & \text{if } s_u^{(o)} \leq \hat{\lambda}_{u-1} \\ \frac{1}{p} \mathbf{1}(s_u^{(o)} > \lambda), & \text{w.p. } p \text{ if } s_u^{(o)} > \hat{\lambda}_{u-1} \\ 0, & \text{w.p. } 1 - p \text{ if } s_u^{(o)} > \hat{\lambda}_{u-1} \end{cases}. \quad (6)$$

Next, we show that the above estimator $\widehat{\text{FPR}}(\lambda, t)$ is indeed an unbiased of false positive rate $\text{FPR}(\lambda)$.

Lemma 2. $\widehat{\text{FPR}}(\lambda, t)$ as defined in equation (6) is an unbiased estimate of the true $\text{FPR}(\lambda)$, i.e., $\mathbb{E}[\widehat{\text{FPR}}(\lambda, t)] = \text{FPR}(\lambda)$.

Proof. Let $i_t^{(o)}$ be the indicator variable denoting whether $s_t^{(o)}$ was sampled using importance sampling (i.e. $i_t^{(o)} = 1$) or not (i.e. $i_t^{(o)} = 0$). Denote the pair as $r_t^{(o)} = (s_t^{(o)}, i_t^{(o)})$ for brevity.

$$\begin{aligned}
 \mathbb{E}_{r_t^{(o)}, r_{t-1}^{(o)}, \dots, r_1^{(o)}}[\widehat{\text{FPR}}(\lambda, t)] &= \frac{1}{N_t^{(o)}} \sum_{u=1}^{N_t^{(o)}} \mathbb{E}_{r_u^{(o)} | r_{u-1}^{(o)}, \dots, r_1^{(o)}}[Z_u(\lambda)] \\
 &= \frac{1}{N_t^{(o)}} \sum_{u=1}^{N_t^{(o)}} \mathbb{E}_{r_u^{(o)} | \hat{\lambda}_{u-1}}[Z_u(\lambda)] \\
 &= \frac{1}{N_t^{(o)}} \sum_{u=1}^{N_t^{(o)}} \mathbb{E}_{(s_u^{(o)}, i_u^{(o)}) | \hat{\lambda}_{u-1}}[Z_u(\lambda)] \\
 &= \frac{1}{N_t^{(o)}} \sum_{u=1}^{N_t^{(o)}} \mathbb{E}_{s_u^{(o)} | \hat{\lambda}_{u-1}}[\mathbb{E}_{i_u^{(o)} | s_u^{(o)}, \hat{\lambda}_{u-1}}[Z_u(\lambda)]] \\
 &= \frac{1}{N_t^{(o)}} \sum_{u=1}^{N_t^{(o)}} \mathbb{E}_{s_u^{(o)} | \hat{\lambda}_{u-1}}[\mathbf{1}(s_u^{(o)} > \lambda)] \\
 &= \frac{1}{N_t^{(o)}} \sum_{u=1}^{N_t^{(o)}} \text{FPR}(\lambda) \\
 &= \text{FPR}(\lambda)
 \end{aligned}$$

□

Having an unbiased estimator solves one part of the problem. In addition we need confidence intervals on this estimate that are valid for anytime and for the choices of λ considered. Due to the dependence between the samples we cannot directly apply similar results developed for quantile estimation in the i.i.d. setting (Howard & Ramdas, 2022). Fortunately, part of this problem has been addressed in (Balsubramani, 2015), where they provide anytime valid confidence intervals when the estimators form a martingale sequence. We restate this result in the following lemma 3 and then building upon this result, in the next lemma 4 we derive such confidence intervals for our setting.

Lemma 3. ((Balsubramani, 2015)) Let \overline{M}_t be a martingale and suppose $|\overline{M}_t - \overline{M}_{t-1}| \leq \rho_t$ for constants $\{\rho_t\}_{t \geq 1}$, let $m_0 = \min_{t \geq 1} |\overline{M}_t|$. Fix any $\delta \in (0, 1)$, and let $t_0 = \min\{u : \sum_{i=1}^u \rho_i^2 \geq 173 \log(\frac{4}{\delta})\}$ then,

$$\mathbb{P}\left(\exists t \geq t_0 : |\overline{M}_t| \geq \sqrt{3 \left(\sum_{i=1}^t \rho_i^2 \right) \left(2 \log \log \left(\frac{3 \sum_{i=1}^t \rho_i^2}{m_0} \right) + \log \left(\frac{2}{\delta} \right) \right)}\right) \leq \delta \quad (7)$$

Proof. This lemma is a restatement of theorem 4 in (Balsubramani, 2015). For proof details please see (Balsubramani, 2015). □

In the next lemma we show that the sums of $Z_u(\lambda)$ form a martingale sequence, allowing us to apply the results from the above lemma (3) and then we generalize it to all λ in some finite set.

Lemma 4. (Anytime valid confidence intervals on FPR) Let $X_t^{(o)} = \{x_1^{(o)}, \dots, x_{N_t^{(o)}}^{(o)}\}$ be the samples drawn from D_{ood} till round t and let $S_t^{(o)} = \{s_1^{(o)}, \dots, s_{N_t^{(o)}}^{(o)}\}$ be the scores of these points, let $c_t = 1 - \beta_t + \frac{\beta_t}{p^2}$, $\beta_t = \frac{N_t^{(o,p)}}{N_t^{(o)}}$ and $N_t^{(o,p)}$ is the

number of points sampled using importance sampling until time t and $\nu \in (0, 1)$ is a discretization parameter set by the user. Let $\Lambda = \{\Lambda_{\min}, \Lambda_{\min} + \nu, \dots, \Lambda_{\max}\}$. Let $n_0 = \min\{u : c_u N_u^{(o)} \geq 173 \log(\frac{4}{\delta})\}$ and t_0 be such that $N_{t_0}^{(o)} \geq n_0$, then for any $\delta \in (0, 1)$,

$$\mathbb{P}\left(\exists t \geq t_0 : \sup_{\lambda \in \Lambda} \widehat{FPR}(\lambda, t) - FPR(\lambda) \geq \psi(t, \delta)\right) \leq \delta \quad (8)$$

for,

$$\psi(t, \delta) = \sqrt{\frac{3c_t}{N_t^{(o)}} \left[2 \log \log \left(\frac{3c_t N_t^{(o)}}{2} \right) + \log \left(\frac{2|\Lambda|}{\delta} \right) \right]} \quad (9)$$

Proof. First we show that we have a martingale sequence as follows,

Let $M_t(\lambda) = \sum_{u=1}^{N_t^{(o)}} Z_u(\lambda)$, and consider the centered random variables,

$$\bar{M}_t(\lambda) = M_t(\lambda) - \mathbb{E}[M_t(\lambda)] \quad \text{and} \quad \bar{Z}_t(\lambda) = Z_t(\lambda) - \text{FPR}(\lambda)$$

Let \mathcal{F}_t be the σ -algebra of events till time t i.e. $(s_1^{(o)}, i_1^{(o)}), \dots, (s_{t-1}^{(o)}, i_{t-1}^{(o)}), (s_t^{(o)}, i_t^{(o)})$.

It is easy to see that $\mathbb{E}[\bar{M}_t] \leq \frac{1}{p} < \infty$ and \bar{M}_t is \mathcal{F}_t -measurable for all $t > 1$. Further, we can see,

$$\mathbb{E}[\bar{M}_t(\lambda) | \mathcal{F}_{t-1}] = \mathbb{E}[\bar{Z}_t(\lambda) + \bar{M}_{t-1}(\lambda) | \mathcal{F}_{t-1}] = \mathbb{E}[\bar{Z}_t(\lambda) | \mathcal{F}_{t-1}] + \mathbb{E}[\bar{M}_{t-1}(\lambda) | \mathcal{F}_{t-1}] = \bar{M}_{t-1}(\lambda)$$

Since, $\mathbb{E}[\bar{Z}_t(\lambda) | \mathcal{F}_{t-1}] = 0$ and $\mathbb{E}[\bar{M}_{t-1}(\lambda) | \mathcal{F}_{t-1}] = \bar{M}_{t-1}(\lambda)$. Thus we have that \bar{M}_t is a martingale sequence. Further, we also have the following,

$$|\bar{M}_t(\lambda) - \bar{M}_{t-1}(\lambda)| \leq \begin{cases} 1 & \text{if } i_t^{(o)} = 0 \\ \frac{1}{p} & \text{if } i_t^{(o)} = 1 \end{cases}$$

Let $\beta_t \in (0, 1)$ be the fraction of OOD points sampled using probability p till round t . Let $N_t^{(o)}$ be the total number of points OOD points sampled till round t and $N_t^{(o,p)}$ be the points sampled from importance sampling, then $\beta_t = \frac{N_t^{(o,p)}}{N_t^{(o)}}$.

Let $c_t = 1 - \beta_t + \frac{\beta_t}{p^2}$. We know p and the number of points sampled with importance sampling, without importance sampling we know β_t, c_t are at time t . Applying lemma 3 we get the following result for a given λ ,

$$\mathbb{P}\left(\exists t \geq t_0 : \bar{M}_t(\lambda) \geq \sqrt{3(c_t N_t^{(o)}) \left(2 \log \log (3c_t N_t^{(o)}) + \log \left(\frac{2}{\delta} \right) \right)}\right) \leq \delta \quad (10)$$

$$\mathbb{P}\left(\exists t \geq t_0 : \widehat{FPR}(\lambda, t) - FPR(\lambda, t) \geq \sqrt{\frac{3c_t}{N_t^{(o)}} \left(2 \log \log (3c_t N_t^{(o)}) + \log \left(\frac{2}{\delta} \right) \right)}\right) \leq \delta \quad (11)$$

Doing the union bound for the failure probability over all $\lambda \in \Lambda$, (where $|\Lambda| < \infty$) gives us the result. \square

Our performance guarantees in the main theorem 2 are based on $\psi(t, \delta)$ becoming smaller than certain values. In the next lemma we derive bound on $N_t^{(o)}$ such that $\psi(t, \delta)$ is at most μ and use it in the proof of the main theorem 2.

Lemma 5. Let $\psi(t, \delta) = \sqrt{\frac{3c_t}{N_t^{(o)}} \left(2 \log \log (3c_t N_t^{(o)}) + \log \left(\frac{2|\Lambda|}{\delta} \right) \right)}$, and let there be a constant C_0 and time T_0 , such that $\beta_t \leq C_0 p^2$ for all $t \geq T_0$ (worst case $T_0 = 1$ and $C_0 = 1/p^2$). Then $\psi(t, \delta) \leq \mu$ for any $t > T_\mu > T_0$ such that $N_{T_\mu}^{(o)} = \frac{10(C_0+1)}{\mu^2} \log \left(\frac{|\Lambda|}{\delta} \log \left(\frac{5(C_0+1)}{\mu} \right) \right)$.

Proof. First we write a simplified form of $\psi(t, \delta)$ for all $t > T_0$ as follows,

$$\psi(t, \delta) = \sqrt{\frac{3(C_0 + 1)}{N_t^{(o)}} \left(2 \log \log \left(3(C_0 + 1) N_t^{(o)} \right) + \log \left(\frac{2|\Lambda|}{\delta} \right) \right)}$$

In the above equation we used the bound on $\beta_t \leq C_0 p^2$ in the equation $c_t = 1 - \beta_t + \beta_t/p^2$ leading to $c_t \leq C_0 + 1$, Now, for brevity let $a_1 = 3(C_0 + 1)$ and $a_2 = 2|\Lambda|$ and rewrite $\psi(t, \delta)$ as follows,

$$\psi^2(t, \delta) = \frac{a_1}{N_t^{(o)}} \left(2 \log \log \left(a_1 N_t^{(o)} \right) + \log \left(\frac{a_2}{\delta} \right) \right) \leq \frac{2a_1}{N_t^{(o)}} \left(\log \left(\frac{a_2}{\delta} \log \left(a_1 N_t^{(o)} \right) \right) \right)$$

We want to find $N_t^{(o)}$ such that $\psi^2(t, \delta) \leq \mu^2$. It is difficult to directly invert this function. To get a bound on $N_t^{(o)}$ we first assume the following form for it with unknown constants $b_1, b_2, b_3 > 0$ and then figure out the constants by simplifying $\psi^2(N_t^{(o)})$ and constraining it to be at most μ^2 .

$$\text{Let } N_{T_\mu}^{(o)} = \frac{b_1}{\mu^2} \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)$$

$$\begin{aligned} \psi^2(T_\mu, \delta) &\leq \frac{2a_1}{N_{T_\mu}^{(o)}} \log \left[\frac{a_2}{\delta} \log(a_1 N_{T_\mu}^{(o)}) \right] \\ &= \frac{2a_1 \mu^2}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\frac{a_2}{\delta} \log \left\{ \frac{a_1 b_1}{\mu^2} \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right) \right\} \right] \\ &\stackrel{a}{\leq} \frac{2a_1 \mu^2}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\frac{a_2}{\delta} \log \left\{ \frac{a_1 b_1}{\mu^2} \log \left(\frac{a_2 b_2}{b_3 \delta \mu} \right) \right\} \right] \\ &= \frac{2a_1 \mu^2}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\frac{a_2}{\delta} \log \left\{ \frac{a_1 b_1}{\mu^2} \log \left(\frac{a_2 b_2}{b_3 \delta \mu} \right) \right\} \right] \\ &\stackrel{b}{\leq} \frac{2a_1 \mu^2}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\frac{a_2}{\delta} \log \left\{ \frac{a_1 b_1 a_2 b_2}{\mu^2 b_3 \delta \mu} \right\} \right] \\ &= \frac{2a_1 \mu^2}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\frac{a_2}{\delta} \log \left\{ \frac{a_1 b_1 a_2 b_2}{\mu^2 b_3 \delta \mu} \right\} \right] \\ &= \frac{2a_1 \mu^2}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\frac{a_2}{\delta} \log \left\{ \frac{a_1 b_1 a_2}{b_3 b_2^2 \delta} \left(\frac{b_2}{\mu} \right)^3 \right\} \right] \\ &\stackrel{c}{\leq} \frac{2a_1 \mu^2}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\frac{a_2}{\delta} \frac{a_1 b_1 a_2}{b_3 b_2^2 \delta} \log \left\{ \left(\frac{b_2}{\mu} \right)^3 \right\} \right] \\ &= \frac{2a_1 \mu^2}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\frac{a_2}{\delta} \frac{3a_1 b_1 a_2}{b_3 b_2^2 \delta} \log \left(\frac{b_2}{\mu} \right) \right] \\ &= \frac{2a_1 \mu^2}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\left(\frac{a_2}{b_3 \delta} \right)^2 \frac{3a_1 b_1 b_3}{b_2^2} \log \left(\frac{b_2}{\mu} \right) \right] \\ &\stackrel{d}{\leq} \frac{2a_1 \mu^2 \frac{3a_1 b_1 b_3}{b_2^2}}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\left(\frac{a_2}{b_3 \delta} \right)^2 \log \left(\frac{b_2}{\mu} \right) \right] \\ &\stackrel{e}{\leq} \frac{2a_1 \mu^2 \cdot 2 \frac{3a_1 b_1 b_3}{b_2^2}}{b_1 \log \left(\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right)} \log \left[\frac{a_2}{b_3 \delta} \log \left(\frac{b_2}{\mu} \right) \right] \end{aligned}$$

$$= \frac{12\mu^2 a_1^2 b_3}{b_2^2}.$$

The inequalities $\textcolor{red}{a}, \textcolor{teal}{b}$ follow from $\log(x) \leq x$ for any $x > 0$.

The inequality $\textcolor{blue}{c}$ comes from $\log(ax) \leq a \log(x)$ for any $a > 2, x > 2$. We use $a = \frac{a_1 b_1 a_2}{b_3 b_2^2 \delta}$ and $x = \left(\frac{b_2}{\mu}\right)^3$, this enforces the following constraints,

$$\frac{b_2}{\mu} > 2^{1/3} \quad (12)$$

$$\frac{a_1 b_1 a_2}{b_3 b_2^2 \delta} > 2 \quad (13)$$

For $\textcolor{blue}{d}$ we again use $\log(ax) \leq a \log(x)$ with $a = \frac{3a_1 b_1 b_3}{b_2^2}$ and $x = \left(\frac{a_2}{b_3 \delta}\right)^2 \log\left(\frac{b_2}{\mu}\right)$, this enforces the following constraints,

$$\frac{3a_1 b_1 b_3}{b_2^2} > 2 \quad (14)$$

$$\left(\frac{a_2}{b_3 \delta}\right)^2 \log\left(\frac{b_2}{\mu}\right) > 2 \quad (15)$$

Lastly, $\textcolor{red}{e}$ follows by using $\log(x^a y) \leq a \log(xy)$ for any $x > 0, a > 1, y > 1$. For this we use $x = \frac{a_2}{b_3 \delta}$ and $y = \log\left(\frac{b_2}{\mu}\right)$, leading the following constraints,

$$\log\left(\frac{b_2}{\mu}\right) > 1 \quad (16)$$

For $\psi^2(T_\mu) \leq \mu^2$, we need

$$12a_1^2 b_3 \leq b_2^2 \quad (17)$$

Let $b_3 = 2, b_1 = 10a_1, b_2 = 5a_1$ then the constraints 12, 13, 14, 15, 16 and 17 are satisfied (when $|\Lambda| \geq 10$) for any $\mu \in (0, 1), \delta \in (0, 1)$. Thus we have,

$$\psi(T_\mu, \delta) \leq \mu \text{ for } N_{T_\mu} = \frac{10(C_0 + 1)}{\mu^2} \log\left(\frac{|\Lambda|}{\delta} \log\left(\frac{5(C_0 + 1)}{\mu}\right)\right) \quad (18)$$

□

Theorem 2. Let $\alpha \in (0, 1)$ and $\delta \in (0, 1)$. Let $x_t \stackrel{i.i.d.}{\sim} (1 - \gamma)\mathcal{D}_{id} + \gamma\mathcal{D}_{ood}$. If Algorithm 1 uses the optimization problem (P2) to find the thresholds with the upper confidence term $\psi(N_t^{(o)}, \delta)$ given by equation (5), then, with probability at least $1 - \delta$,

1. (Controlled FPR) For all $t \geq 1$, $FPR(\hat{\lambda}_t) \leq \alpha$.
2. (Time to reach feasibility) Let T_f be such that $N_{T_f}^{(o)} \geq \frac{C_1}{\alpha^2} \log\left(\frac{C_2}{\delta} \log\left(\frac{C_3}{\alpha}\right)\right)$ then for any $t \geq T_f$ the algorithm will find a feasible threshold, $\hat{\lambda}_t$ such that $\widehat{FPR}(\hat{\lambda}_t) + \psi(N_t^{(o)}) \leq \alpha$.
3. (Time to reach η -Optimality) Let T_{opt} be such that $N_{T_{opt}}^{(o)} \geq \frac{4C_1}{\eta^2} \log\left(\frac{C_2}{\delta} \log\left(\frac{2C_3}{\eta}\right)\right)$ and $\widehat{FPR}(\hat{\lambda}_{T_{opt}}) \geq [FPR(\lambda^*) - \eta/2, \alpha]$, then then for any $t \geq T_{opt}$ if the $\hat{\lambda}_t$ satisfy the η -Optimality condition in definition 1.

Proof. Controlled FPR: This follows from the fact the algorithm uses $\psi(t, \delta)$ that are valid for all t and the choices of λ it considers.

Time to reach feasibility:

Applying lemma 5 with $\mu = \alpha$ gives bound on N_{T_f} with $C_1 = 10(C_0 + 1), C_2 = |\Lambda|, C_3 = 5(C_0 + 1)$.

Time to reach η -optimality

We know, $\text{FPR}(\lambda^*) = \alpha$, and it is given that $\widehat{\text{FPR}}(\hat{\lambda}_t) \in [\text{FPR}(\lambda^*) - \eta/2, \alpha]$

$$\text{FPR}(\hat{\lambda}_t) \in [\text{FPR}(\lambda^*) - \eta/2 - \psi(t, \delta), \alpha]$$

this means $\text{FPR}(\hat{\lambda}_t) \geq \text{FPR}(\lambda^*) - \eta/2 - \psi(t, \delta)$

$$\text{FPR}(\lambda^*) - \text{FPR}(\hat{\lambda}_t) \leq \eta/2 + \psi(t, \delta)$$

If $\psi(t, \delta) \leq \eta/2$ we have, $\text{FPR}(\lambda^*) - \text{FPR}(\hat{\lambda}_t) \leq \eta$. Thus applying we want to find t for which $\psi(t, \delta) = \eta/2$. Applying lemma 5 with $\mu = \eta/2$ gives bound on N_{T_f} with $C_1 = 40(C_0 + 1)$, $C_2 = |\Lambda|$, $C_3 = 10(C_0 + 1)$. \square

This concludes the proofs of the main results. Next, we present additional experiments on synthetic and real datasets.

8. Additional Experiments and Details

In the simulations we study the effect of changing γ , using different window sizes and the case when the In-Distribution shifts. For the real data experiments we study the performance of the methods under different settings with different scoring functions on CIFAR-10 and CIFAR-100 as In-Distribution datasets.

8.1. Searching for constants in LIL-Heuristic

As we saw in the main paper the theoretical LIL bound in 5 has constants that can be pessimistic in practice. We get around this by using a LIL-Heuristic bound which has the same form as in equation (5) but with different constants in particular we consider the form in equation **LIL-Heuristic**. We find the constants C_1, C_2, C_3 using a simulation on estimating the bias of a coin with different constants and picking the ones so that the observed failure probability is below 5%.

$$C_1 \sqrt{\frac{c_t}{N_t^{(o)}} \left(\log \log (C_2 c_t N_t^{(o)}) + \log \left(\frac{C_3}{\delta} \right) \right)}. \quad (\text{LIL-Heuristic})$$

Specifically, we keep $C_3 = 1$, and run for $\delta = 0.01, 0.05, 0.1, 0.2, 0.3, 0.4$ with varying C_1 from 0.1 to 0.9 and C_2 from 1.5 to 4.75. For each choice of δ, C_1, C_2 , we toss an unbiased coin (mean $p = 0.5$) for $T = 10k$ times. For each choice of $t = 1, 2, \dots, T$, we compute the empirical mean of the coin \hat{p} and define it as a failure if $p \notin [\hat{p} - \psi(t), \hat{p} + \psi(t)]$. We run this process for 100 times and compute the average failure probability for each choice of $t = 1, 2, \dots, T$. We then pick the constant so that the observed average probability is below 5%. Throughout the paper, we use $C_1 = 0.5$ and $C_2 = 0.75$.

8.2. Additional Simulations

In addition to the main paper, we run the following simulation to better understand the performance of the system. We use $\alpha = 0.05$, $\delta = 0.2$, and importance sampling probability $p = 0.2$ through all the simulations.

1. **Effect of window size.** We have shown window size = $10k$ for the main paper. Here we also show the window size = $30k$ or $50k$. When the window size increases, we see that the empirical FPR with different UCBs is closer to the 5%, but it takes a longer time to react to a distribution shift. To be consistent with the main paper, we simulate the OOD and ID scores using a mixture of two Gaussians $\mathcal{N}_{id}(\mu = 5.5, \sigma = 4)$ and $\mathcal{N}_{ood}(\mu = -6, \sigma = 4)$ with $\gamma = 0.1$. To simulate distribution change we change the OOD distribution to $\mathcal{N}_{ood}(\mu = -5, \sigma = 4)$ at time $t = 55k$. The result is in figure 7.
2. **Effect of γ .** We show that changing the mixing ratio γ of ID and OOD samples does not affect the control of FPR. We show two settings for changing γ . First, the gamma changes from 0.1 to 1 when $t = 55k$. Second, we show γ gradually changes from 0.1 to 1 starting from $t = 20k$ and ends at $t = 80k$, with a step size of 0.1. We see that our system is able to control FPR with the change of γ . We simulate the OOD and ID scores using a mixture of two Gaussians $\mathcal{N}_{id}(\mu = 5.5, \sigma = 4)$ and $\mathcal{N}_{ood}(\mu = -6, \sigma = 4)$. The result is in figure 6.
3. **In-Distribution shift.** We show that our system is able to control FPR when the ID distribution is shifted. We simulate the OOD and ID scores using a mixture of two Gaussians $\mathcal{N}_{id}(\mu = 5.5, \sigma = 4)$ and $\mathcal{N}_{ood}(\mu = -5, \sigma = 4)$ with $\gamma = 0.1$. To simulate distribution change we change the ID distribution to $\mathcal{N}_{id}(\mu = 5, \sigma = 4)$ at time $t = 55k$. The result is in figure 8.

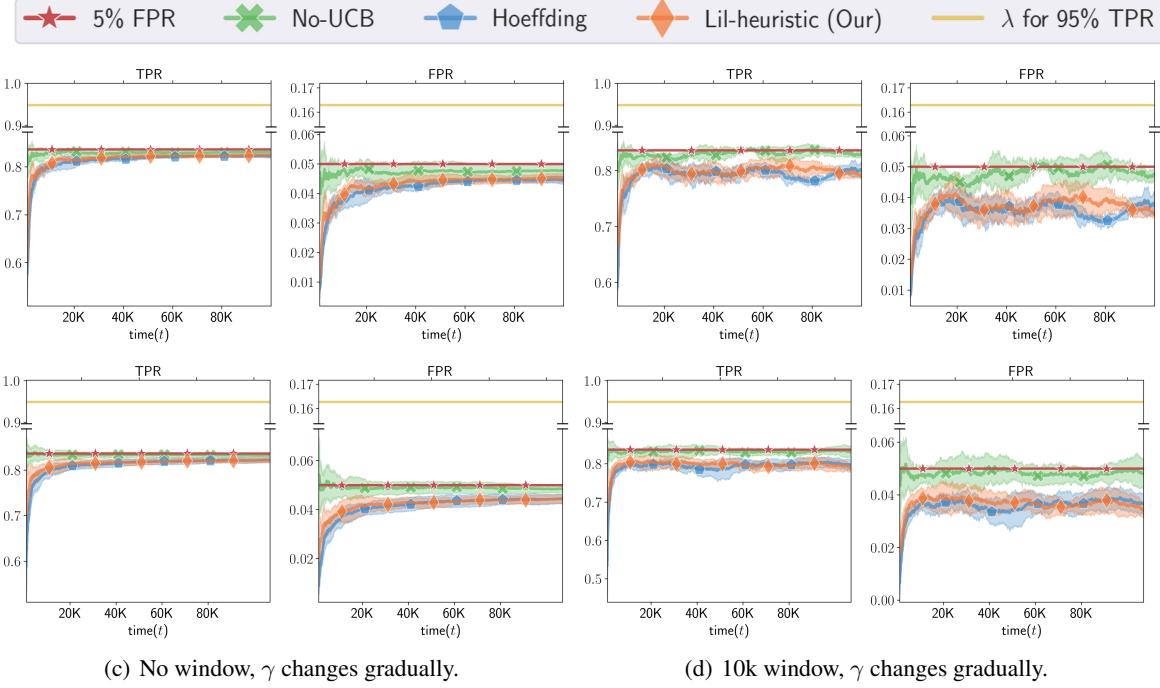
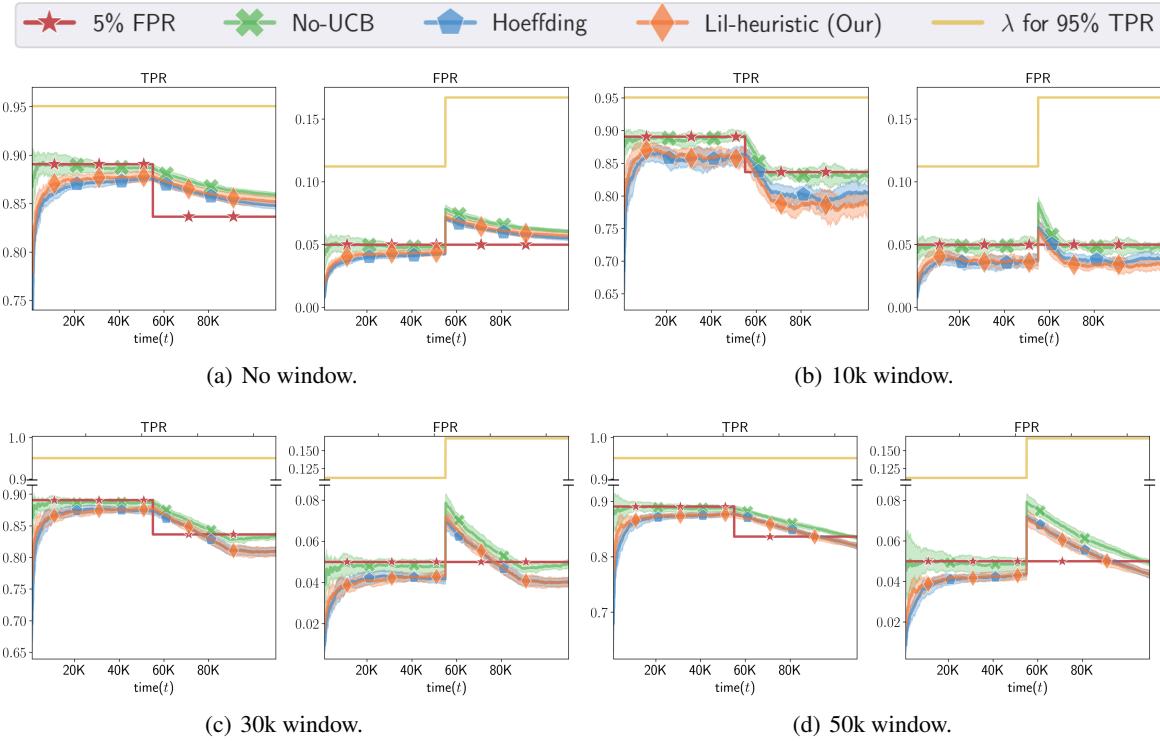

 Figure 6. Changing γ in the synthetic data.


Figure 7. Changing window size with the synthetic data.

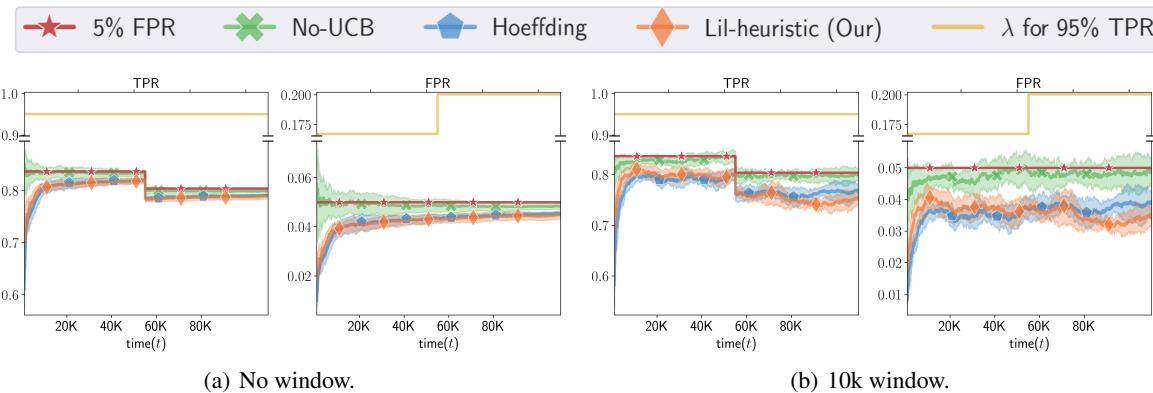


Figure 8. Changing ID distribution in synthetic data.

8.3. Additional Real OOD Datasets Experiments

We run our proposed system on different OOD scoring methods. We use CIFAR-10 or CIFAR-100 as ID datasets. We run the experiment with different window sizes and distribution shifts. In the distribution shift setting, if not specified, we use MNIST, SVHN, and Texture as the first mixture of OOD datasets, and TinyImageNet, Places365, and CIFAR-10/100 as the second mixture of OOD datasets by default. We use $\alpha = 0.05$, $\delta = 0.2$, and importance sampling probability $p = 0.2$ through all the experiments. We use a pre-trained Resnet-50 model for SSD method, and Resnet-18 for the rest of the methods.

1. **ODIN**: ODIN (Liang et al., 2017) takes the soft-max score from DNNs, and scales the score with temperature. A gradient-based input perturbation is also used for better performance. We choose temperature 1000 and input perturbation noise 0.0014, as discussed in (Liang et al., 2017). Please see figures 19 and 20 for the results with this score.
 2. **Mahalanobis Distance**: For a given point x , the Mahalanobis Distance (MDS) based score is its MD from the closest class conditional distribution. We use the MD-based score as given in (Lee et al., 2018) for detecting OOD and adversarial samples. They compute the scores using representations from various layers of DNNs and combine them to get a better scoring function. We choose input perturbation noise to be 0.0014. Please see figures 13 and 14 for the results with this score.
 3. **Energy Score**: This score was proposed in (Liu et al., 2020) and it is well aligned with the probability density of the samples, with low energy implying ID and high energy implying OOD. We choose the temperature parameter to be 1. Please see figures 11 and 12 for the results with this score.
 4. **SSD**. It is based on computing the Mahalanobis distance in the feature space of the model trained on the unlabeled in-distribution data using self-supervised learning. We use the official implementation of (Sehwag et al., 2021). We train a Resnet-50 on ID datasets using a contrastive self-supervised learning loss, SimCLR (Chen et al., 2020). When calculating the distance-based OOD scores, we use one unsupervised clustering center as the only center for ID distribution for both CIFAR-10 and CIFAR-100. Please see figures 17 and 18 for the results with this score.
 5. **Virtual-logit Match**. Virtual-logit Match (VIM) (Wang et al., 2022) combines the class-agnostic score from feature space and ID class-dependent logits. Specifically, an additional logit representing the virtual OOD class is generated from the residual of the feature against the principal space and then matched with the original logits by a constant scaling. We set the dimension of the principal space $D = 100$. Please see figures 15 and 16 for the results with this score.
 6. **K-Nearest-Neighborhood**. unlike other methods that impose a strong distributional assumption of the underlying feature space, the KNN-based method (Sun et al., 2022) explores the efficacy of non-parametric nearest-neighbor distance for OOD detection. The distance between the test sample and its k-nearest training IN sample will be used as the score for a threshold based OOD detection. We choose neighbor number $k = 50$. Please see figures 9 and 10 for the results with KNN scores.

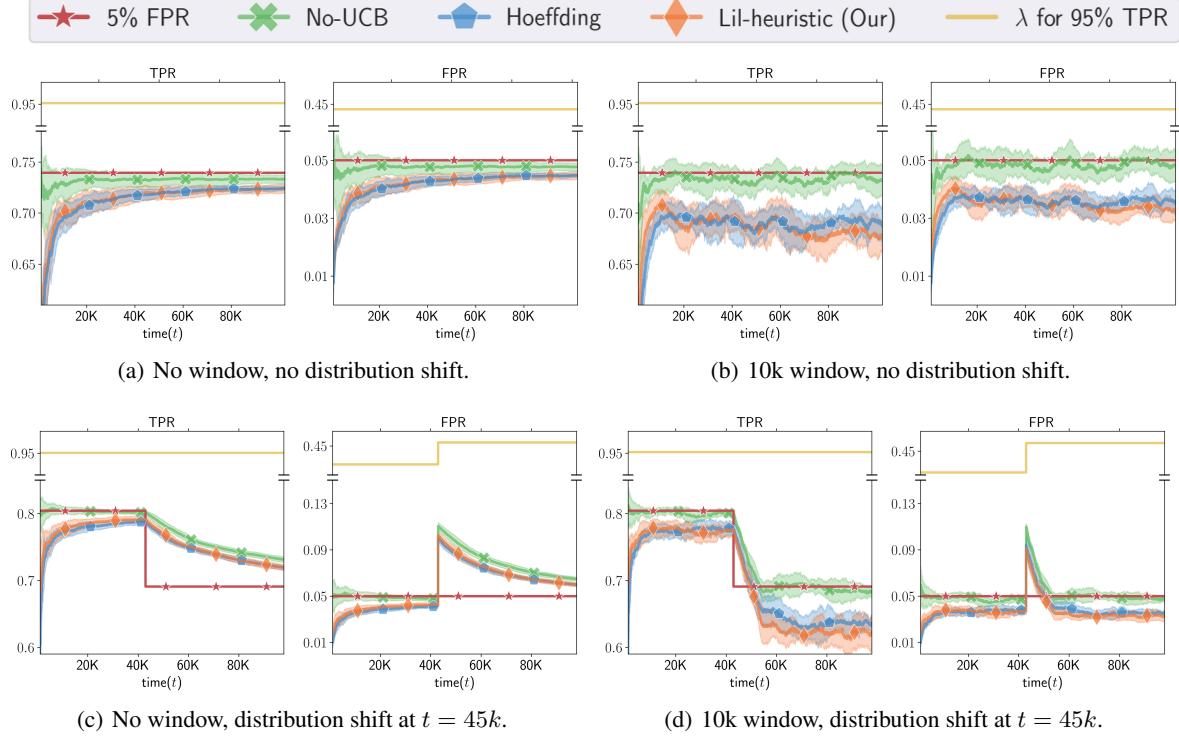


Figure 9. Results on the KNN method with Cifar-10 as the ID dataset.

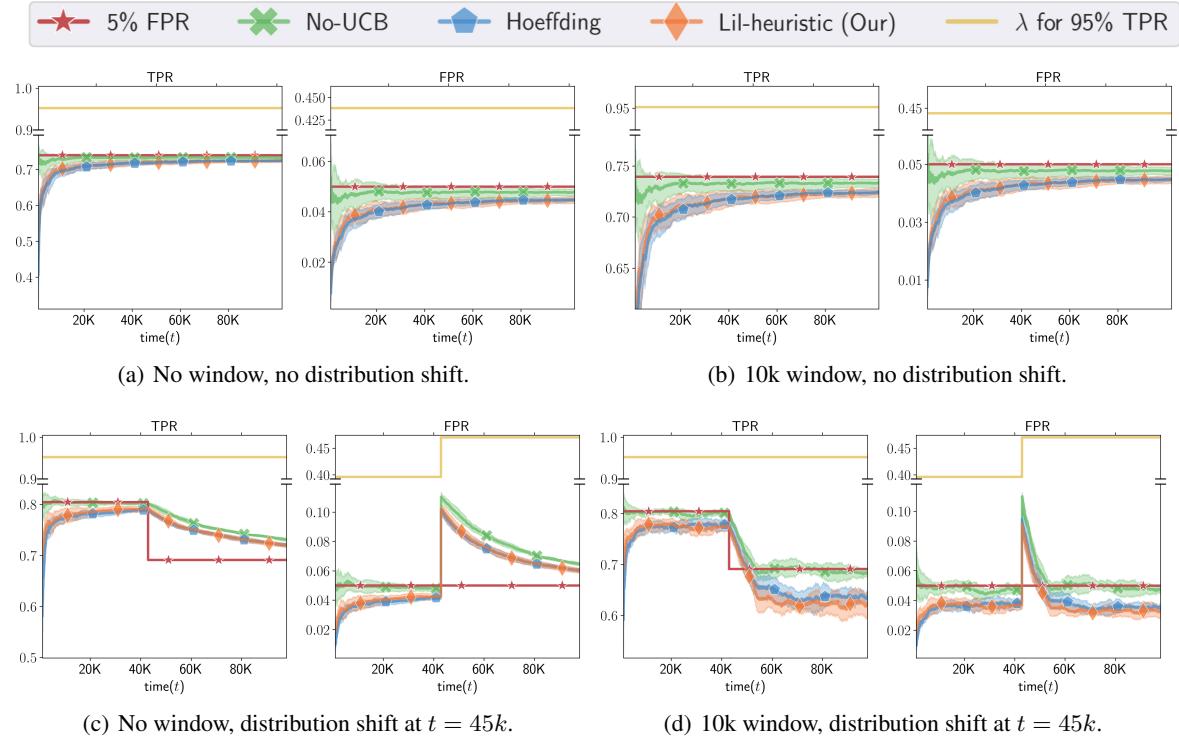


Figure 10. Results on the KNN scores with Cifar-100 as the ID dataset.

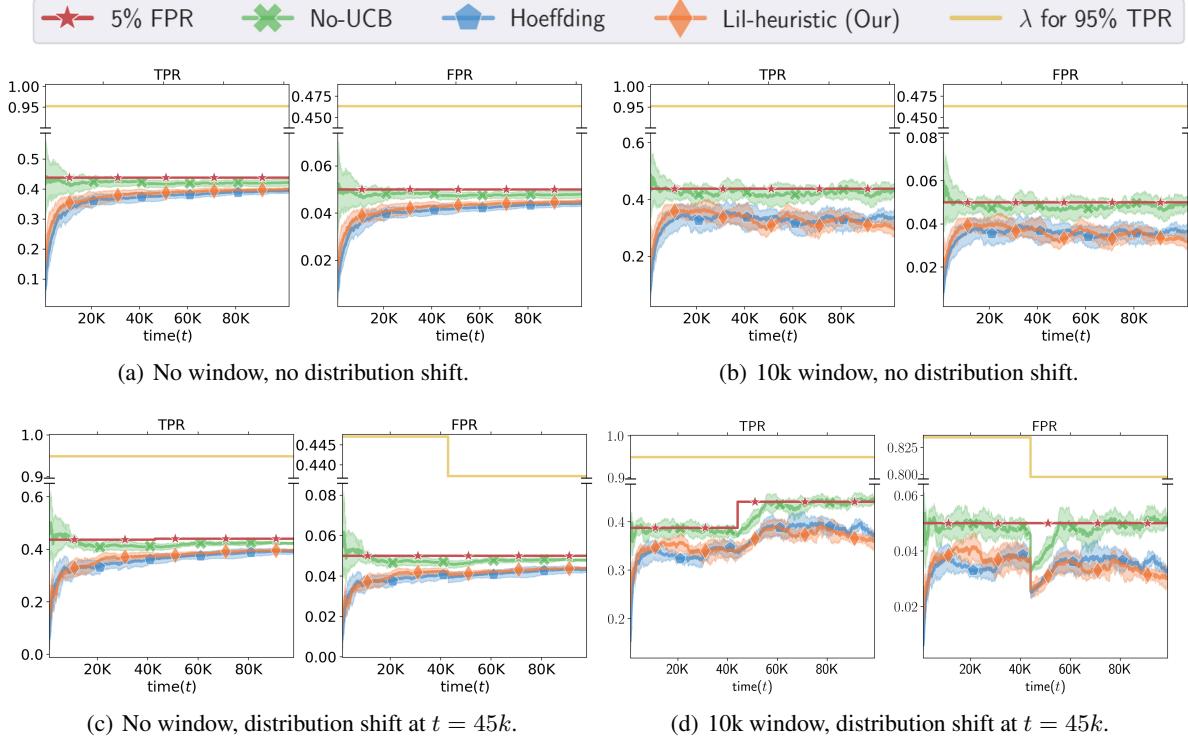


Figure 11. Results on the Energy Based Score (EBO) method with Cifar-10 as the ID dataset.

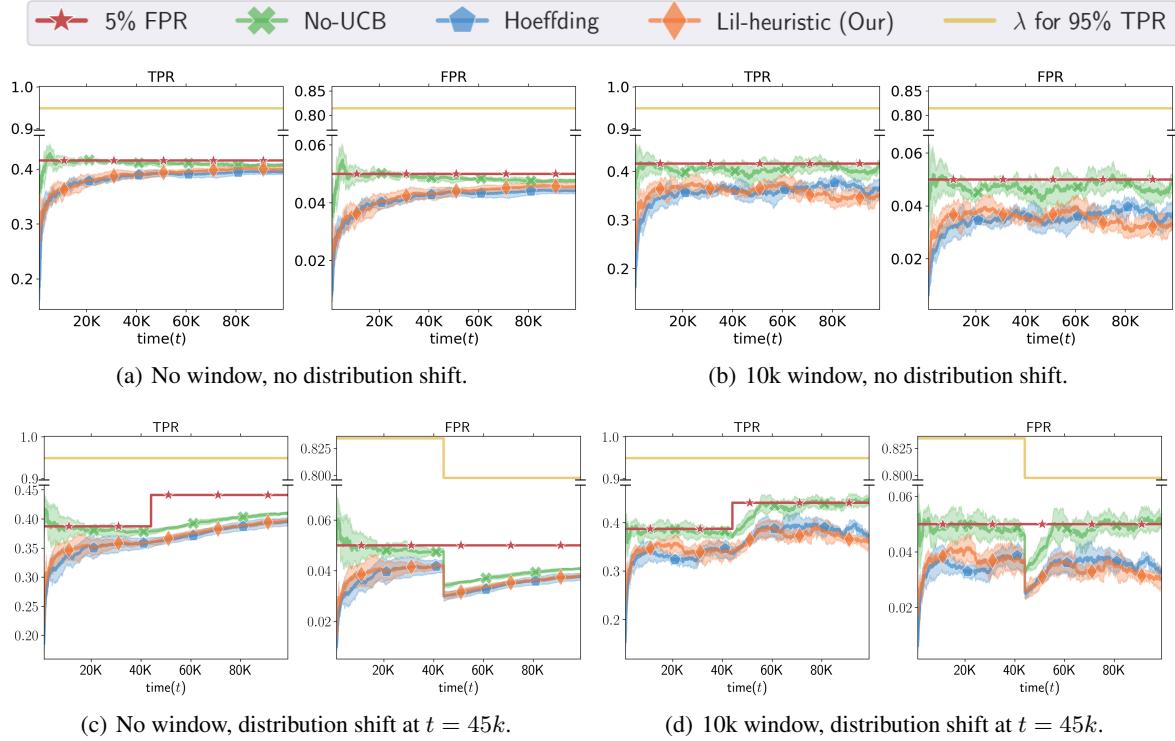


Figure 12. Results on the EBO scores with Cifar-100 as the ID dataset.

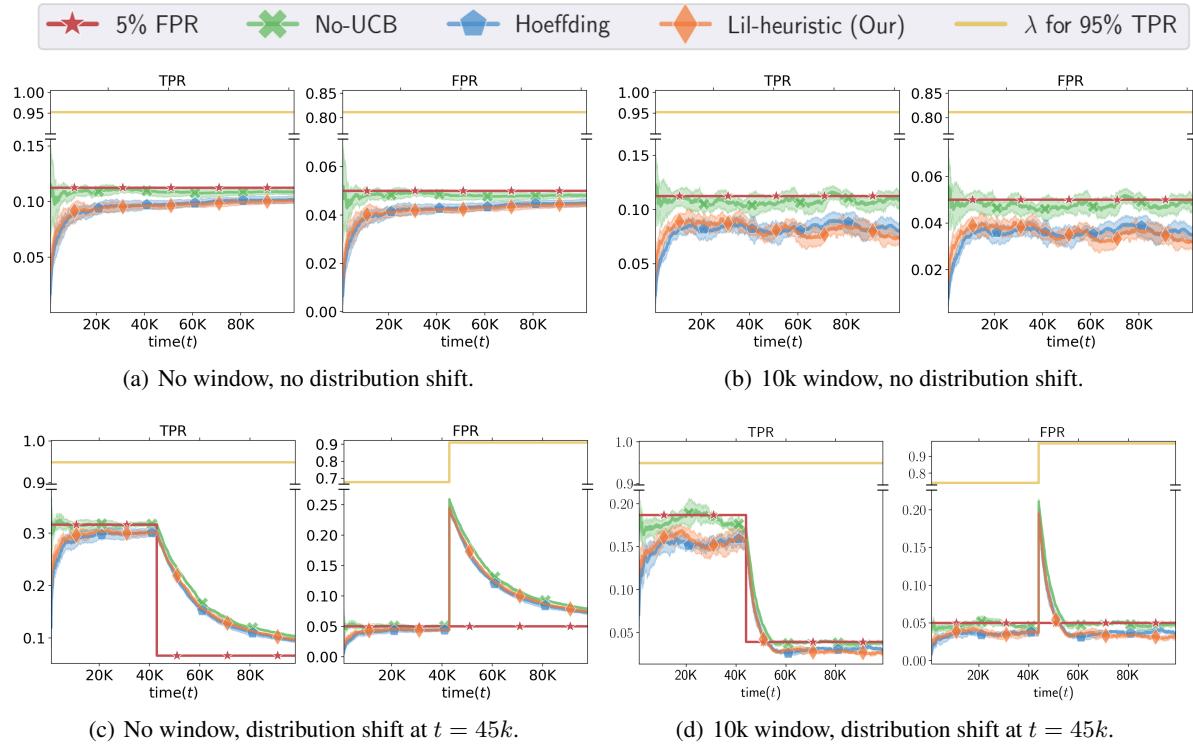


Figure 13. Results on the Mahalanobis distance (MDS) method with Cifar-10 as the ID dataset.

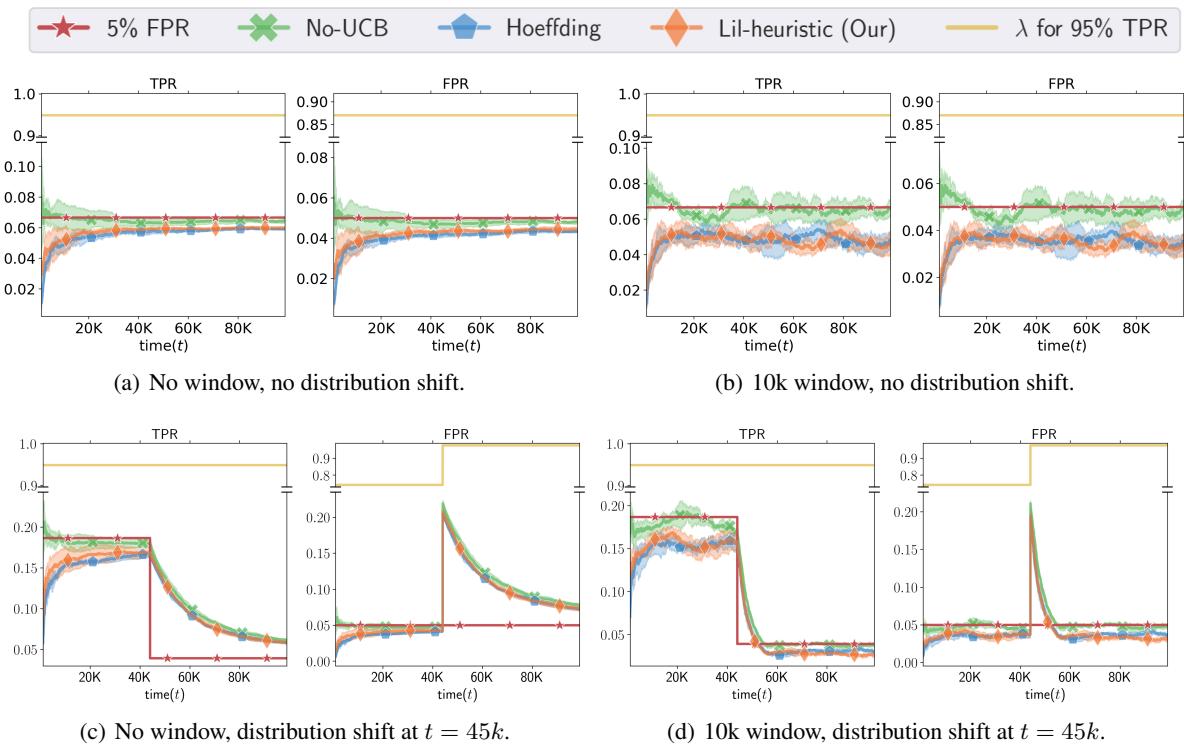


Figure 14. Results on the MDS scores with Cifar-100 as the ID dataset.

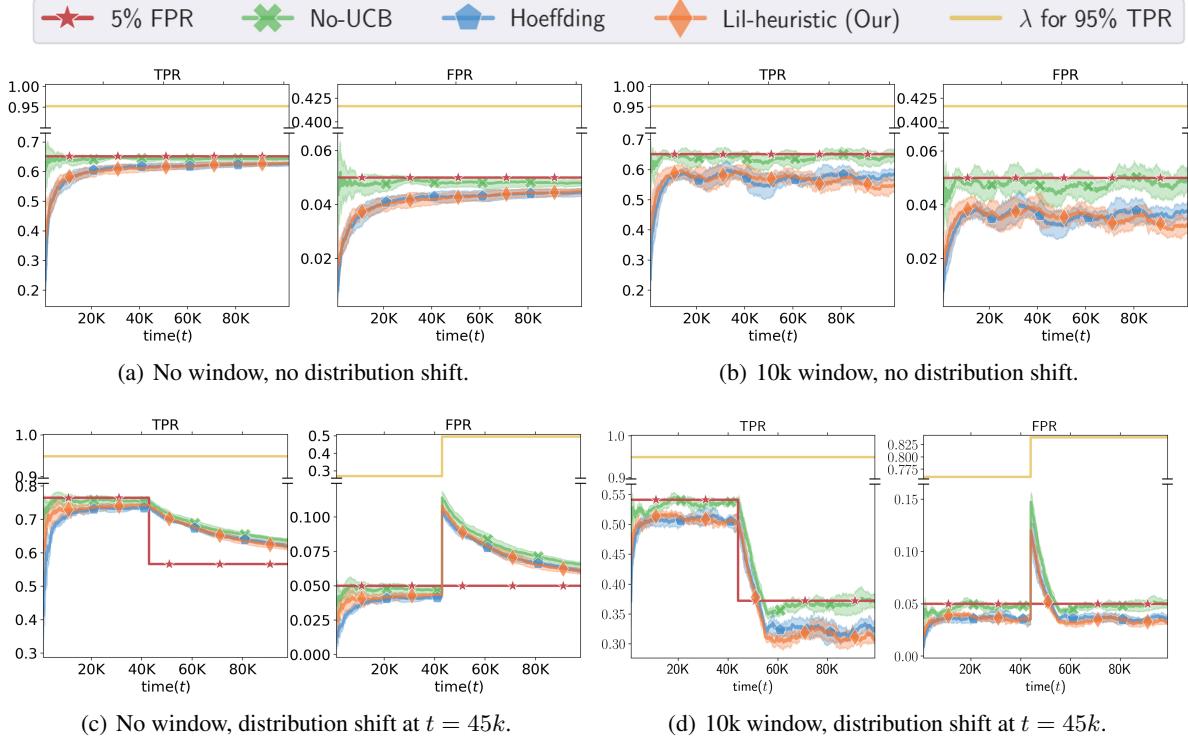


Figure 15. Results on the Virtual-logit Match (VIM) method with Cifar-10 as the ID dataset.

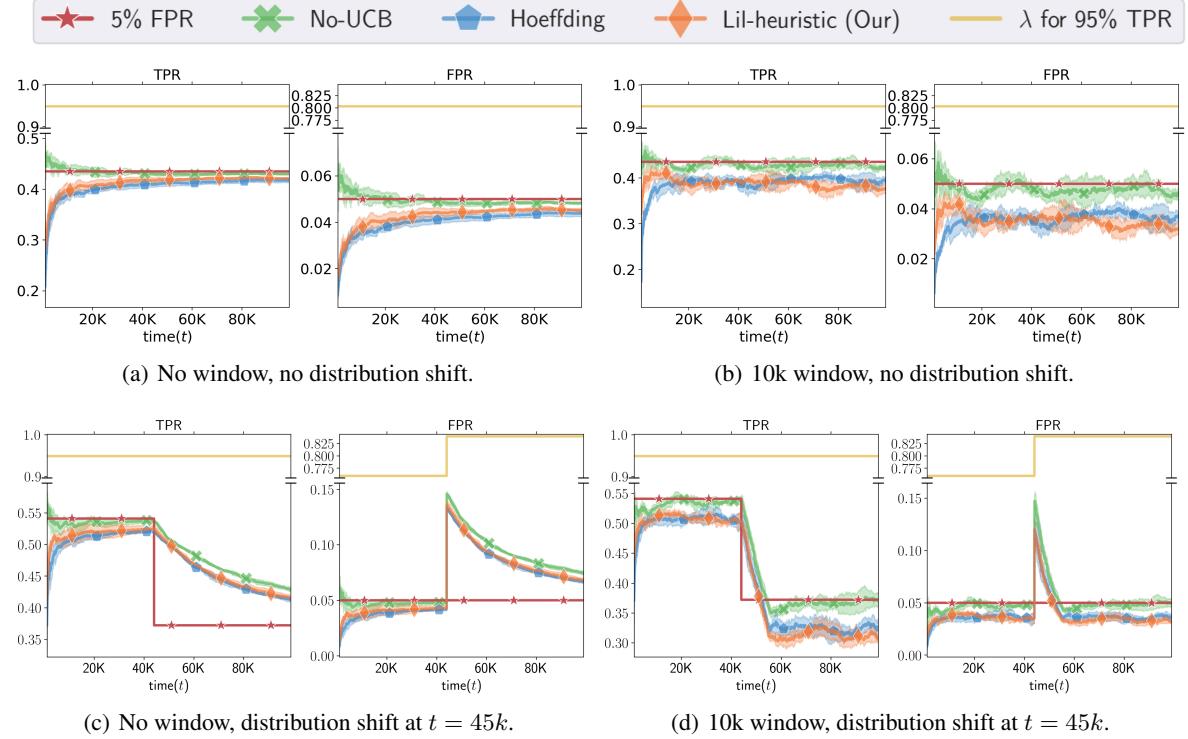


Figure 16. Results on the VIM scores with Cifar-100 as the ID dataset.

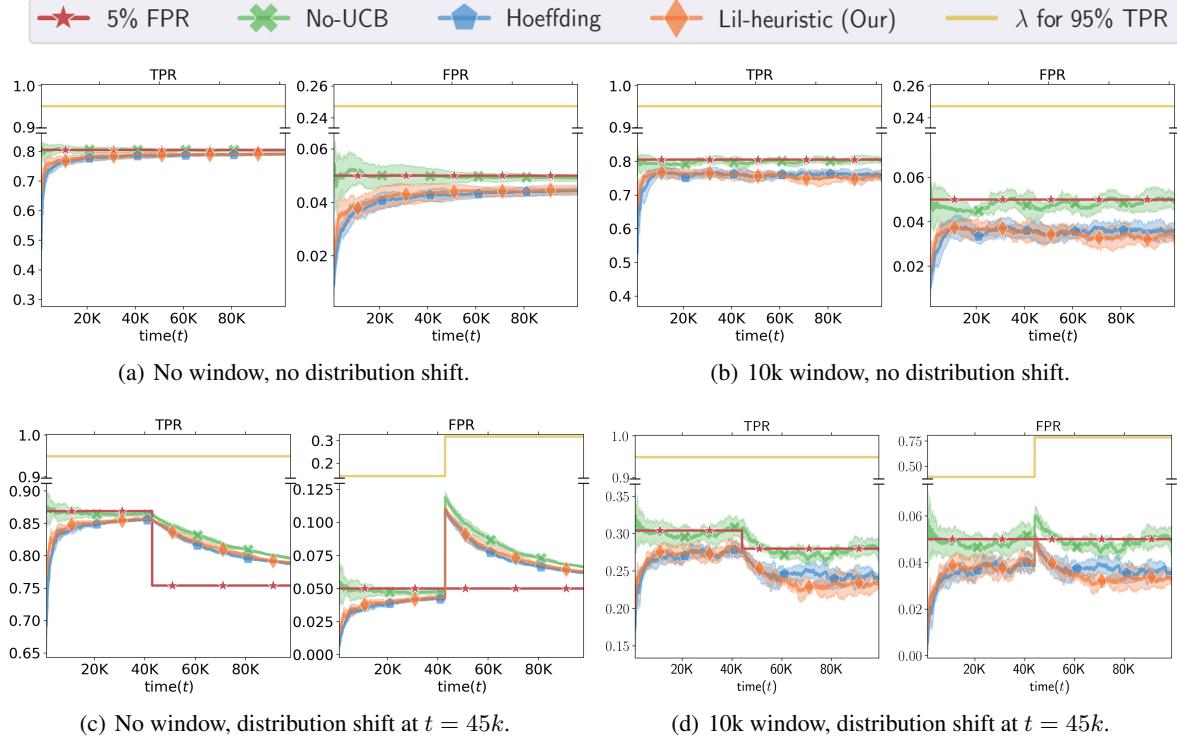


Figure 17. Results on the SSD method with Cifar-10 as the ID dataset.

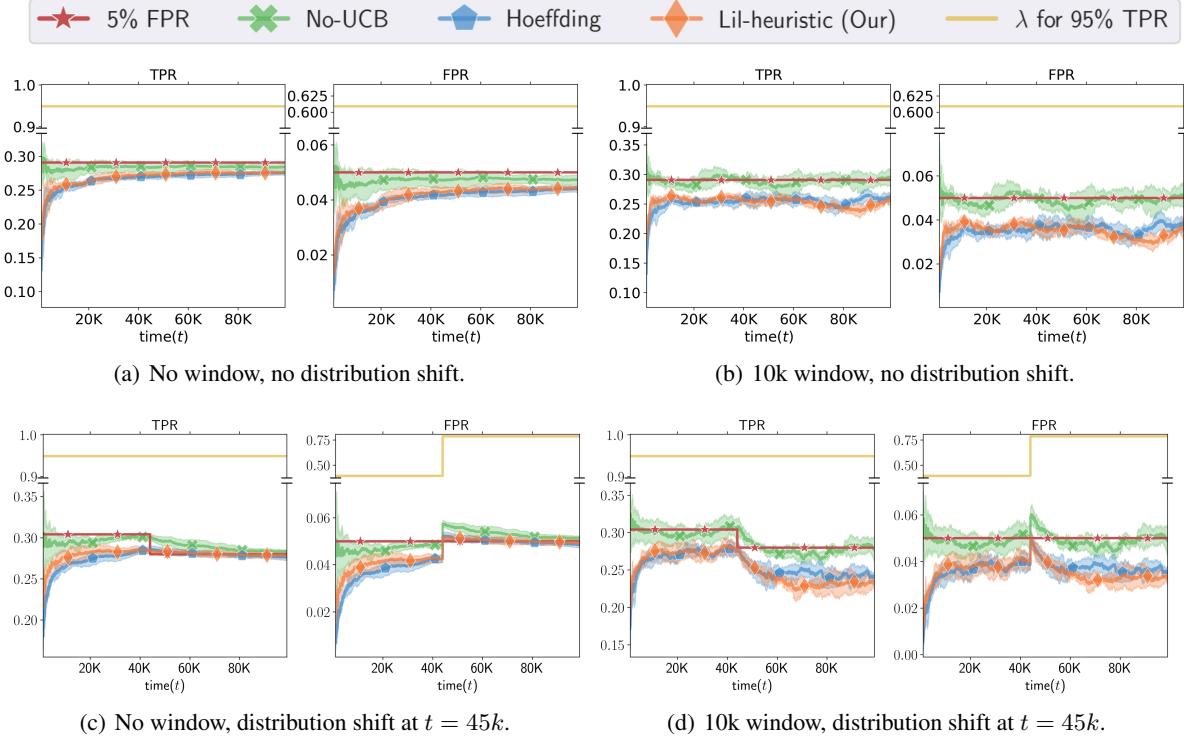


Figure 18. Results on the SSD scores with Cifar-100 as the ID dataset.

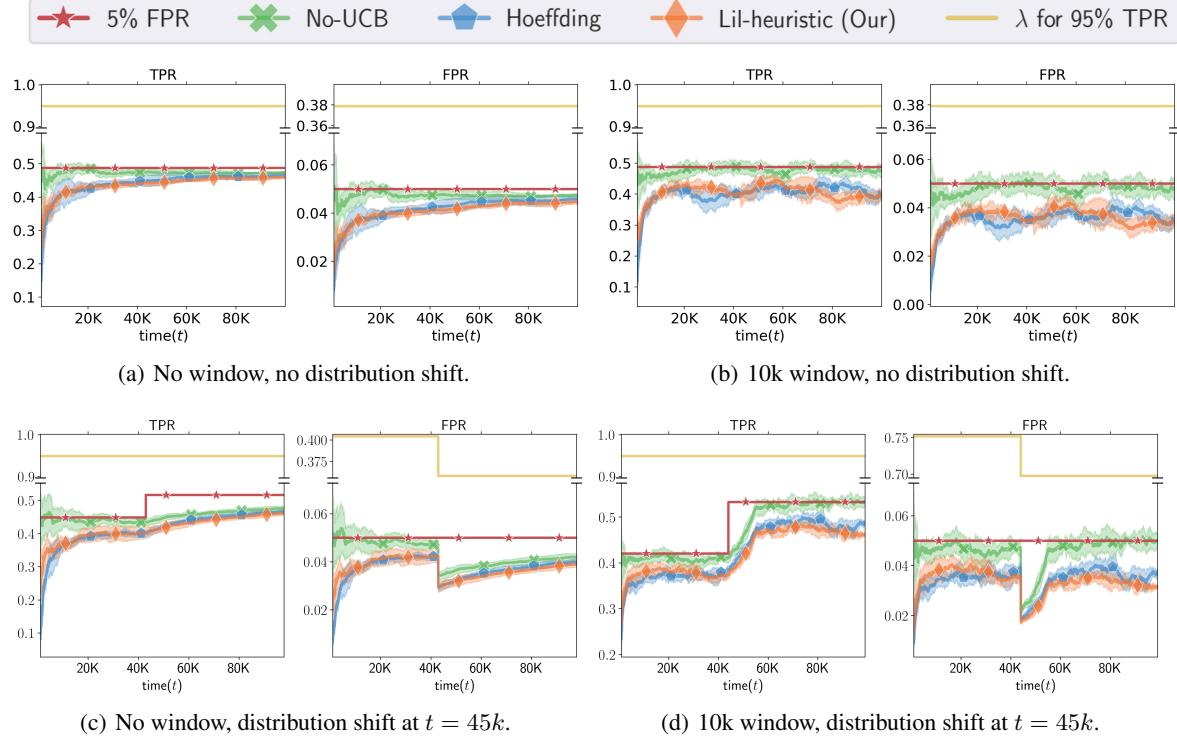


Figure 19. Results on the ODIN method with Cifar-10 as the ID dataset.

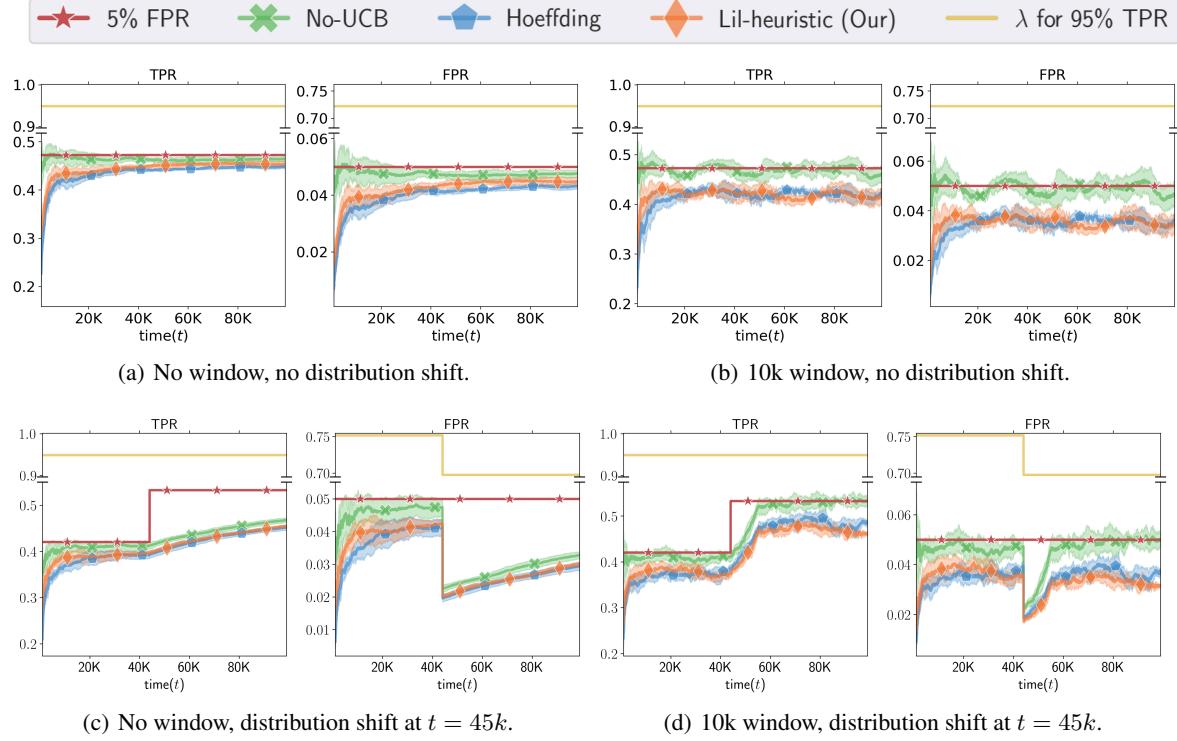


Figure 20. Results on the ODIN scores with Cifar-10 as the ID dataset.

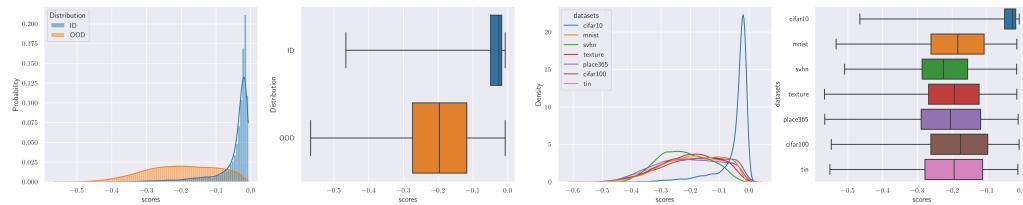


Figure 21. Scores distribution for KNN with CIFAR-10 as In-Distribution.

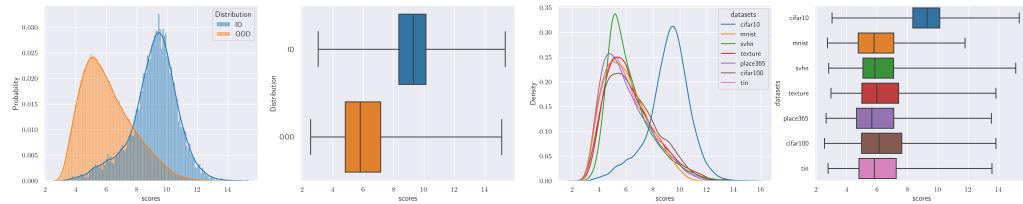


Figure 22. Scores distribution for EBO with CIFAR-10 as In-Distribution.

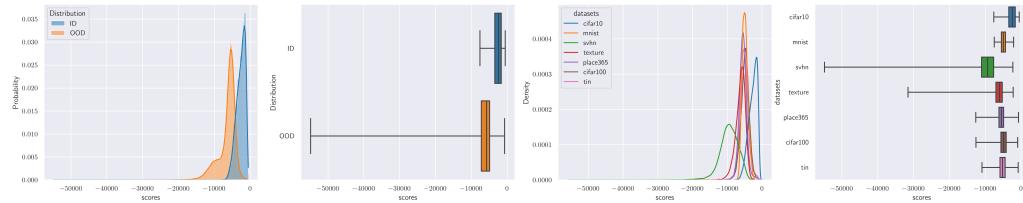


Figure 23. Scores distribution for SSD with CIFAR-10 as In-Distribution.

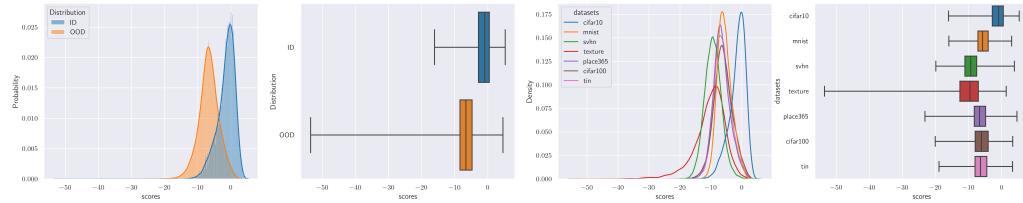


Figure 24. Scores distribution for VIM with CIFAR-10 as In-Distribution.

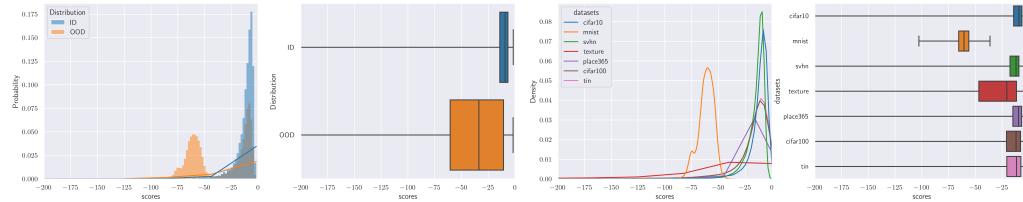


Figure 25. Scores distribution for MDS with CIFAR-10 as In-Distribution.

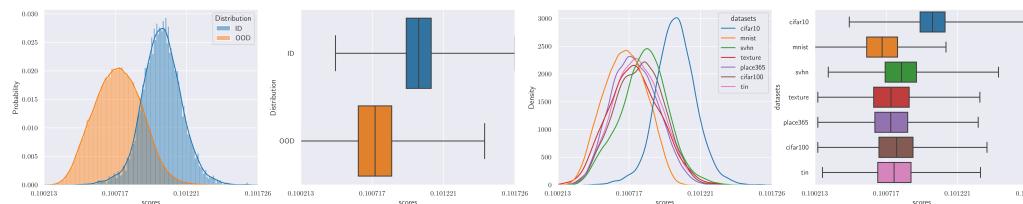


Figure 26. Scores distribution for ODIN with CIFAR-10 as In-Distribution.

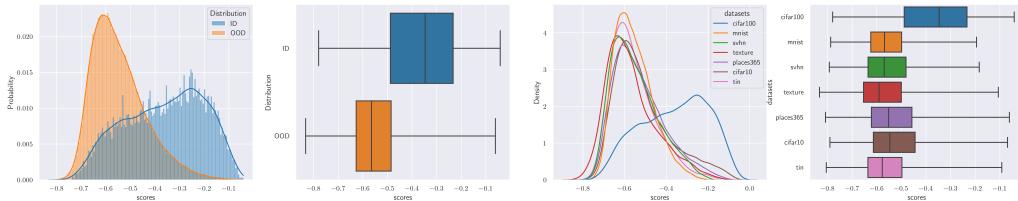


Figure 27. Scores distribution for KNN with cifar-100 as In-Distribution.

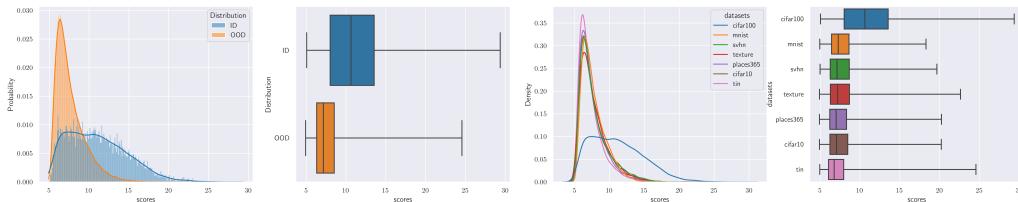


Figure 28. Scores distribution for EBO with cifar-100 as In-Distribution.

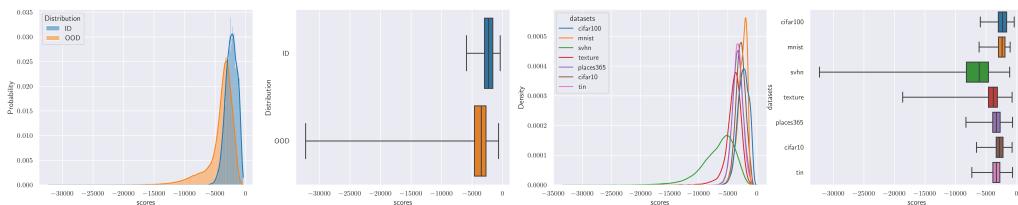


Figure 29. Scores distribution for SSD with cifar-100 as In-Distribution.

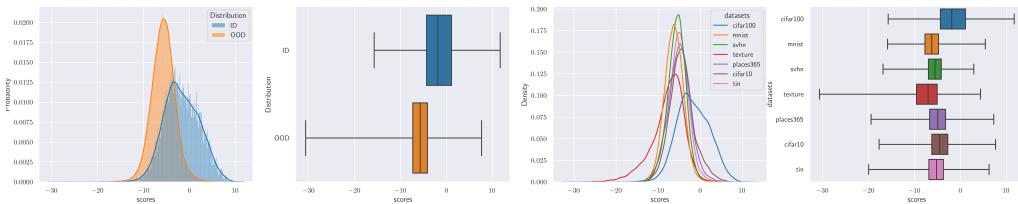


Figure 30. Scores distribution for VIM with cifar-100 as In-Distribution.

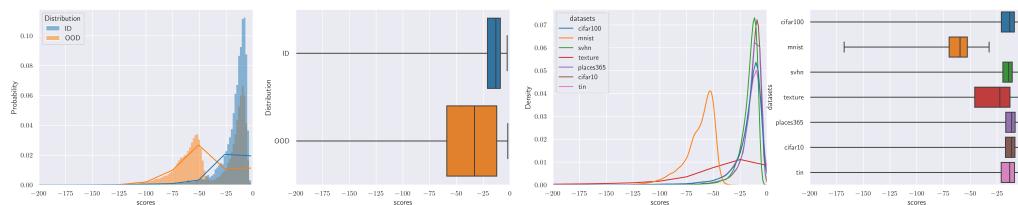


Figure 31. Scores distribution for MDS with cifar-100 as In-Distribution.

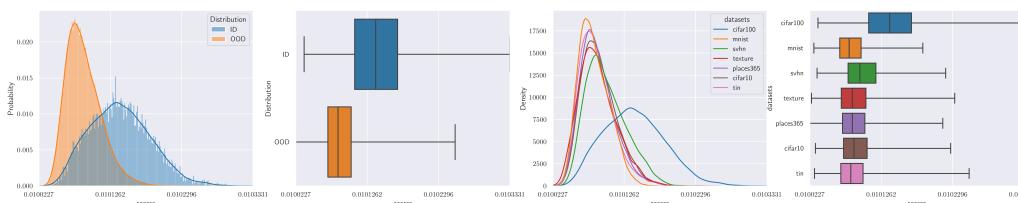


Figure 32. Scores distribution for ODIN with cifar-100 as In-Distribution.