

# A MEDIAN PERSPECTIVE ON UNLABELED DATA FOR OUT-OF-DISTRIBUTION DETECTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Out-of-distribution (OOD) detection plays a crucial role in ensuring the robustness and reliability of machine learning systems deployed in real-world applications. Recent approaches have explored the use of unlabeled data, showing potential for enhancing OOD detection capabilities. However, effectively utilizing unlabeled in-the-wild data remains challenging due to the mixed nature of both in-distribution (InD) and OOD samples. The lack of a distinct set of OOD samples complicates the task of training an optimal OOD classifier. In this work, we introduce Medix, a novel framework designed to identify potential outliers from unlabeled data using the median operation. We use the median because it provides a stable estimate of the central tendency, as an OOD detection mechanism, due to its robustness against noise and outliers. Using these identified outliers, along with labeled InD data, we train a robust OOD classifier. From a theoretical perspective, we derive error bounds that demonstrate Medix achieves a low error rate. Empirical results further substantiate our claims, as Medix outperforms existing methods across the board in open-world settings, confirming the validity of our theoretical insights.

## 1 INTRODUCTION

Deploying machine learning models in real-world applications often exposes them to challenges related to safety and reliability, particularly due to the presence of out-of-distribution (OOD) data. These OOD samples, stemming from unknown categories, should not be predicted by the model. However, neural networks are inherently vulnerable and typically lack the necessary mechanisms to detect and appropriately handle OOD inputs in practice (Nguyen et al., 2015).

Identifying OOD samples during inference is critical yet inherently not easy, as models are not exposed to unknown distributions during training and, therefore, cannot reliably distinguish OOD from in-distribution (InD) data. To address this challenge, recent approaches (Katz-Samuels et al., 2022a; Du et al., 2024a) have explored leveraging additional “in-the-wild” data to improve OOD detection. Specifically, Katz-Samuels et al. (2022a) introduced a method that uses unlabeled wild data for regularizing model training, while still focusing on classifying labeled InD data. The advantage of using such unlabeled wild data lies in its availability—being easily collectible once a model is deployed in its operating environment. This approach allows the model to better capture the true distribution of OOD data encountered during test time, leading to a more robust OOD detection.

However, leveraging unlabeled wild data presents significant challenges due to the complex mixture of InD and OOD data. The absence of a distinct and clean set of OOD samples complicates the development of robust OOD detection methods, especially since the OOD detector model only encounters data drawn from this mixed distribution, without knowledge of whether each sample is from the InD or OOD category. At present, the problem remains underexplored, with substantial opportunities for further advancement. Moreover, few studies establish a formal theoretical foundation, and to the best of our knowledge, Du et al. (2024a) is the only work that provides such a foundation for the “in-the-wild” setting.

Meanwhile, recent studies have demonstrated the effectiveness of median-based approaches in data pruning (Acharya et al., 2024). Prompted by these developments, this paper aims to answer the following question:

How can median-based methods leverage unlabeled wild data to facilitate OOD detection with theoretical guarantees?

In an attempt to provide an affirmative answer to this question, we introduce a novel median-centric perspective for OOD detection. Specifically, we propose a median-based optimization framework and develop an algorithm, Medix, that effectively identifies OOD samples from wild data with low error rate. We provide a theoretical foundation that guarantees minimal error, which is further validated through experiments, demonstrating the robustness and efficiency of our algorithm. We are one of the few studies that provide such a theoretical foundation for the unlabeled “in-the-wild” setting. We show that median-based filtering is robust for outlier detection in unlabeled mixtures by bounding the fraction of InD samples incorrectly flagged as outliers. This fraction is controlled by two effects: the contamination effect, which quantifies the impact of OOD points and remains manageable as long as OOD proportion is below 50%; and the concentration effect, which leverages the sub-Gaussian nature of InD gradients to bound deviations of InD points from their mean gradient. Together, these effects ensure that, with high probability, the number of misclassified InD points remains small, demonstrating the method’s effectiveness even in the worst-case scenario.

We benchmark our approach against two categories of methods: (1) those trained solely on InD data, and (2) those trained with both InD data and an auxiliary unlabeled dataset. On CIFAR-100 (Krizhevsky et al., 2009), Medix demonstrates a significant improvement over the strong baseline KNN+ (Sun et al., 2022), outperforming it by an average of 40.98% in terms of FPR95. Unlike approaches such as Outlier Exposure (Hendrycks et al., 2019), which rely on a clean, auxiliary unlabeled dataset (i.e. they make a strong distributional assumption that the auxiliary data is completely separable from the InD data), Medix achieves superior results without such assumptions, offering greater flexibility. Compared to WOODS (Katz-Samuels et al., 2022a), Medix reduces the average FPR95 by 1.32% on CIFAR-100 and 2.60% on CIFAR-10.

Our key contributions are as follows:

- C1) We propose Medix, a median-centric greedy approach that filters outliers from the unlabeled wild data and then trains an OOD detector on the identified outliers and the InD samples. Our main contribution is the filtering stage.
- C2) We establish theoretical guarantees for the robustness of median-based filtering in identifying both inliers and outliers within unlabeled mixtures. Specifically, we prove that the misclassification rates for both InD samples flagged as outliers and OOD samples retained as inliers are tightly controlled by two effects: the contamination effect, which remains bounded as long as the proportion of OOD samples is below 50%, and the concentration effect, which ensures stability in the gradient behavior of InD samples.
- C3) We conduct an extensive evaluation of Medix across eleven InD-OOD pairs, comparing its performance against 20 competitive baselines. Our results demonstrate that Medix outperforms all the baselines, achieving superior performance across the board. We show via experiments that Medix outlier extraction achieves a low error rate (e.g. only 12.5%; see Figure 2), corroborating our theoretical findings.

## 2 PRELIMINARIES AND PROBLEM SETUP

In this section, we first provide an overview of the OOD detection problem, and then formally define the data setup, model architecture, loss functions, and the learning goal.

**Labeled In-Distribution Data.** Consider the input space  $\mathcal{X}$  and the label space  $\mathcal{Y} = \{1, \dots, K\}$ , which together define the structure of InD data. A labeled dataset  $\mathcal{S}_{\text{in}} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  is generated by sampling  $n$  pairs independently and identically distributed (i.i.d.) from  $\mathbb{P}_{XY}$ , an unknown joint probability distribution over  $\mathcal{X} \times \mathcal{Y}$ . The marginal distribution of  $\mathbb{P}_{XY}$  on  $\mathcal{X}$  is denoted as  $\mathbb{P}_{\text{in}}$ , representing the underlying distribution of InD inputs. We use  $\mathcal{S}_{\text{in}}$  to train an InD model.

**Out-of-distribution detection.** We address a practical scenario where the model is trained using labeled InD data but is later deployed in environments that may contain OOD inputs from classes not represented in the training data, i.e., for some label  $y \notin \mathcal{Y}$ . The model is expected to abstain from making predictions for such OOD inputs. At inference time, the primary objective is to determine whether a given input belongs to the InD distribution or arises from an OOD source.

**Unlabeled wild data.** One of the primary obstacles in OOD detection is the scarcity of labeled OOD samples. The potential sample space of OOD data can be very large, making the collection of labeled examples both costly and impractical. To address this, we introduce unlabeled wild data,  $\mathcal{S}_{\text{wild}} = \{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_m\}$ , into our learning framework to better mimic real-world scenarios as proposed by Katz-Samuels et al. (2022a). Wild data is a blend of InD and OOD samples and can be readily collected during the deployment phase of a pre-trained model on  $\mathcal{S}_{\text{in}}$ . Similar to Du et al. (2024a); Katz-Samuels et al. (2022a), we adopt the Huber contamination model (Huber, 1964) to characterize the marginal distribution of the wild data

$$\mathbb{P}_{\text{wild}} := (1 - \pi)\mathbb{P}_{\text{in}} + \pi\mathbb{P}_{\text{out}}, \quad \pi \in (0, 1], \quad (1)$$

where  $\pi$  denotes the contamination proportion and  $\mathbb{P}_{\text{out}}$  captures the OOD distribution over  $\mathcal{X}$ . We note that the scenario where  $\pi = 0$  corresponds to the absence of OOD samples, rendering the problem trivial.

**Models and Loss Functions.** Let  $f_\phi : \mathcal{X} \rightarrow \mathbb{R}^K$  represent the InD classifier parameterized by  $\phi \in \Phi$ , where  $\Phi$  denotes the parameter space for this classifier. The output of  $f_\phi$  corresponds to a soft probability distribution over the  $K = |\mathcal{Y}|$  InD classes. The loss function for the labeled InD data is defined as  $\ell : \mathbb{R}^K \times \mathcal{Y} \rightarrow \mathbb{R}$ . For OOD detection, we introduce a separate classifier  $g_\theta : \mathcal{X} \rightarrow \mathbb{R}$ , parameterized by  $\theta \in \Theta$ , with  $\Theta$  as the parameter space. The binary loss function associated with  $g_\theta$  is denoted as  $\ell_b(g_\theta(x), y_b)$ , where  $y_b \in \mathcal{Y}_b := \{y_+, y_-\}$ . Here,  $y_+ > 0$  represents the InD class, while  $y_- < 0$  corresponds to the OOD class.

**Learning objective.** Our learning framework is designed to simultaneously train the OOD detector  $g_\theta$  and the multi-class classifier  $f_\phi$ , leveraging both the InD data  $\mathbb{P}_{\text{in}}$  and the wild data  $\mathbb{P}_{\text{wild}}$ . During testing, we evaluate the performance using the following metrics:

$$\begin{aligned} \downarrow \text{FPR}(g_\theta) &= \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{out}}^{\text{test}}} [\mathbb{I}\{g_\theta(\mathbf{x}) = \text{in}\}], \\ \uparrow \text{TPR}(g_\theta) &= \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{\text{in}}} [\mathbb{I}\{g_\theta(\mathbf{x}) = \text{in}\}], \\ \uparrow \text{Acc}(f_\phi) &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{P}_{\mathcal{X}\mathcal{Y}}} [\mathbb{I}\{f_\phi(\mathbf{x}) = y\}], \end{aligned}$$

where  $\mathbb{I}\{\cdot\}$  denotes the indicator function, and  $\mathbb{P}_{\text{out}}^{\text{test}}$  represents the OOD test data distribution.

### 3 METHOD: MEDIAN-CENTRIC FRAMEWORK FOR OOD DETECTION

In this section, we present a novel learning paradigm, termed **Medix**, designed for OOD detection by harnessing the power of unlabeled wild data. Our framework overcomes the limitations of conventional approaches that rely exclusively on InD data and is particularly well-suited for applications in open-world environments, where models are often confronted with previously unseen inputs. The Medix framework is composed of two integral stages: **1) Outlier Extraction:** A filtering process that isolates candidate OOD samples from the unlabeled wild data (explained in Section 3.1), and **2) Detector Training:** train a binary OOD detector using both InD data and the outlier candidates identified in the previous step (explained in Section 3.2). For stage 2, we follow the protocol introduced by Du et al. (2024a). As we will demonstrate in the subsequent sections, this two-step methodology not only facilitates the effective extraction of OOD data from the unlabelled wild data but also establishes a robust foundation for deploying machine learning models in dynamic, open-world scenarios.

#### 3.1 EXTRACTING CANDIDATE OUTLIERS FROM THE WILD DATA

To isolate potential outliers from the wild mixture  $\mathcal{S}_{\text{wild}}$ , our framework leverages an optimization-based approach that exploits the gradients of the model parameters. These gradients are derived from a classification model,  $f_\phi$ , which is trained solely on the InD dataset  $\mathcal{S}_{\text{in}}$ . The detailed methodology for this process is formally outlined below.

**Reference gradient estimation from InD data.** The first step in our proposed framework is to estimate a reference gradient using the InD dataset  $\mathcal{S}_{\text{in}}$ . This is achieved by training a classifier  $f_\phi$  on  $\mathcal{S}_{\text{in}}$  through empirical risk minimization (ERM) as follows

$$\phi_{\mathcal{S}_{\text{in}}} \in \arg \min_{\phi \in \Phi} \mathcal{L}_{\mathcal{S}_{\text{in}}}(f_\phi), \quad \text{where } \mathcal{L}_{\mathcal{S}_{\text{in}}}(f_\phi) = \frac{1}{n} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{S}_{\text{in}}} \ell(f_\phi(\mathbf{x}_i), y_i), \quad (2)$$

where  $\phi_{\mathcal{S}_{\text{in}}}$  denotes the learned parameters. Once the classifier has been trained, we compute the mean gradient  $\bar{\nabla}_{\text{in}}$  as the average of the gradients of the loss function with respect to the model parameters

over the InD data:

$$\bar{\nabla}_{\text{in}} = \frac{1}{n} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{S}_{\text{in}}} \nabla \ell(f_{\phi_{\mathcal{S}_{\text{in}}}}(\mathbf{x}_i), y_i). \quad (3)$$

In our approach,  $\bar{\nabla}_{\text{in}}$  serves as the reference gradient, enabling the quantification of deviations for other data points relative to this reference.

**Motivation.** We hypothesize that increasing the number of OOD samples in the wild dataset,  $\mathcal{S}_{\text{wild}}$ , will lead to a greater deviation from the average InD gradient,  $\bar{\nabla}_{\text{in}}$ . To test this hypothesis, we design an initial experiment using CIFAR-10 (Krizhevsky et al., 2009) as the InD dataset and SVHN (Netzer et al., 2011) as the OOD dataset. Specifically,  $\mathcal{S}_{\text{wild}}$  consists of 10,000 samples drawn from CIFAR-10, ensuring that these samples are disjoint from the training set used to train the model  $\phi_{\mathcal{S}_{\text{in}}}$ , which we leverage to compute  $\bar{\nabla}_{\text{in}}$ . We incrementally add SVHN OOD samples to  $\mathcal{S}_{\text{wild}}$  and track the behavior of the  $L_2$ -norm deviation between  $\bar{\nabla}_{\text{in}}$  and the element-wise median (EWM) of the gradients of the wild dataset as follows:

$$\left\| \bar{\nabla}_{\text{in}} - \text{EWM} \left( \left\{ \nabla \ell \left( f_{\phi_{\mathcal{S}_{\text{in}}}}(\tilde{\mathbf{x}}_i), \hat{y}_{\tilde{\mathbf{x}}_i} \right) \right\}_{i \in \mathcal{S}_{\text{wild}}} \right) \right\|.$$

Our results, depicted in Figure 1, reveal a clear and monotonic increase in the  $L_2$ -norm deviation, supporting our hypothesis. This observation serves as a key motivation for the method we introduce in the subsequent section. Notably, the stopping criterion for our algorithm is derived from this monotonically increasing behavior, where we terminate the algorithm when the  $L_2$ -norm deviation between consecutive iterations drops below a threshold  $\epsilon$ ; we will explain this method in detail in the following section.

**Filtering potential outliers from unlabeled wild data.** Motivated by the results in Figure 1, we formulate the following optimization problem to identify the outlier subset  $\mathcal{S}_{\text{out}}^*$  in  $\mathcal{S}_{\text{wild}}$ :

$$\mathcal{S}_{\text{in}}^* = \arg \min_{\mathcal{S} \subseteq \mathcal{S}_{\text{wild}}} \left\| \bar{\nabla}_{\text{in}} - \text{EWM}(G_{\mathcal{S}}) \right\|, \quad \text{where } G_{\mathcal{S}} = \left\{ \nabla \ell \left( f_{\phi_{\mathcal{S}_{\text{in}}}}(\tilde{\mathbf{x}}_i), \hat{y}_{\tilde{\mathbf{x}}_i} \right) \right\}_{i \in \mathcal{S}}. \quad (4)$$

Here  $\text{EWM}(\cdot)$  denotes the element-wise median function, and  $\hat{y}_{\tilde{\mathbf{x}}_i}$  represents the predicted label for a wild sample  $\tilde{\mathbf{x}}_i$ . For notational simplicity, we define  $G_{\mathcal{S}} = \{ \nabla \ell(\tilde{\mathbf{x}}_i) \}_{i \in \mathcal{S}}$ . The above optimization problem aims to identify a subset  $\mathcal{S}$  in  $\mathcal{S}_{\text{wild}}$  that minimizes the distance between the EWM of the gradients and the average gradient  $\bar{\nabla}_{\text{in}}$ . According to Figure 1, such a subset, denoted by  $\mathcal{S}_{\text{in}}^*$ , may well represent the InD data in  $\mathcal{S}_{\text{wild}}$ , in which case  $\mathcal{S}_{\text{out}}^* = \mathcal{S}_{\text{wild}} \setminus \mathcal{S}_{\text{in}}^*$  may capture the OOD data in  $\mathcal{S}_{\text{wild}}$ .

Solving the optimization problem in equation 4 can be computationally prohibitive, especially as the size of  $\mathcal{S}_{\text{wild}}$  increases. To address this, we propose a greedy approximation based on a leave-one-out approach, as outlined in Algorithm 1. The algorithm implements an iterative procedure for outlier detection from a wild dataset, leveraging deviations in gradient information relative to the InD dataset.

The algorithm begins by computing the EWM of the wild data gradients at each iteration, which serves as a reference for comparing against the average gradient of the InD data  $\bar{\nabla}_{\text{in}}$ . We denote by  $d_t$ , the  $L_2$  distance between the average InD gradient  $\bar{\nabla}_{\text{in}}$  and the EWM gradients of the data left in the wild set, represented by  $\mathcal{S}$ , i.e.,  $d_t = \|\text{EWM}(G_{\mathcal{S}}) - \bar{\nabla}_{\text{in}}\|$ . The algorithm then iteratively identifies samples in  $\mathcal{S}$  that incur the most significant drop in the  $L_2$  distance with  $d_t$  when removed from  $\mathcal{S}$  as OOD. Specifically, having data samples  $\mathcal{S}$  left, we remove each sample  $i \in \mathcal{S}$  and compute the EWM as  $\text{EWM}(G_{\mathcal{S} \setminus \{i\}})$  and the distance to  $\bar{\nabla}_{\text{in}}$  as  $\|\text{EWM}(G_{\mathcal{S} \setminus \{i\}}) - \bar{\nabla}_{\text{in}}\|$ . We then find the drop in distance  $\delta_i = d_t - \|\text{EWM}(G_{\mathcal{S} \setminus \{i\}}) - \bar{\nabla}_{\text{in}}\|$  and find the  $k$  samples with the largest drop and identify them as OOD. The algorithm repeats until there is no significant drop in  $\delta_i$  or a maximum number of iterations is reached. The convergence criterion is based on the change in the  $L_2$  distance between two iterations, which must fall below a predefined  $\epsilon$  threshold; this criterion is inspired by the monotonically increasing trend that we observed in our preliminary experiment in Figure 1, i.e.,

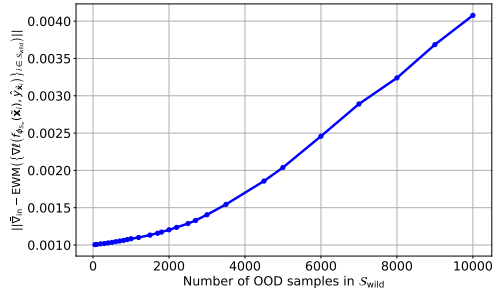


Figure 1: Distance deviation as we increase OOD samples in  $\mathcal{S}_{\text{wild}}$ .

as a substantial number of OOD samples are removed, the  $L_2$  distance gradually decreases to a small value, signaling the point at which the algorithm should halt to avoid erroneously identifying InD samples as outliers, thus preventing any degradation in performance.

---

**Algorithm 1** Iterative Outlier Detection via **Medix**


---

**Require:**  $\bar{\nabla}_{\text{in}}, G_i = \nabla \ell(f_{\phi_{\mathcal{S}_{\text{in}}}(\tilde{\mathbf{x}}_i), \hat{y}_{\tilde{\mathbf{x}}_i})}$ , maximum iterations  $T$ , hyperparameters  $\epsilon, k$

- 1: Initialize  $\mathcal{S} \leftarrow \mathcal{S}_{\text{wild}}$  (wild set),  $\mathcal{O} \leftarrow \emptyset$  (outliers),  $t \leftarrow 0$  (iteration),  $d_t \leftarrow 0$ ,  $\delta_{\text{max}} \leftarrow \infty$  (deviations),  $\mathcal{I}_k \leftarrow \emptyset$
- 2: **while**  $t \leq T$  or  $|\delta_{\text{max}}| > \epsilon$  **do**
- 3:      $\mathcal{O} \leftarrow \mathcal{O} \cup \{\tilde{\mathbf{x}}_i : i \in \mathcal{I}_k\}$  ▷ Add outliers to set  $\mathcal{O}$
- 4:      $d_t \leftarrow \|\text{EWM}(G_{\mathcal{S}}) - \bar{\nabla}_{\text{in}}\|$  ▷ Compute L2 deviation
- 5:     **for each**  $i \in \mathcal{S}$  **do**
- 6:          $\delta_i \leftarrow d_t - \|\text{EWM}(G_{\mathcal{S} \setminus \{i\}}) - \bar{\nabla}_{\text{in}}\|$
- 7:     **end for**
- 8:      $\mathcal{I}_k \leftarrow \text{indices}(\text{top-}k(\{\delta_i\}_{i \in \mathcal{S}}))$  ▷ Select top- $k$  indices
- 9:      $\mathcal{S} \leftarrow \mathcal{S} \setminus \mathcal{I}_k$  ▷ Remove outliers from  $\mathcal{S}$
- 10:      $\delta_{\text{max}} \leftarrow \max_{i \in \mathcal{S}} \{\delta_i\}$ ,  $t \leftarrow t + 1$
- 11: **end while**
- 12: **return**  $\mathcal{S}_{\text{out}} = \mathcal{O}$  ▷ Return detected outliers

---

### 3.2 TRAINING THE OOD DETECTOR WITH CANDIDATE OUTLIERS

After identifying the candidate outlier set  $\hat{\mathcal{S}}_{\text{out}}$  from the wild data, we proceed to train the OOD detector  $g_\theta$  designed to maximize the distinction between InD and candidate outlier data following the protocol in Du et al. (2024a). The objective function explicitly enforces separability at the decision boundary (thresholded at 0), assigning positive outputs for labeled InD samples  $\mathbf{x} \in \mathcal{S}_{\text{in}}$  and negative outputs for candidate outliers  $\tilde{\mathbf{x}} \in \hat{\mathcal{S}}_{\text{out}}$ . Specifically, the loss function is defined as:

$$\begin{aligned} \mathcal{L}_{\mathcal{S}_{\text{in}}, \hat{\mathcal{S}}_{\text{out}}}(g_\theta) &= \mathcal{L}_{\mathcal{S}_{\text{in}}}^+(g_\theta) + \mathcal{L}_{\hat{\mathcal{S}}_{\text{out}}}^-(g_\theta), \\ \mathcal{L}_{\mathcal{S}_{\text{in}}}^+(g_\theta) &= \mathbb{E}_{\mathbf{x} \in \mathcal{S}_{\text{in}}} \mathbb{I}\{g_\theta(\mathbf{x}) \leq 0\}, \\ \mathcal{L}_{\hat{\mathcal{S}}_{\text{out}}}^-(g_\theta) &= \mathbb{E}_{\tilde{\mathbf{x}} \in \hat{\mathcal{S}}_{\text{out}}} \mathbb{I}\{g_\theta(\tilde{\mathbf{x}}) > 0\}. \end{aligned} \quad (5)$$

To address the challenge posed by the non-differentiable nature of the 0/1 loss, a binary loss based on a differentiable sigmoid function is employed as a smooth surrogate, ensuring tractable optimization while retaining alignment with the original objective. The learning framework incorporates the InD risk as defined in Equation 2, thereby safeguarding the predictive accuracy for InD samples. This unified optimization approach significantly bolsters the generalization performance of  $g_\theta$ , enabling robust detection of OOD samples drawn from  $\mathbb{P}_{\text{out}}$ , as corroborated by our results.

## 4 THEORETICAL ANALYSIS

We now present the theoretical guarantees of Medix’s filtering stage. The following theorems provide provable upper bounds on the misclassification rates for both InD and OOD points. Together, they demonstrate the two-sided robustness of our EWM filtering. For detailed proofs, see Appendix C.

**Theorem 4.1** (Inlier Misclassification Bound). *Assume that the gradients of InD points in  $\mathcal{S}_{\text{wild}}$  are i.i.d., and each coordinate is sub-Gaussian with variance proxy  $\sigma^2$ . Let  $\epsilon = \sigma \sqrt{2 \log(2dm_{\text{in}})}$ , and fix any confidence level  $\delta \in (0, 1)$ . Then, with probability at least  $1 - \delta$ , the inlier misclassification rate of the EWM filtering rule satisfies:*

$$\text{ERR}_{\text{in}} \leq \underbrace{\frac{1}{m_{\text{in}}} + 2\sqrt{\frac{\log(1/\delta)}{2m_{\text{in}}}}}_{\text{Concentration term}} + \underbrace{\frac{\pi}{2(1-\pi)}}_{\text{Contamination term}}.$$

Theorem 4.1 reveals that the inlier misclassification rate is governed by two key effects: a *concentration term*, which vanishes with larger sample size and captures statistical fluctuations in the InD gradients, and a *contamination term*, which quantifies the influence of OOD points on the EWM. Notably, the bound remains controlled as long as the contamination ratio  $\pi < 0.5$ , underscoring the robustness of median-based filtering to moderate OOD presence. We now turn to the complementary question: *how many true outliers fail to be flagged and are mistakenly retained?* The next theorem answers this by bounding the OOD misclassification rate.

**Theorem 4.2** (Outlier Misclassification Bound). Assume that the gradients of OOD points in  $\mathcal{S}_{\text{wild}}$  are i.i.d., and each coordinate is sub-Gaussian with variance proxy  $\sigma_{\text{out}}^2$ . Suppose the mean OOD gradient  $\mu_{\text{out}}$  satisfies a separation condition:

$$\|\mu_{\text{out}} - \bar{\nabla}_{\text{in}}\|_2 \geq \Delta\sqrt{d},$$

for some  $\Delta > 0$ . Then, for any tolerance  $\epsilon \in (0, \Delta)$  and confidence level  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ , the outlier misclassification rate satisfies:

$$\text{ERR}_{\text{out}} \leq \underbrace{2d \exp\left(-\frac{(\Delta - \epsilon)^2}{2\sigma_{\text{out}}^2}\right)}_{\text{Separation term}} + \underbrace{\sqrt{\frac{\log(1/\delta)}{2m_{\text{out}}}}}_{\text{Concentration term}} + \underbrace{\frac{1 - \pi}{2\pi}}_{\text{Contamination term}}.$$

Theorem 4.2 complements the previous result by bounding the fraction of OOD points mistakenly retained as InD. Together, Theorems 4.1 and 4.2 provide a two-sided guarantee for the robustness of the Medix filtering method. The inlier and outlier misclassification rates are governed by a balance between the following three effects:

**Contamination Effect.** For Theorem 4.1, even when the wild dataset contains a significant fraction  $\pi$  of OOD points, the median-based criterion remains stable as long as  $\pi < 0.5$ . This is reflected in the  $\pi/[2(1 - \pi)]$  term of the bound, which increases with  $\pi$  and encodes the risk that OOD gradients may skew the EWM. Conversely, Theorem 4.2 includes a symmetric penalty  $\frac{1 - \pi}{2\pi}$ , quantifying the difficulty of isolating OOD points when they are underrepresented.

**Concentration Effect.** Both bounds exploit the sub-Gaussian nature of the gradient coordinates. For InD data, this ensures that most gradients concentrate near  $\bar{\nabla}_{\text{in}}$ , keeping the median stable even in finite samples. For OOD points, concentration around  $\mu_{\text{out}}$  allows us to quantify the risk that a misaligned OOD sample slips past the filter. These concentration effects decay as  $1/\sqrt{m}$ , where  $m$  is the number of samples from each class.

**Separation Effect.** Unique to Theorem 4.2, this effect quantifies how far the OOD mean gradient must lie from the InD mean gradient in order to reliably reject OOD samples. The exponential term  $\exp(-(\Delta - \epsilon)^2/2\sigma_{\text{out}}^2)$  captures this trade-off: the more separated the distributions, the less likely it is for OOD gradients to fool the filter.

Theorems 4.1 and 4.2 jointly establish that, with high probability, median filtering achieves robust separation of InD and OOD samples in unlabeled data mixtures. The fraction of InD samples misclassified as outliers is bounded by contamination and concentration effects, while the fraction of OOD points incorrectly retained is governed by an exponential separation term, a concentration bound, and a reverse contamination effect. These results provide rigorous theoretical assurance that Medix minimizes both types of errors under mild assumptions.

*Remark 4.3. Sub-Gaussian Assumption for Gradient Coordinates.* The assumption that gradients of InD samples in each coordinate are sub-Gaussian is crucial for deriving concentration bounds and ensuring robustness in statistical estimation. Further empirical evidence as shown in Figure 4a confirms this assumption, where we observe that the histogram of gradient values for InD samples is bell-shaped and concentrated around the mean. This indicates that the gradients exhibit light tails, a defining characteristic of sub-Gaussian random variables, in which variables have exponentially decaying tails, meaning extreme gradient values are rare. This aligns with the observed behavior in the histogram. Furthermore, Figure 4b compares the empirical quantiles of the gradient distribution with the theoretical quantiles of a Gaussian distribution in a Q-Q plot (Wilk & Gnanadesikan, 1968). The close alignment of points with the 45-degree reference line demonstrates that the empirical distribution of gradients indeed closely resembles a Gaussian distribution. Since Gaussian random variables are a subset of sub-Gaussian variables, this supports the sub-Gaussian assumption. We also provide a looser version of these bounds, which **removes the sub-Gaussian assumption** and holds under merely bounded second moments. This version is presented in Theorem C.3 (Appendix C.3). While the rates degrade, the core robustness guarantee of Medix still holds in this broader setting.

## 5 EXPERIMENTS

This section will demonstrate the efficacy of Medix across various InD-OOD dataset pairs, benchmarking it against 20 widely-used baselines. All experiments are performed on hardware equipped with NVIDIA A100-SXM4-80GB GPUs. We provide the necessary code to reproduce our results.

## 5.1 MODELS, DATASETS, AND BASELINES

**Datasets.** We use the same experimental protocol as Katz-Samuels et al. (2022a), which introduced the problem of learning OOD detectors with wild data, enabling a fair comparison to prior work. Specifically, we use CIFAR-10 and CIFAR-100 as InD datasets ( $\mathbb{P}_{\text{in}}$ ). For OOD testing, we select a suite of natural image datasets, including PLACES365 (Zhou et al., 2017), SVHN, TEXTURES (Cimpoi et al., 2014), and LSUN-RESIZE & LSUN-C (Yu et al., 2015) as OOD datasets ( $\mathbb{P}_{\text{out}}$ ). To simulate wild data ( $\mathbb{P}_{\text{wild}}$ ), we combine a subset of InD data ( $\mathbb{P}_{\text{in}}$ ) with the OOD data ( $\mathbb{P}_{\text{out}}$ ) under a default mixing parameter  $\pi = 0.5$ . For example, when using PLACES365 as an OOD test set, we construct a wild mixture by combining CIFAR with PLACES365 as wild data and test on PLACES365 as the OOD set. This procedure is repeated across all OOD datasets and baselines. The InD CIFAR dataset is split into two halves: the first 25,000 images to train  $\phi_{\mathcal{S}_{\text{in}}}$ , while the remaining images to generate the wild mixture  $\mathcal{S}_{\text{wild}}$ . For gradient computation, we use the penultimate layer weights, as these have been shown to be particularly informative for OOD detection (Huang et al., 2021a).

**Evaluation metrics.** We evaluate using three standard metrics: (1) False Positive Rate (FPR95 $\downarrow$ ) of OOD samples when the True Positive Rate of InD samples is 95%, (2) Area Under the Receiver Operating Characteristic Curve (AUROC $\uparrow$ ), and (3) InD Classification Accuracy (InD Acc $\uparrow$ ).

**Baselines.** Our comparison encompasses a diverse set of competitive baselines, categorized based on whether they are trained using only InD data or both InD and wild data. For the methods trained exclusively on InD data ( $\mathbb{P}_{\text{in}}$ ), we compare Medix against a variety of established OOD detection methods, including Maximum Softmax Probability (MSP) (Hendrycks & Gimpel, 2016), ODIN (Liang et al., 2018), Mahalanobis Distance (Lee et al., 2018b), Energy Score (Liu et al., 2020c), ReAct (Sun et al., 2021), DICE (Sun & Li, 2022), KNN Distance (Sun et al., 2022), and ASH (Djurisic et al., 2022); these methods use a model trained with softmax cross-entropy loss. Additionally, we also compare against methods based on contrastive loss, such as CSI (Tack et al., 2020b) and KNN+ (Sun et al., 2022), for a more comprehensive comparison. For methods that leverage both InD and wild data, we compare against Outlier Exposure (OE) (Hendrycks et al., 2019) and energy-regularization learning (Liu et al., 2020c), which regularize the training by promoting lower confidence or higher energy on outlier data. We also include a comparison with WOODS (Katz-Samuels et al., 2022a), which introduced the concept of wild unlabeled data and utilizes it for OOD detection through a constrained optimization approach. Finally, we included more recent baselines, including CONJ (Peng et al., 2024) and DRL (Zhang et al., 2024), to provide a more thorough evaluation.

## 5.2 EXPERIMENTAL SETUP

In line with WOODS (Katz-Samuels et al., 2022a), we employ a Wide ResNet architecture (Zagoruyko, 2016) with 40 layers and a width factor of 2 for the InD classifier  $\phi_{\mathcal{S}_{\text{in}}}$ . It is trained using stochastic gradient descent with a momentum of 0.9, weight decay of 0.0005, and an initial learning rate of 0.1. Training is performed for 100 epochs with cosine learning rate decay, a batch size of 128, and a dropout rate of 0.3. Hyperparameters  $\epsilon$  and  $k$  used in the proposed method Medix are selected from the sets  $\{5e-5, 5e-4, 5e-3, 5e-2\}$  and  $\{4k, 7k, 10k, 20k\}$ , respectively, taking into account dataset sizes and with the objective of maximizing OOD performance. For the OOD classifier  $g_{\theta}$ , we initialize it with the pre-trained InD classifier  $\phi_{\mathcal{S}_{\text{in}}}$  and add a linear layer that performs binary classification using the penultimate-layer features. The learning rate for this classifier is set to 0.001, and fine-tuning is done for 100 epochs as outlined in Equation equation 5. We combine the binary classification loss with the InD classification loss, assigning a weight of 10 to the binary classification component. All other training parameters remain the same as those used for training  $\phi_{\mathcal{S}_{\text{in}}}$ .

## 5.3 RESULTS

We present our main results in Table 2 on CIFAR-100, where Medix remarkably outperforms all OOD detection baselines. The results highlight the following key observations: (1) Methods trained on both InD and wild data significantly outperform those trained exclusively on InD data. Medix reduces the FPR95 by 52.31% on PLACES365 and 38.24% on TEXTURES compared to KNN+, demonstrating the effectiveness of incorporating in-the-wild data for model regularization. (2) Medix further outperforms competitive methods utilizing wild data ( $\mathbb{P}_{\text{wild}}$ ): On CIFAR-100, Medix achieves an average FPR95 of 5.42%, which represents a 1.32% improvement over WOODS. Additionally, Medix maintains a competitive InD accuracy of 73.33%. This slight difference can be attributed to the fact that our method is trained on 25,000 labeled InD samples, while baseline methods, which do not leverage wild data, use the full CIFAR-100 training set of 50,000 samples. We present the results for CIFAR-10 in Table 1, where Medix surpasses all baseline methods. Medix outperforms

Table 1: OOD detection performance comparison of Medix and baselines on CIFAR-10 as InD data. Performance averaged over five runs; best results are highlighted in **bold**.

Methods	OOD Datasets												InD ACC $\uparrow$
	SVHN		PLACES365		LSUN-C		LSUN-RESIZE		TEXTURES		Average		
	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$	
	Using $\mathbb{P}_m$ only												
MSP	48.49	91.89	59.48	88.20	30.80	95.65	52.15	91.37	59.28	88.50	50.04	91.12	94.84
ODIN	33.35	91.96	57.40	84.49	15.52	97.04	26.62	94.57	49.12	84.97	36.40	90.61	94.84
Mahalanobis	12.89	97.62	68.57	84.61	39.22	94.15	42.62	93.23	15.00	97.33	35.66	93.34	94.84
Energy	35.59	90.96	40.14	89.89	8.26	98.35	27.58	94.24	52.79	85.22	32.87	91.73	94.84
KNN	24.53	95.96	25.29	95.69	25.55	95.26	27.57	94.71	50.90	89.14	30.77	94.15	94.84
ReAct	40.76	89.57	41.44	90.44	14.38	97.21	33.63	93.58	53.63	86.59	36.77	91.48	94.84
DICE	35.44	89.65	46.83	86.69	6.32	98.68	28.93	93.56	53.62	82.20	34.23	90.16	94.84
ASH	6.51	98.65	48.45	88.34	0.90	99.73	4.96	98.92	24.34	95.09	17.03	96.15	94.84
CSI	17.30	97.40	34.95	93.64	1.95	99.55	12.15	98.01	20.45	95.93	17.36	96.91	94.17
KNN+	2.99	99.41	24.69	94.84	2.95	99.39	11.22	97.98	9.65	98.37	10.30	97.99	93.19
	Using $\mathbb{P}_m$ and $\mathbb{P}_{\text{out}}$												
OE	1.13	99.53	19.48	94.88	1.91	98.16	0.54	98.84	7.75	98.56	6.16	97.99	94.12
Energy (w/ OE)	5.24	98.72	14.66	96.18	2.35	99.30	4.85	98.62	10.51	97.10	7.52	97.98	94.24
WOODS	0.17	99.91	10.19	98.05	0.31	99.14	0.11	99.38	6.21	98.13	3.40	98.92	94.74
Medix	<b>0.06</b>	<b>99.98</b>	<b>2.98</b>	<b>99.10</b>	<b>0.01</b>	<b>99.98</b>	<b>0.01</b>	<b>99.98</b>	<b>0.96</b>	<b>99.66</b>	<b>0.80</b>	<b>99.74</b>	<b>93.58</b>
(Ours)	$\pm 0.01$	$\pm 0.01$	$\pm 0.29$	$\pm 0.09$	$\pm 0.01$	$\pm 0.00$	$\pm 0.01$	$\pm 0.01$	$\pm 0.13$	$\pm 0.06$	$\pm 0.09$	$\pm 0.03$	$\pm 0.64$

Table 2: OOD detection performance comparison of Medix and baselines on CIFAR-100 as InD data. Performance averaged over five runs; best results are highlighted in **bold**.

Methods	OOD Datasets												InD ACC $\uparrow$
	SVHN		PLACES365		LSUN-C		LSUN-RESIZE		TEXTURES		Average		
	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$	FPR95 $\downarrow$	AUROC $\uparrow$	
	Using $\mathbb{P}_m$ only												
MSP	84.59	71.44	82.84	73.78	66.54	83.79	82.42	75.38	83.29	73.34	79.94	75.55	75.96
ODIN	84.66	67.26	87.88	71.63	55.55	87.73	71.96	81.82	79.27	73.45	75.86	76.38	75.96
Mahalanobis	57.52	86.01	88.83	67.87	91.18	69.69	21.23	96.00	39.39	90.57	59.63	82.03	75.96
Energy	85.82	73.99	80.56	75.44	35.32	93.53	79.47	79.23	79.41	76.28	72.12	79.69	75.96
KNN	66.38	83.76	79.17	71.91	70.96	83.71	77.83	78.85	88.00	67.19	76.47	77.08	75.96
ReAct	74.33	88.04	81.33	74.32	39.30	91.19	79.86	73.69	67.38	82.80	68.44	82.01	75.96
DICE	88.35	72.58	81.61	75.07	26.77	94.74	80.21	78.50	76.29	76.07	70.65	79.39	75.96
ASH	21.36	94.28	68.37	71.22	15.27	95.65	68.18	85.42	40.87	92.29	42.81	87.77	75.96
CSI	64.70	84.97	82.25	73.63	38.10	92.52	91.55	63.42	74.70	92.66	70.26	81.44	69.90
KNN+	32.21	93.74	68.30	75.31	40.37	86.13	44.86	88.88	46.26	87.40	46.40	86.29	73.78
	Using $\mathbb{P}_m$ and $\mathbb{P}_{\text{out}}$												
OE	2.86	99.05	40.21	88.75	4.13	99.05	1.25	99.38	22.86	94.63	14.26	96.17	73.38
Energy (w/ OE)	2.71	99.34	34.82	90.05	3.27	99.18	2.54	99.23	30.16	94.76	14.70	96.51	72.76
WOODS	0.17	99.80	21.87	93.73	0.48	99.61	1.24	99.54	9.95	95.97	6.74	97.73	73.91
Medix	<b>0.16</b>	<b>99.96</b>	<b>15.99</b>	<b>95.23</b>	<b>0.13</b>	<b>99.98</b>	<b>0.83</b>	<b>99.83</b>	<b>8.02</b>	<b>97.79</b>	<b>5.42</b>	<b>98.96</b>	<b>73.33</b>
(Ours)	$\pm 0.02$	$\pm 0.00$	$\pm 0.66$	$\pm 0.14$	$\pm 0.06$	$\pm 0.02$	$\pm 0.36$	$\pm 0.06$	$\pm 0.75$	$\pm 0.30$	$\pm 0.37$	$\pm 0.10$	$\pm 0.83$

WOODS by 7.21% on PLACES365 and 5.25% on TEXTURES in terms of FPR95, demonstrating its effectiveness in detecting OOD samples.

**A representative visual example of Medix.** We further investigate the performance of Algorithm 1 (Outlier Extraction) in extracting OOD samples from wild data  $\mathcal{S}_{\text{wild}}$ . To visualize this, we design an experiment using 2-dimensional synthetic data. This simulation is designed to be simple to facilitate better understanding. We generate the InD data by sampling from three multivariate Gaussian distributions, corresponding to three classes. The mean vectors are set to  $[-2, 0]$ ,  $[2, 0]$ , and  $[0, 2\sqrt{3}]$ , respectively. The covariance matrix for all three classes is fixed at  $0.25 \cdot I$ . For the OOD data, we use a multivariate Gaussian distribution  $\mathcal{N}([20, 2\sqrt{3}], 0.25 \cdot I)$ . In Figure 2, we observe that our method successfully identifies outliers, the majority of which align with the ground truth. Medix successfully flags 87.5% of actual OOD samples as outliers, underscoring its robustness in outlier extraction.

**Additional studies and insights.** Due to space constraints, we defer additional experiments and insights to Appendix A, which include (1) ablation studies on hyperparameter selection and sensitivity analysis (Appendix A.2), showing Medix’s strong robustness to hyperparameters, (2) a comparison between EWM and geometric median (Appendix A.1), showing that EWM is more sensitive to distributional shifts, making it more effective and reliable choice for filtering in our method, (3) a comprehensive comparison with methods employing competitive contrasting learning objectives (Appendix A.3), demonstrating that Medix outperforms even these competitive baselines by a significant margin, (4) evaluation on large-scale data under the complex unseen OOD setting, where  $\mathbb{P}_{\text{out}}^{\text{test}} \neq \mathbb{P}_{\text{out}}$  (Appendix A.4), demonstrating that Medix outperforms baselines by a significant margin, (5) evaluating computation and memory efficiency of Medix (Appendix A.6) (6) evaluating the impact of pseudo-label quality (Appendix A.5), showing that our method is resilient to noisy or low-confidence labels, and (4) a comparison with semi-supervised open-set recognition methods (Appendix A.7).

## 6 RELATED WORK

In recent years, there has been a growing interest in OOD detection (Fort et al., 2021; Yang et al., 2024; Fang et al., 2022; Zhu et al., 2022; Yang et al., 2022; Wang et al., 2022c; Galil et al., 2023; Djuricic et al., 2023; Tao et al., 2023; Zheng et al., 2023; Wang et al., 2022b; 2023b; Uppaal et al., 2023; Zhu et al., 2023; Bai et al., 2023; Ming & Li, 2024; Zhang et al., 2023; Ghosal et al., 2024). One

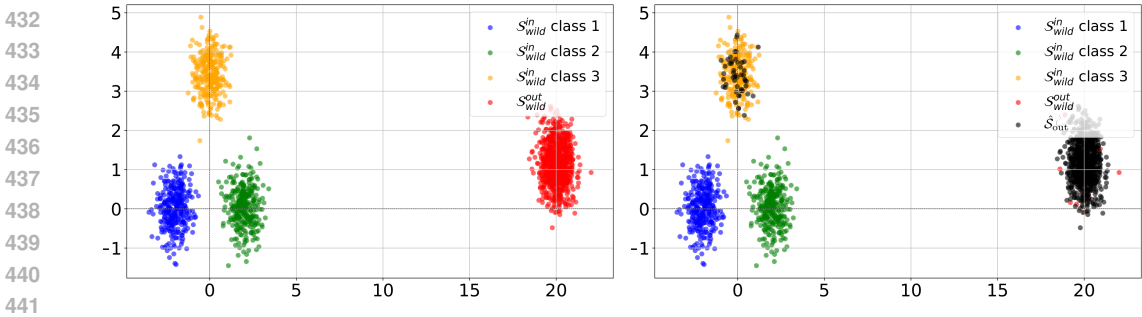


Figure 2: Example of Medix applied to unlabeled wild data. (a) Setup of the InD data  $\mathcal{S}_{\text{wild}}^{\text{in}}$  and OOD data  $\mathcal{S}_{\text{wild}}^{\text{out}}$  in the wild, with inliers sampled from three multivariate Gaussian distributions. (b) Outliers  $\hat{\mathcal{S}}_{\text{out}}$  filtered by Medix (in black), with an error rate of  $\hat{\mathcal{S}}_{\text{out}}$  containing InD data  $\mathcal{S}_{\text{wild}}^{\text{in}}$  is only 12.5%.

approach to detect OOD data uses scoring functions to assess data distribution, including distance-based methods (Lee et al., 2018a; Tack et al., 2020a; Ren et al., 2021; Sehwan et al., 2021; Sun et al., 2022; Du et al., 2022a; Ming et al., 2023; Ren et al., 2022), gradient-based score (Huang et al., 2021b), energy-based score (Liu et al., 2020b; Wang et al., 2021; Wu et al., 2023), confidence-based approaches (Bendale & Boult, 2016; Hendrycks & Gimpel, 2017; Liang et al., 2018), and Bayesian methods (Gal & Ghahramani, 2016; Lakshminarayanan et al., 2017; Maddox et al., 2019; Malinin & Gales, 2019; Wen et al., 2020; Kristiadi et al., 2020).

Another approach to OOD detection involves using regularization techniques during the training phase (Malinin & Gales, 2018; Geifman & El-Yaniv, 2019; Hein et al., 2019; Meinke & Hein, 2020; Jeong & Kim, 2020; Liu et al., 2020a; Van Amersfoort et al., 2020; Yang et al., 2021; Wei et al., 2022; Du et al., 2022b; 2023; Wang et al., 2023a). For example, regularization techniques can be applied to the model to either reduce its confidence (Lee et al., 2017; Hendrycks et al., 2019) or increase its energy (Liu et al., 2020b; Du et al., 2022c; Ming et al., 2022) on the OOD data. Most of these regularization methods assume the availability of an auxiliary OOD dataset.

Several studies (Zhou et al., 2021; Katz-Samuels et al., 2022b; He et al., 2023) have relaxed the assumption of using only labeled data by incorporating unlabeled wild data (Katz-Samuels et al., 2022a; Geng et al., 2025), though they did not propose a clear mechanism for outlier detection. In contrast, Du et al. (2024a;b) introduced an explicit outlier filtering method, but their thresholding technique differs fundamentally from ours, as we utilize a new median-centric approach to detect the outliers. Additionally, Katz-Samuels et al. (2022a); Du et al. (2024a) operate under the assumption of batch-level mixing, where each batch has a set ratio of InD and OOD samples. However, with large outsourced datasets, such structured mixing is not available—data is mixed randomly across the dataset. Our method addresses this by enabling dataset-level mixing without relying on batch-level structure. Many studies also leverage positive-unlabeled learning, which trains classifiers using positive and/or unlabeled data (Letouzey et al., 2000; Hsieh et al., 2015; Plessis et al., 2015; Niu et al., 2016; Gong et al., 2018; Chapel et al., 2020; Garg et al., 2021; Xu & Denil, 2021; Garg et al., 2022; Zhao et al., 2022; Acharya et al., 2022). However, a key distinction from our approach is that these methods focus solely on differentiating  $\mathbb{P}_{\text{out}}$  from  $\mathbb{P}_{\text{in}}$ , without simultaneously training an OOD classifier. Additionally, we propose a median-centric method to identify outliers in unlabeled data, with provably low error rates.

## 7 CONCLUSIONS

In this work, we introduced Medix, a novel median-centric framework for OOD detection that leverages unlabeled in-the-wild data. Using the inherent robustness of median operation, Medix effectively filters outliers from mixed unlabeled data, enabling the training of a reliable OOD detector. Our theoretical analysis established provable bounds on the inlier misclassification rate, demonstrating that Medix maintains robustness even under significant OOD contamination (up to 50%), with errors controlled by sub-Gaussian concentration and contamination effects. We also provided complementary theoretical limits on the rate of OOD misclassification to accurately isolate OOD samples under clear separation conditions. Empirical validation across diverse benchmarks showcased Medix’s performance superiority over 20 baselines: it reduced the average FPR95 by 40.98% compared to strong baselines like KNN+ and outperformed state-of-the-art methods such as WOODS and DRL, while achieving an outlier extraction error rate as low as 12.5%.

## REFERENCES

- 486  
487  
488 Anish Acharya, Sujay Sanghavi, Li Jing, Bhargav Bhushanam, Dhruv Choudhary, Michael Rabbat,  
489 and Inderjit Dhillon. Positive unlabeled contrastive learning. *arXiv preprint arXiv:2206.01206*,  
490 2022.
- 491 Anish Acharya, Inderjit S Dhillon, and Sujay Sanghavi. Geometric median (gm) matching for robust  
492 data pruning. *arXiv preprint arXiv:2406.17188*, 2024.
- 493 Haoyue Bai, Gregory Canal, Xuefeng Du, Jeongyeol Kwon, Robert D Nowak, and Yixuan Li. Feed  
494 two birds with one scone: Exploiting wild data for both out-of-distribution generalization and  
495 detection. In *International Conference on Machine Learning*, pp. 1454–1471. PMLR, 2023.
- 496  
497 Abhijit Bendale and Terrance E Boult. Towards open set deep networks. In *Proceedings of the IEEE*  
498 *conference on computer vision and pattern recognition*, pp. 1563–1572, 2016.
- 499 Laetitia Chapel, Mokhtar Z. Alaya, and Gilles Gasso. Partial optimal transport with applications on  
500 positive-unlabeled learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin  
501 (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 2903–2913. Curran  
502 Associates, Inc., 2020.
- 503 Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describ-  
504 ing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern*  
505 *recognition*, pp. 3606–3613, 2014.
- 506  
507 Andrija Djurisic, Nebojsa Bozanic, Arjun Ashok, and Rosanne Liu. Extremely simple activation  
508 shaping for out-of-distribution detection. *arXiv preprint arXiv:2209.09858*, 2022.
- 509 Andrija Djurisic, Nebojsa Bozanic, Arjun Ashok, and Rosanne Liu. Extremely simple activation  
510 shaping for out-of-distribution detection. In *The Eleventh International Conference on Learning*  
511 *Representations*, 2023.
- 512  
513 Xuefeng Du, Gabriel Gozum, Yifei Ming, and Yixuan Li. Siren: Shaping representations for detecting  
514 out-of-distribution objects. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and  
515 A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 20434–20449.  
516 Curran Associates, Inc., 2022a.
- 517 Xuefeng Du, Xin Wang, Gabriel Gozum, and Yixuan Li. Unknown-aware object detection: Learning  
518 what you don’t know from videos in the wild. In *Proceedings of the IEEE/CVF Conference on*  
519 *Computer Vision and Pattern Recognition*, pp. 13678–13688, 2022b.
- 520 Xuefeng Du, Zhaoning Wang, Mu Cai, and Sharon Li. Towards unknown-aware learning with virtual  
521 outlier synthesis. In *International Conference on Learning Representations*, 2022c.
- 522  
523 Xuefeng Du, Yiyao Sun, Jerry Zhu, and Yixuan Li. Dream the impossible: Outlier imagination with  
524 diffusion models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine  
525 (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 60878–60901. Curran  
526 Associates, Inc., 2023.
- 527 Xuefeng Du, Zhen Fang, Ilias Diakonikolas, and Yixuan Li. How does unlabeled data provably help  
528 out-of-distribution detection? In *The Twelfth International Conference on Learning Representa-*  
529 *tions*, 2024a.
- 530 Xuefeng Du, Reshmi Ghosh, Robert Sim, Ahmed Salem, Vitor Carvalho, Emily Lawton, Yixuan Li,  
531 and Jack W Stokes. Vlmguard: Defending vlms against malicious prompts via unlabeled data.  
532 *arXiv preprint arXiv:2410.00296*, 2024b.
- 533  
534 Yue Fan, Anna Kukleva, Dengxin Dai, and Bernt Schiele. Ssb: Simple but strong baseline for  
535 boosting performance of open-set semi-supervised learning. In *Proceedings of the IEEE/CVF*  
536 *International Conference on Computer Vision (ICCV)*, pp. 16068–16078, October 2023.
- 537 Zhen Fang, Yixuan Li, Jie Lu, Jiahua Dong, Bo Han, and Feng Liu. Is out-of-distribution de-  
538 tection learnable? In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh  
539 (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 37199–37213. Curran  
Associates, Inc., 2022.

- 540 Stanislav Fort, Jie Ren, and Balaji Lakshminarayanan. Exploring the limits of out-of-distribution  
541 detection. *Advances in Neural Information Processing Systems*, 34:7068–7081, 2021.  
542
- 543 Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model  
544 uncertainty in deep learning. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings*  
545 *of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine*  
546 *Learning Research*, pp. 1050–1059, New York, New York, USA, 20–22 Jun 2016. PMLR.
- 547 Ido Galil, Mohammed Dabbah, and Ran El-Yaniv. A framework for benchmarking class-out-of-  
548 distribution detection and its application to imagenet. In *The Eleventh International Conference on*  
549 *Learning Representations*, 2023.  
550
- 551 Saurabh Garg, Yifan Wu, Alexander J Smola, Sivaraman Balakrishnan, and Zachary Lipton. Mixture  
552 proportion estimation and pu learning:a modern approach. In M. Ranzato, A. Beygelzimer,  
553 Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing*  
554 *Systems*, volume 34, pp. 8532–8544. Curran Associates, Inc., 2021.
- 555 Saurabh Garg, Sivaraman Balakrishnan, and Zachary Lipton. Domain adaptation under open set label  
556 shift. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in*  
557 *Neural Information Processing Systems*, volume 35, pp. 22531–22546. Curran Associates, Inc.,  
558 2022.  
559
- 560 Yonatan Geifman and Ran El-Yaniv. SelectiveNet: A deep neural network with an integrated  
561 reject option. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th*  
562 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning*  
563 *Research*, pp. 2151–2159. PMLR, 09–15 Jun 2019.
- 564 Chuanxing Geng, Qifei Li, Xinrui Wang, Dong Liang, Songcan Chen, and Pong C Yuen. Lod:  
565 Loss-difference ood detection by intentionally label-noisifying unlabeled wild data. *arXiv preprint*  
566 *arXiv:2505.12952*, 2025.  
567
- 568 Soumya Suvra Ghosal, Yiyu Sun, and Yixuan Li. How to overcome curse-of-dimensionality for  
569 out-of-distribution detection? In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
570 volume 38, pp. 19849–19857, 2024.
- 571 Tieliang Gong, Guangtao Wang, Jieping Ye, Zongben Xu, and Ming Lin. Margin based pu learning.  
572 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.  
573
- 574 Jun-Yi Hang and Min-Ling Zhang. Binary decomposition: A problem transformation perspective for  
575 open-set semi-supervised learning. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian  
576 Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st*  
577 *International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning*  
578 *Research*, pp. 17505–17518. PMLR, 21–27 Jul 2024.
- 579 Rundong He, Rongxue Li, Zhongyi Han, Xihong Yang, and Yilong Yin. Topological structure  
580 learning for weakly-supervised out-of-distribution detection. In *Proceedings of the 31st ACM*  
581 *International Conference on Multimedia*, pp. 4858–4866, 2023.  
582
- 583 Matthias Hein, Maksym Andriushchenko, and Julian Bitterwolf. Why relu networks yield high-  
584 confidence predictions far away from the training data and how to mitigate the problem. In  
585 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 41–50,  
586 2019.  
587
- 588 Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution  
589 examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- 590 Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution  
591 examples in neural networks. In *International Conference on Learning Representations*, 2017.  
592
- 593 Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier  
exposure. In *International Conference on Learning Representations*, 2019.

- 594 Cho-Jui Hsieh, Nagarajan Natarajan, and Inderjit Dhillon. Pu learning for matrix completion. In  
595 Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine*  
596 *Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2445–2453, Lille, France,  
597 07–09 Jul 2015. PMLR.
- 598  
599 Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional  
600 shifts in the wild. *Advances in Neural Information Processing Systems*, 34:677–689, 2021a.
- 601  
602 Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional  
603 shifts in the wild. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.),  
604 *Advances in Neural Information Processing Systems*, 2021b.
- 605  
606 Peter J Huber. Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, 35  
(1):73–101, 1964.
- 607  
608 Taewon Jeong and Heeyoung Kim. Ood-maml: Meta-learning for few-shot out-of-distribution  
609 detection and classification. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin  
610 (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 3907–3916. Curran  
611 Associates, Inc., 2020.
- 612  
613 Julian Katz-Samuels, Julia B Nakhleh, Robert Nowak, and Yixuan Li. Training ood detectors in their  
614 natural habitats. In *International Conference on Machine Learning*, pp. 10848–10865. PMLR,  
2022a.
- 615  
616 Julian Katz-Samuels, Julia B Nakhleh, Robert Nowak, and Yixuan Li. Training OOD detectors  
617 in their natural habitats. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari,  
618 Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine*  
619 *Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 10848–10865. PMLR,  
17–23 Jul 2022b.
- 620  
621 Sungyeon Kim, Dongwon Kim, Minsu Cho, and Suha Kwak. Proxy anchor loss for deep metric  
622 learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp.  
623 3235–3244, 2020. doi: 10.1109/CVPR42600.2020.00330.
- 624  
625 Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. Being bayesian, even just a bit, fixes  
626 overconfidence in relu networks. In *International conference on machine learning*, pp. 5436–5446.  
PMLR, 2020.
- 627  
628 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 629  
630 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive  
631 uncertainty estimation using deep ensembles. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach,  
632 R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing*  
633 *Systems*, volume 30. Curran Associates, Inc., 2017.
- 634  
635 Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. Training confidence-calibrated classifiers for  
636 detecting out-of-distribution samples. *arXiv preprint arXiv:1711.09325*, 2017.
- 637  
638 Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting  
639 out-of-distribution samples and adversarial attacks. In S. Bengio, H. Wallach, H. Larochelle,  
640 K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing*  
641 *Systems*, volume 31. Curran Associates, Inc., 2018a.
- 642  
643 Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting  
644 out-of-distribution samples and adversarial attacks. *Advances in neural information processing*  
645 *systems*, 31, 2018b.
- 646  
647 Fabien Letouzey, François Denis, and Rémi Gilleron. Learning from positive and unlabeled examples.  
In *International Conference on Algorithmic Learning Theory*, pp. 71–85. Springer, 2000.
- Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image  
detection in neural networks. In *International Conference on Learning Representations*, 2018.

- 648 Jeremiah Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax Weiss, and Balaji Lakshminarayanan.  
649 Simple and principled uncertainty estimation with deterministic deep learning via distance aware-  
650 ness. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural*  
651 *Information Processing Systems*, volume 33, pp. 7498–7512. Curran Associates, Inc., 2020a.
- 652 Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection.  
653 In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural*  
654 *Information Processing Systems*, volume 33, pp. 21464–21475. Curran Associates, Inc., 2020b.
- 655 Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection.  
656 *Advances in neural information processing systems*, 33:21464–21475, 2020c.
- 657 Wesley J Maddox, Pavel Izmailov, Timur Garipov, Dmitry P Vetrov, and Andrew Gordon Wilson.  
658 A simple baseline for bayesian uncertainty in deep learning. In H. Wallach, H. Larochelle,  
659 A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information*  
660 *Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 661 Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. In S. Bengio,  
662 H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in*  
663 *Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- 664 Andrey Malinin and Mark Gales. Reverse kl-divergence training of prior networks: Improved  
665 uncertainty and adversarial robustness. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-  
666 Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32.  
667 Curran Associates, Inc., 2019.
- 668 Alexander Meinke and Matthias Hein. Towards neural networks that provably know when they don't  
669 know. In *International Conference on Learning Representations*, 2020.
- 670 Yifei Ming and Yixuan Li. How does fine-tuning impact out-of-distribution detection for vision-  
671 language models? *International Journal of Computer Vision*, 132(2):596–609, 2024.
- 672 Yifei Ming, Ying Fan, and Yixuan Li. POEM: Out-of-distribution detection with posterior sampling.  
673 In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato  
674 (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of  
675 *Proceedings of Machine Learning Research*, pp. 15650–15665. PMLR, 17–23 Jul 2022.
- 676 Yifei Ming, Yiyou Sun, Ousmane Dia, and Yixuan Li. How to exploit hyperspherical embed-  
677 dings for out-of-distribution detection? In *The Eleventh International Conference on Learning*  
678 *Representations*, 2023. URL <https://openreview.net/forum?id=aEFaE0W5pAd>.
- 679 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al.  
680 Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep*  
681 *learning and unsupervised feature learning*, volume 2011, pp. 4. Granada, 2011.
- 682 Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence  
683 predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision*  
684 *and pattern recognition*, pp. 427–436, 2015.
- 685 Gang Niu, Marthinus Christoffel Du Plessis, Tomoya Sakai, Yao Ma, and Masashi Sugiyama. Theo-  
686 retical comparisons of positive-unlabeled learning against positive-negative learning. *Advances in*  
687 *neural information processing systems*, 29, 2016.
- 688 Bo Peng, Yadan Luo, Yonggang Zhang, Yixuan Li, and Zhen Fang. Conjnorm: Tractable density  
689 estimation for out-of-distribution detection. In *The Twelfth International Conference on Learning*  
690 *Representations*, 2024. URL <https://openreview.net/forum?id=1pSL2cXWoz>.
- 691 Marthinus Du Plessis, Gang Niu, and Masashi Sugiyama. Convex formulation for learning from  
692 positive and unlabeled data. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd*  
693 *International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning*  
694 *Research*, pp. 1386–1394, Lille, France, 07–09 Jul 2015. PMLR.

- 702 Jie Ren, Stanislav Fort, Jeremiah Liu, Abhijit Guha Roy, Shreyas Padhy, and Balaji Lakshmi-  
703 narayanan. A simple fix to mahalanobis distance for improving near-ood detection. *arXiv preprint*  
704 *arXiv:2106.09022*, 2021.
- 705  
706 Jie Ren, Jiaming Luo, Yao Zhao, Kundan Krishna, Mohammad Saleh, Balaji Lakshminarayanan, and  
707 Peter J Liu. Out-of-distribution detection and selective generation for conditional language models.  
708 In *The Eleventh International Conference on Learning Representations*, 2022.
- 709 Kuniaki Saito, Donghyun Kim, and Kate Saenko. Openmatch: Open-set semi-supervised  
710 learning with open-set consistency regularization. In M. Ranzato, A. Beygelzimer,  
711 Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural In-*  
712 *formation Processing Systems*, volume 34, pp. 25956–25967. Curran Associates, Inc.,  
713 2021. URL [https://proceedings.neurips.cc/paper\\_files/paper/2021/](https://proceedings.neurips.cc/paper_files/paper/2021/file/dalle8cd1811acb79ccf0fd62cd58f86-Paper.pdf)  
714 [file/dalle8cd1811acb79ccf0fd62cd58f86-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/dalle8cd1811acb79ccf0fd62cd58f86-Paper.pdf).
- 715 Vikash Sehwal, Mung Chiang, and Prateek Mittal. {SSD}: A unified framework for self-supervised  
716 outlier detection. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=v5gjXpmR8J>.
- 717  
718  
719 Yiyou Sun and Yixuan Li. Dice: Leveraging sparsification for out-of-distribution detection. In  
720 *European Conference on Computer Vision*, pp. 691–708. Springer, 2022.
- 721 Yiyou Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations.  
722 *Advances in Neural Information Processing Systems*, 34:144–157, 2021.
- 723  
724 Yiyou Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest  
725 neighbors. In *International Conference on Machine Learning*, pp. 20827–20840. PMLR, 2022.
- 726 Jihoon Tack, Sangwoo Mo, Jongheon Jeong, and Jinwoo Shin. Csi: Novelty detec-  
727 tion via contrastive learning on distributionally shifted instances. In H. Larochelle,  
728 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural In-*  
729 *formation Processing Systems*, volume 33, pp. 11839–11852. Curran Associates, Inc.,  
730 2020a. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/8965f76632d7672e7d3cf29c87ecaa0c-Paper.pdf)  
731 [file/8965f76632d7672e7d3cf29c87ecaa0c-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/8965f76632d7672e7d3cf29c87ecaa0c-Paper.pdf).
- 732  
733 Jihoon Tack, Sangwoo Mo, Jongheon Jeong, and Jinwoo Shin. Csi: Novelty detection via contrastive  
734 learning on distributionally shifted instances. *Advances in neural information processing systems*,  
735 33:11839–11852, 2020b.
- 736 Leitian Tao, Xuefeng Du, Jerry Zhu, and Yixuan Li. Non-parametric outlier synthesis. In *The*  
737 *Eleventh International Conference on Learning Representations*, 2023.
- 738 Rheeeya Uppaal, Junjie Hu, and Yixuan Li. Is fine-tuning needed? pre-trained language models  
739 are near perfect for out-of-domain detection. In Anna Rogers, Jordan Boyd-Graber, and Naoaki  
740 Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational*  
741 *Linguistics (Volume 1: Long Papers)*, pp. 12813–12832, Toronto, Canada, July 2023. Association  
742 for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.717.
- 743  
744 Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a  
745 single deep deterministic neural network. In *International conference on machine learning*, pp.  
746 9690–9700. PMLR, 2020.
- 747 Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual-  
748 logit matching. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
749 *recognition*, pp. 4921–4930, 2022a.
- 750 Haoran Wang, Weitang Liu, Alex Bocchieri, and Yixuan Li. Can multi-label classification networks  
751 know what they don’t know? *Advances in Neural Information Processing Systems*, 34:29074–  
752 29087, 2021.
- 753  
754 Qizhou Wang, Feng Liu, Yonggang Zhang, Jing Zhang, Chen Gong, Tongliang Liu, and Bo Han.  
755 Watermarking for out-of-distribution detection. *Advances in Neural Information Processing*  
*Systems*, 35:15545–15557, 2022b.

- 756 Qizhou Wang, Zhen Fang, Yonggang Zhang, Feng Liu, Yixuan Li, and Bo Han. Learning to augment  
757 distributions for out-of-distribution detection. In A. Oh, T. Naumann, A. Globerson, K. Saenko,  
758 M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36,  
759 pp. 73274–73286. Curran Associates, Inc., 2023a.
- 760 Qizhou Wang, Junjie Ye, Feng Liu, Quanyu Dai, Marcus Kalander, Tongliang Liu, Jianye HAO,  
761 and Bo Han. Out-of-distribution detection with implicit outlier transformation. In *The Eleventh  
762 International Conference on Learning Representations*, 2023b.
- 763 Yu Wang, Jingjing Zou, Jingyang Lin, Qing Ling, Yingwei Pan, Ting Yao, and Tao Mei. Out-of-  
764 distribution detection via conditional kernel independence model. In S. Koyejo, S. Mohamed,  
765 A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing  
766 Systems*, volume 35, pp. 36411–36425. Curran Associates, Inc., 2022c.
- 767 Zerun Wang, Liuyu Xiang, Lang Huang, Jiafeng Mao, Ling Xiao, and Toshihiko Yamasaki. Sco-  
768 match: Alleviating overtrusting in open-set semi-supervised learning. In *European Conference on  
769 Computer Vision*, pp. 217–233. Springer, 2024.
- 770 Hongxin Wei, Renchunzi Xie, Hao Cheng, Lei Feng, Bo An, and Yixuan Li. Mitigating neural network  
771 overconfidence with logit normalization. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba  
772 Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference  
773 on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 23631–  
774 23644. PMLR, 17–23 Jul 2022.
- 775 Yeming Wen, Dustin Tran, and Jimmy Ba. Batchensemble: an alternative approach to efficient  
776 ensemble and lifelong learning. In *International Conference on Learning Representations*, 2020.
- 777 M. B. Wilk and R. Gnanadesikan. Probability plotting methods for the analysis of data. *Biometrika*,  
778 55(1):1–17, 1968. ISSN 00063444, 14643510.
- 779 Qitian Wu, Yiting Chen, Chenxiao Yang, and Junchi Yan. Energy-based out-of-distribution detection  
780 for graph neural networks. In *The Eleventh International Conference on Learning Representations*,  
781 2023.
- 782 Danfei Xu and Misha Denil. Positive-unlabeled reward learning. In *Conference on Robot Learning*,  
783 pp. 205–219. PMLR, 2021.
- 784 Jinggang Yang, Haoqi Wang, Litong Feng, Xiaopeng Yan, Huabin Zheng, Wayne Zhang, and  
785 Ziwei Liu. Semantically coherent out-of-distribution detection. In *Proceedings of the IEEE/CVF  
786 International Conference on Computer Vision*, pp. 8301–8309, 2021.
- 787 Jinggang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi  
788 Wang, Guangyao Chen, Bo Li, Yiyu Sun, et al. Openood: Benchmarking generalized out-of-  
789 distribution detection. *Advances in Neural Information Processing Systems*, 35:32598–32611,  
790 2022.
- 791 Jinggang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection:  
792 A survey. *International Journal of Computer Vision*, 132(12):5635–5662, 2024.
- 793 Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun:  
794 Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv  
795 preprint arXiv:1506.03365*, 2015.
- 796 Sergey Zagoruyko. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- 797 Jingyang Zhang, Jinggang Yang, Pengyun Wang, Haoqi Wang, Yueqian Lin, Haoran Zhang, Yiyu  
798 Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, et al. Openood v1. 5: Enhanced benchmark for  
799 out-of-distribution detection. *arXiv preprint arXiv:2306.09301*, 2023.
- 800 Yonggang Zhang, Jie Lu, Bo Peng, Zhen Fang, and Yiu ming Cheung. Learning to shape in-  
801 distribution feature space for out-of-distribution detection. In *The Thirty-eighth Annual Confer-  
802 ence on Neural Information Processing Systems*, 2024. URL [https://openreview.net/  
803 forum?id=1Du3mMP5YN](https://openreview.net/forum?id=1Du3mMP5YN).

810 Yunrui Zhao, Qianqian Xu, Yangbangyan Jiang, Peisong Wen, and Qingming Huang. Dist-pu:  
811 Positive-unlabeled learning from a label distribution perspective. In *Proceedings of the IEEE/CVF*  
812 *Conference on Computer Vision and Pattern Recognition*, pp. 14461–14470, 2022.

813  
814 Haotian Zheng, Qizhou Wang, Zhen Fang, Xiaobo Xia, Feng Liu, Tongliang Liu, and Bo Han.  
815 Out-of-distribution detection learning with unreliable out-of-distribution sources. *Advances in*  
816 *Neural Information Processing Systems*, 36:72110–72123, 2023.

817 Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10  
818 million image database for scene recognition. *IEEE transactions on pattern analysis and machine*  
819 *intelligence*, 40(6):1452–1464, 2017.

820  
821 Zhi Zhou, Lan-Zhe Guo, Zhazhan Cheng, Yu-Feng Li, and Shiliang Pu. Step: Out-of-distribution  
822 detection in the presence of limited in-distribution labeled data. In M. Ranzato, A. Beygelzimer,  
823 Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing*  
824 *Systems*, volume 34, pp. 29168–29180. Curran Associates, Inc., 2021.

825 Jianing Zhu, Yu Geng, Jiangchao Yao, Tongliang Liu, Gang Niu, Masashi Sugiyama, and Bo Han. Di-  
826 versified outlier exposure for out-of-distribution detection via informative extrapolation. *Advances*  
827 *in Neural Information Processing Systems*, 36:22702–22734, 2023.

828  
829 Yao Zhu, YueFeng Chen, Chuanlong Xie, Xiaodan Li, Rong Zhang, Hui Xue', Xiang Tian, bolun  
830 zheng, and Yaowu Chen. Boosting out-of-distribution detection with typical features. In S. Koyejo,  
831 S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information*  
832 *Processing Systems*, volume 35, pp. 20758–20769. Curran Associates, Inc., 2022.

833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

## A ADDITIONAL STUDIES

### A.1 ELEMENT-WISE MEDIAN VS GEOMETRIC MEDIAN

Another pertinent question to ask is why we chose the element-wise median over the geometric median (Acharya et al., 2024) for outlier detection in Medix. To answer this question, we conducted a preliminary experiment using CIFAR-10 as the InD dataset and SVHN as the OOD dataset. Specifically,  $\mathcal{S}_{\text{wild}}$  consists of 5k samples drawn from CIFAR-10, ensuring that these samples are disjoint from the training set used to train the model  $\phi_{\mathcal{S}_{\text{in}}}$ , which we leverage to compute  $\bar{\nabla}_{\text{in}}$ . We incrementally add SVHN OOD samples to  $\mathcal{S}_{\text{wild}}$  and track the behavior of the  $L_2$ -norm deviation between  $\bar{\nabla}_{\text{in}}$  and the element-wise median of the gradients of the wild dataset as well as the proportion of OOD samples removed. The results, shown in the Figure 3, demonstrate that for the same number of OOD samples in the wild dataset, the element-wise median identified a significantly higher proportion of OOD samples as outliers compared to the geometric median. This indicates that the element-wise median is more sensitive to distributional shifts, making it a more effective and reliable choice for filtering in our method.

Table 3: Effect of hyperparameters  $\epsilon$  and  $k$  on OOD detection.

Method	FPR95↓	$\epsilon$	$k$
DICE	88.35	–	–
ASH	21.36	–	–
CSI	64.70	–	–
KNN+	32.21	–	–
OE	2.86	–	–
Energy	2.71	–	–
Medix	0.16	0.005	20000
Medix	0.20	0.0005	20000
Medix	0.68	0.005	10000

### A.2 HYPERPARAMETER SELECTION AND SENSITIVITY ANALYSIS

We conducted an ablation study to assess the sensitivity of Medix to the hyperparameters  $\epsilon$  and  $k$ . For this experiment, we employed CIFAR-100 as the InD dataset and SVHN as the OOD dataset. As shown in Table 3, Medix exhibits strong robustness to variations in the values of  $\epsilon$  and  $k$ . In particular, Medix achieves a remarkably low FPR95 of 0.16 when using  $\epsilon = 0.005$  and  $k = 20000$ . Even when the values of  $\epsilon$  and  $k$  are varied—e.g., reducing  $\epsilon$  to 0.0005 or halving  $k$  to 10000—the performance remains competitive (FPR95 of 0.20 and 0.68, respectively), demonstrating a graceful degradation rather than a sharp drop. The results indicate that while optimal hyperparameter selection is important, Medix maintains its effectiveness across a variety of hyperparameter choices, surpassing the baselines. This insensitivity to exact hyperparameter tuning makes Medix a reliable choice in real-world deployments, where exhaustive tuning may not be feasible.

### A.3 COMPREHENSIVE COMPARISON WITH RECENT COMPETITIVE METHODS

We have extended our comparisons to include competitive baselines that employ diverse contrasting learning objectives, such as SSD+ (Sehwag et al., 2021), ProxyAnchor (Kim et al., 2020), and CIDER Ming et al. (2023). Although these methods use contrasting learning objectives, which we do not, we included them for a comprehensive comparison. Additionally, we included more recent methods, such as CONJ (Peng et al., 2024), DRL (Zhang et al., 2024), Vim (Wang et al., 2022a), and VOS (Du et al., 2022c), to provide a more thorough evaluation. For this experiment, we use CIFAR-100 as InD

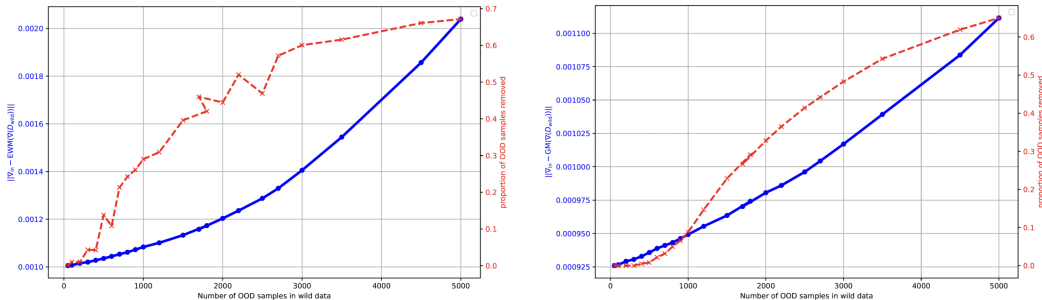


Figure 3: Comparison of element-wise median (EWM) and geometric median (GM).

Table 4: OOD detection performance comparison of Medix and recent competitive baselines on CIFAR-100 as InD data. Performance averaged over five runs; best results are highlighted in **bold**.

Methods	OOD Datasets									
	SVHN		PLACES365		LSUN-C		TEXTURES		AVERAGE	
	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑	FPR95↓	AUROC↑
MSP	78.89	79.80	84.38	74.21	83.47	75.28	86.51	72.53	83.31	75.46
Mahalanobis	87.09	80.62	84.63	73.89	84.15	79.43	61.72	84.87	79.40	79.70
ODIN	70.16	84.88	82.16	75.19	76.36	80.10	85.28	75.23	78.49	78.85
Energy	66.91	85.25	81.41	76.37	59.77	86.69	79.01	79.96	71.78	82.07
ReAct	50.93	88.75	83.55	73.10	64.02	80.31	64.40	81.95	65.72	81.03
KNN	46.25	90.39	82.08	75.44	60.85	85.61	62.39	83.95	62.89	83.85
Vim	73.42	84.62	85.34	69.34	86.96	69.74	74.56	76.23	80.07	74.98
VOS	43.24	82.80	76.85	78.63	73.61	84.69	57.57	87.31	62.82	83.36
CSI	44.53	92.65	79.08	76.27	75.58	83.78	61.61	86.47	65.20	84.79
ProxyAnchor	87.21	82.43	70.10	79.84	37.19	91.68	65.64	84.99	65.04	84.74
SSD+	31.19	94.19	77.74	79.90	79.39	85.18	66.63	86.18	63.74	86.36
KNN+	39.23	92.78	80.74	77.58	48.99	89.30	57.15	88.35	56.53	87.00
ASH	52.96	90.19	72.62	76.38	75.18	76.52	56.17	86.75	64.23	82.46
CIDER	23.09	95.16	79.63	73.43	16.16	96.33	43.87	90.42	40.69	88.84
CONJ	46.19	90.44	80.81	75.83	60.45	85.90	62.13	83.77	62.40	83.99
DRL	20.15	94.07	76.64	77.55	16.97	94.63	31.97	92.09	36.43	89.59
Medix	<b>0.48</b>	<b>99.87</b>	<b>24.52</b>	<b>93.42</b>	<b>0.73</b>	<b>99.84</b>	<b>8.99</b>	<b>97.92</b>	<b>8.68</b>	<b>97.76</b>
(Ours)	$\pm 0.05$	$\pm 0.02$	$\pm 1.76$	$\pm 0.36$	$\pm 0.08$	$\pm 0.03$	$\pm 0.72$	$\pm 0.24$	$\pm 0.65$	$\pm 0.16$

dataset and train the ResNet-34 model<sup>1</sup> following the setup in Ming et al. (2023); Zhang et al. (2024). The model is trained using stochastic gradient descent with momentum 0.9, and weight decay  $10^{-4}$  for 500 epochs. The initial learning rate is 0.5 with cosine scheduling and the batch size is 512. We use the checkpoints provided by Ming et al. (2023)<sup>2</sup>. The results, as shown in Table 4, demonstrate that Medix outperforms even these competitive baselines by a significant margin, including those trained with contrastive learning objectives, achieving an average FPR95 of 8.68% and an average AUROC of 97.76%. Furthermore, Medix provides a provably low error rate, offering theoretical guarantees that most alternative approaches do not.

#### A.4 EVALUATION ON LARGE-SCALE, COMPLEX UNSEEN OOD SETTINGS

We next investigate whether Medix can handle large-scale wild OOD data under the unseen OOD setting. This setting is more challenging for two main reasons: (1) it leverages a large-scale wild OOD dataset ( $\mathbb{P}_{\text{out}}$ ), increasing both data complexity and computational demand, and (2) the test-time OOD distribution ( $\mathbb{P}_{\text{out}}^{\text{test}}$ ) is intentionally chosen to be distributionally different from the wild OOD data used during training, i.e.,  $\mathbb{P}_{\text{out}}^{\text{test}} \neq \mathbb{P}_{\text{out}}$ , which better reflects realistic and challenging deployment scenarios. We used CIFAR-100 as the labeled InD data, 300K Random Images dataset Hendrycks et al. (2019) as the unlabeled wild OOD data and tested on the SVHN dataset as the test OOD data. This setup introduces a greater level of complexity because the large-scale wild OOD data (300K Random Images) is significantly different from the OOD test data (SVHN). To evaluate the performance of our approach, we compare it against baselines that also leverage wild data in training. The results in the Table 5 highlight the superior performance of our method, Medix, compared to the baselines, with a significantly lower FPR95 ( $41.29 \pm 1.2$ ) and higher AUROC ( $87.25 \pm 0.6$ ), showing its effectiveness in distinguishing between InD and OOD data.

Table 5: Comparison of OOD detection performance on large-scale unseen OOD data.

Method	FPR ↓	AUROC ↑
OE	$68.80 \pm 2.8$	$82.89 \pm 1.1$
Energy (w/ OE)	$69.81 \pm 2.4$	$85.59 \pm 1.0$
WOODS	$69.41 \pm 2.7$	$86.76 \pm 0.8$
Medix (ours)	<b><math>41.29 \pm 1.2</math></b>	<b><math>87.25 \pm 0.9</math></b>

<sup>1</sup>Since ResNet-34 is the standard model across these works, we adopt the same model to ensure consistency and facilitate a fair comparison.

<sup>2</sup><https://github.com/deeplearning-wisc/cider>

#### 972 A.5 EVALUATING THE IMPACT OF PSEUDO-LABEL QUALITY ON OOD FILTERING AND 973 ROBUSTNESS 974

975 Another important question to ask is whether the quality of pseudo-labels  $\hat{y}_{\tilde{x}_i}$ , particularly in terms of  
976 softmax confidence, affects the performance of OOD filtering, and whether a simple pre-filtering step  
977 that removes low-confidence pseudo-labels could improve the robustness of the model. To answer  
978 this question, we conducted an experiment to evaluate how filtering low-confidence pseudo-labels  
979 affects OOD filtering. For this experiment, we used CIFAR-10 as the InD dataset and LSUN-Resize  
980 as the OOD dataset. Specifically, we computed the softmax probabilities for each pseudo-labeled  
981 sample and discarded those with low-confidence predictions, setting a threshold of 0.6. After filtering,  
982 we found that 15.98% of the samples were removed from the training set.

983 The results showed that there was virtually no difference in performance between the method with  
984 filtering and the method without filtering. Both methods yielded a FPR95 of 0.01%, but with filtering,  
985 AUROC increased slightly from 99.98% to 99.99%. Thus, removing low-confidence pseudo-labels  
986 doesn't significantly impact the robustness of the model or its ability to detect OOD samples, showing  
987 that our method is resilient to noisy or low-confidence labels, and further filtering steps are unlikely  
988 to yield meaningful improvements.

#### 989 A.6 EVALUATING COMPUTATION AND MEMORY EFFICIENCY IN MEDIX FILTERING 990

991 We now turn to the question of whether Medix's filtering phase is computationally feasible and  
992 memory-efficient on large datasets. In response to this question, we conducted profiling experiments  
993 on an NVIDIA A100-SXM4-80GB GPU using approximately 15,000 unlabeled samples for each  
994 InD-OOD pair. For this experiment, we used CIFAR-10 and CIFAR-100 as the InD datasets and  
995 LSUN-Resize as the OOD dataset. As seen from the results in Table 6, the GPU memory usage  
996 remains modest, and the filtering phase completes within a tractable timeframe, even when handling  
997 thousands of unlabeled samples. This demonstrates that the Medix filtering process is computationally  
998 feasible and does not impose excessive memory overhead, even on large datasets.

999 Table 6: Profiling results for Medix filtering on NVIDIA A100 80GB GPU.  
1000

InD-OOD	Wall Clock Time (s)	Peak GPU Memory (MB)	Current GPU Memory (MB)
CIFAR10 – LSUN-Resize	4497.17	99.46	31.74
CIFAR100 – LSUN-Resize	5478.36	99.56	31.79

#### 1006 A.7 COMPARISON WITH SEMI-SUPERVISED OPEN-SET RECOGNITION METHODS 1007

1008 Our work differs from recent semi-supervised open-set recognition methods (Saito et al., 2021; Fan  
1009 et al., 2023; Hang & Zhang, 2024; Wang et al., 2024) in several key ways. For example, in OpenML  
1010 Saito et al. (2021), the main problem is to handle outliers in semi-supervised learning (SSL) when  
1011 training a standard classifier, whereas in our setting, the main challenge is to detect OOD samples  
1012 from unlabeled wild data and train a dedicated OOD detector classifier. In other words, while SSL  
1013 methods (Saito et al., 2021; Fan et al., 2023; Hang & Zhang, 2024; Wang et al., 2024) aim to improve  
1014 classification despite outliers, our approach enhances OOD detection in an open-world setting where  
1015 labeled OOD data is unavailable.

1016 Secondly, these SSL methods aim to train a classifier that is robust to OOD samples in SSL, treating  
1017 outliers as noise, while in our setting, the goal is to explicitly detect OOD samples and train an OOD  
1018 detector. In other words, SSL methods focus on suppressing OOD samples during training to improve  
1019 SSL performance, whereas Medix actively detects and leverages these OOD samples to build a more  
1020 effective and reliable OOD detection system.

1021 Lastly, the two approaches differ in their techniques: OpenMatch is based on consistency regular-  
1022 ization, where the main idea is to enforce consistency across different stochastic transformations of  
1023 the same input. This is mathematically modeled through soft constraints on the classifier's outputs,  
1024 encouraging smooth decision boundaries. On the other hand, Medix relies on median-based filtering  
1025 for outlier detection. The method uses statistical robustness properties of the median, leveraging its  
insensitivity to extreme values to identify potential OOD samples. The theoretical analysis derives

1026 error bounds based on the contamination effect (proportion of OOD samples) and the concentration  
1027 effect (sub-Gaussian behavior of InD gradients), ensuring that OOD detection error remains low.  
1028

## 1029 B BROADER IMPACTS AND LIMITATIONS 1030

1031 As machine learning continues to advance, tackling the challenges of OOD detection has become  
1032 essential for ensuring robust and reliable model performance in real-world applications. We introduce  
1033 a median-centric framework for OOD detection that enhances OOD handling and model safety.  
1034 The impact of our research goes beyond theoretical advancements, with practical applications in  
1035 healthcare, autonomous systems, and finance. By improving OOD detection, we address a key  
1036 challenge in model deployment, fostering greater trust and adoption of machine learning technologies.  
1037 Future work could explore integrating it with generative models for outlier synthesis or adapting it to  
1038 dynamic environments where OOD distributions evolve over time.

1039 A potential limitation of Medix lies in its reliance on a mixture of unlabeled InD and OOD data,  
1040 which, in practice, may be noisy, corrupted, or inconsistently sampled. However, this assumption is  
1041 not unique to our method—it reflects the inherent uncertainty of real-world deployment environments  
1042 and is a common starting point in robust learning theory. To make our assumptions explicit and  
1043 verifiable, we rigorously formalize them in the main theorems and provide a precise geometric  
1044 condition in Remark 4.3 to justify the necessary separation required for reliable filtering. Moreover,  
1045 we extend our guarantees to looser settings in Theorem C.3, removing the need for sub-Gaussian tails  
1046 and demonstrating the resilience of our framework under weaker distributional assumptions. These  
1047 efforts underscore that while Medix does make simplifying assumptions, they are both interpretable  
1048 and relaxable, forming a principled foundation for future improvements in OOD detection under  
1049 realistic noise conditions.

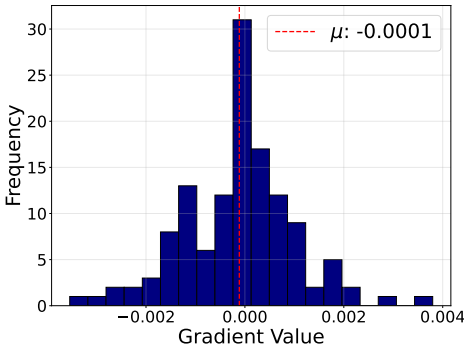
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

C PROOF OF THEOREM

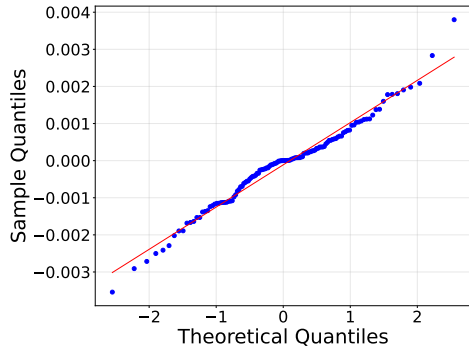
To aid clarity and consistency throughout the theoretical and empirical sections, we summarize the key notation used in this paper in Table 7. This includes definitions for gradient-related quantities, sample sizes, and filtering terms, as well as parameters that appear in our concentration bounds. Unless stated otherwise, we use boldface for vectors,  $\mathcal{X}$  to denote the input space, and assume all gradients are computed with respect to a fixed pretrained model  $f_\phi$ .

Table 7: Notation summary. We use bold symbols for vectors and  $\mathcal{X}$  to denote the input space.

Symbol	Meaning
$\mathcal{X}$	Input space (e.g., images in $[0, 1]^d$ )
$f_\phi$	Model parameterized by $\phi$
$\ell(f_\phi(x))$	Loss function evaluated at input $x$
$\nabla\ell(f_\phi(x))$	Gradient of loss at point $x$
$\mathbf{g}(x)$	Shorthand for $\nabla\ell(f_\phi(x))$ , gradient vector
$\bar{\mathbf{g}}_{\text{in}}$	Empirical mean of gradients over InD points
$m$	Total number of unlabeled samples in $\mathcal{S}_{\text{wild}}$
$m_{\text{in}} = (1 - \pi)m$	Number of InD samples
$m_{\text{out}} = \pi m$	Number of OOD samples
$\pi \in (0, 1)$	OOD contamination ratio
$\mathcal{S}_{\text{wild}}$	Unlabeled mixture of InD and OOD data
$\mathcal{S}_{\text{in}}^*$	Inliers selected by the median filter
$\mathcal{S}_{\text{out}}^*$	OOD points selected by the median filter
$\mathcal{B}_\varepsilon(\bar{\mathbf{g}}_{\text{in}})$	Ball of radius $\varepsilon$ centered at $\bar{\mathbf{g}}_{\text{in}}$
$m_\varepsilon$	Number of InD points inside $\mathcal{B}_\varepsilon$
$m_\varepsilon^*$	Number of selected inliers inside $\mathcal{B}_\varepsilon$
$m_{\text{sep}}^*$	Number of well-separated OODs in selected set
$\text{ERR}_{\text{in}}$	Fraction of InD points misclassified as OOD
$\text{ERR}_{\text{out}}$	Fraction of OOD points misclassified as InD
$\sigma^2$	Sub-Gaussian proxy variance of each gradient coordinate
$\varepsilon$	Threshold for gradient deviation, $\varepsilon = \sigma\sqrt{2\log(2dm_{\text{in}})}$
$\delta \in (0, 1)$	Confidence parameter for concentration bounds
$d$	Dimensionality of gradient vectors



(a) InD gradient histogram.



(b) InD gradient Q-Q plot.

Figure 4: Illustration of InD sample gradients exhibiting sub-Gaussian behavior in each coordinate. (left) Histogram of gradient values (CIFAR-100 InD data) showing concentration around the mean with light tails, consistent with sub-Gaussianity. (right) Q-Q plot comparing empirical quantiles of InD gradients against a theoretical Gaussian distribution, confirming alignment with sub-Gaussianity.

### C.1 BOUND ON THE INLIER MISCLASSIFICATION RATE

**Theorem C.1** (Inlier Misclassification Bound). *Assume that the gradients of InD points in  $\mathcal{S}_{\text{wild}}$  are i.i.d., and each coordinate is sub-Gaussian with variance proxy  $\sigma^2$ . Let  $\epsilon = \sigma\sqrt{2\log(2dm_{\text{in}})}$ , and fix any confidence level  $\delta \in (0, 1)$ . Then, with probability at least  $1 - \delta$ , the inlier misclassification rate of the element-wise median (EWM) filtering rule satisfies:*

$$\text{ERR}_{\text{in}} \leq \underbrace{\frac{1}{m_{\text{in}}} + 2\sqrt{\frac{\log(1/\delta)}{2m_{\text{in}}}}}_{\text{Concentration term}} + \underbrace{\frac{\pi}{2(1-\pi)}}_{\text{Contamination term}}.$$

*Proof.* We structure our argument into two key steps. In **Step 1**, we establish a high-probability concentration bound for the gradients of InD samples around their mean, denoted as  $\bar{\nabla}_{\text{in}}$ . We define an InD sample whose gradient lies within this high-probability region as a *good inlier*. In **Step 2**, we leverage a “swapping” argument on the EWM objective to demonstrate that removing more than a  $(2\eta + \frac{\pi}{2}(1-\pi))$  fraction of *good inliers*—where  $\eta$  represents the fraction of InD samples with gradients outside the high-probability region (i.e., non-good InD samples)—would contradict optimality. Combining these two steps, we derive the stated upper bound on the inlier misclassification rate,  $\text{ERR}_{\text{in}}$ .

**Step 1: Concentration of InD Gradients.** Let  $\nabla\ell(f_\phi(\mathbf{x})) \in \mathbb{R}^d$  be the gradient (with respect to  $\phi$ ) of a loss evaluated at an InD sample  $\mathbf{x}$ . Suppose each coordinate of this gradient is sub-Gaussian with variance proxy  $\sigma^2$ , centered at the same mean vector  $\bar{\nabla}_{\text{in}}$ . That is, for each coordinate  $j \in \{1, \dots, d\}$ ,

$$\mathbb{P}\left(|\nabla\ell(f_\phi(\mathbf{x}))_j - \bar{\nabla}_{\text{in},j}| \geq \epsilon\right) \leq 2 \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right). \quad (6)$$

Define

$$v := \nabla\ell(f_\phi(\mathbf{x})) - \bar{\nabla}_{\text{in}} \in \mathbb{R}^d.$$

It is straightforward to show that

$$\|v\|_2 \leq \sqrt{d} \|v\|_\infty,$$

so if  $\|v\|_2 > \epsilon\sqrt{d}$  then there must exist at least one coordinate  $j$  with  $|v_j| > \epsilon$ .

Hence, combining equation 6 across  $j = 1, \dots, d$  via a union bound yields

$$\mathbb{P}\left(\|\nabla\ell(f_\phi(\mathbf{x})) - \bar{\nabla}_{\text{in}}\|_2 > \epsilon\sqrt{d}\right) \leq 2d \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right). \quad (7)$$

Let  $m_{\text{in}} = (1-\pi)m$  denote the total number of InD points in the set  $\mathcal{S}_{\text{wild}}$ . We define an InD point  $\mathbf{x}$  to be *bad* if

$$\|\nabla\ell(f_\phi(\mathbf{x})) - \bar{\nabla}_{\text{in}}\|_2 > \epsilon\sqrt{d}. \quad (8)$$

Our goal is to show that, with high probability, only a small fraction of these  $m_{\text{in}}$  inliers can be bad.

From our sub-Gaussian assumption and applying the union bound across the  $d$  coordinates, we define a constant  $p \in [0, 1]$  as follows:

$$p := \mathbb{P}\left(\|\nabla\ell(f_\phi(\mathbf{x})) - \bar{\nabla}_{\text{in}}\|_2 > \epsilon\sqrt{d}\right) \leq 2d \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right). \quad (9)$$

Here,  $\mathbf{x}$  is a generic InD sample, and we assume that the samples are i.i.d. across the  $m_{\text{in}}$  inliers.

Let  $X$  be the random variable denoting the number of bad inliers among the  $m_{\text{in}}$  samples. For each inlier  $\mathbf{x}_i$ , define the indicator variable

$$I_i := \begin{cases} 1, & \text{if } \|\nabla\ell(f_\phi(\mathbf{x}_i)) - \bar{\nabla}_{\text{in}}\|_2 > \epsilon\sqrt{d}, \\ 0, & \text{otherwise.} \end{cases} \quad (10)$$

1188 Then,  $X$  can be expressed as  
1189

$$1190 \quad X = \sum_{i=1}^{m_{\text{in}}} I_i. \quad (11)$$

1193 By the linearity of expectation, we have  
1194

$$1195 \quad \mathbb{E}[X] = \sum_{i=1}^{m_{\text{in}}} \mathbb{E}[I_i] \quad (12)$$

$$1196 \quad = \sum_{i=1}^{m_{\text{in}}} \mathbb{P}(\mathbf{x}_i \text{ is bad}) \quad (13)$$

$$1200 \quad \text{from equation 9} \\ 1201 \quad \leq m_{\text{in}} \cdot p \quad (14)$$

$$1202 \quad \leq m_{\text{in}} \cdot 2d \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right). \quad (15)$$

1205 Observing that  $X$  is a sum of  $m_{\text{in}}$  Bernoulli( $p$ ) random variables, each bounded between 0 and 1, we  
1206 can apply Hoeffding's inequality to obtain, for any  $t > 0$ ,  
1207

$$1208 \quad \mathbb{P}(X \geq \mathbb{E}[X] + t) \leq \exp\left(-\frac{2t^2}{m_{\text{in}}}\right). \quad (16)$$

1211 Setting  
1212

$$1213 \quad t := \sqrt{\frac{m_{\text{in}} \log\left(\frac{1}{\delta}\right)}{2}}, \quad (17)$$

1216 we obtain  
1217

$$1218 \quad \mathbb{P}\left(X \geq \mathbb{E}[X] + \sqrt{\frac{m_{\text{in}} \log(1/\delta)}{2}}\right) \leq \delta. \quad (18)$$

1221 With probability at least  $1 - \delta$ , we therefore have  
1222

$$1223 \quad X \leq \mathbb{E}[X] + \sqrt{\frac{m_{\text{in}} \log(1/\delta)}{2}}. \quad (19)$$

1225 Dividing both sides by  $m_{\text{in}}$  yields an upper bound on the fraction of bad inliers:  
1226

$$1227 \quad \frac{X}{m_{\text{in}}} \leq \frac{\mathbb{E}[X]}{m_{\text{in}}} + \sqrt{\frac{\log(1/\delta)}{2m_{\text{in}}}} \quad (20)$$

$$1230 \quad \text{from equation 12} \\ 1231 \quad \leq 2d \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) + \sqrt{\frac{\log(1/\delta)}{2m_{\text{in}}}} \quad (21)$$

$$1232 \quad = \underbrace{\frac{1}{m_{\text{in}}}}_{\text{Concentration term}} + 2 \underbrace{\sqrt{\frac{\log(1/\delta)}{2m_{\text{in}}}}}_{\text{Contamination term}} + \frac{\pi}{2(1-\pi)}. \quad (22)$$

1237 Let us denote this fraction by  $\eta$ . By setting  $\epsilon = \sigma \sqrt{2 \log(2dm_{\text{in}})}$  Thus, with probability at least  
1238  $1 - \delta$ , no more than an  $O\left(d \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) + \sqrt{\frac{\log(1/\delta)}{m_{\text{in}}}}\right)$  fraction of inliers have gradients lying  
1239 outside the radius  $\epsilon \sqrt{d}$  around  $\bar{\nabla}_{\text{in}}$ . These are precisely the *bad inliers*, while the remainder are  
1240 called *good inliers*.  
1241

1242 **Step 2: Contradiction via a Swapping Argument** We demonstrate that if  $\mathcal{S}^*$  excludes an excessive  
 1243 number of InD, specifically when

$$1244 \text{ERR}_{\text{in}} > 2\eta + \frac{\pi}{2(1-\pi)}, \quad (23)$$

1245 then it is possible to construct a new set  $\hat{\mathcal{S}}$  whose EWM is strictly closer to  $\bar{\nabla}_{\text{in}}$  than that of  $\mathcal{S}^*$ ,  
 1246 thereby contradicting the optimality of  $\mathcal{S}^*$ .

1247 Let  $m_\varepsilon$  be the number of indices  $i$  such that either  $x_i$  is OOD or  $x_i$  is an InD point whose gradient  
 1248 lies outside the closed  $\ell_2$  ball

$$1249 B_\varepsilon = \left\{ g \in \mathbb{R}^d : \|g - \bar{\nabla}_{\text{in}}\|_2 \leq \varepsilon\sqrt{d} \right\}. \quad (24)$$

1250 Let  $m_{\text{in}}^*$  be the number of InD points contained in the candidate minimizer  $\mathcal{S}^*$ .

1251 Assume for contradiction that

$$1252 \text{ERR}_{\text{in}} > 2\eta + \frac{\pi}{2(1-\pi)}, \quad (25)$$

1253 Define the swap size as

$$1254 m'_\varepsilon := \min \{ m_\varepsilon, (\text{ERR}_{\text{in}} - \eta)m_{\text{in}} \}. \quad (26)$$

1255 This quantity counts how many good inliers (those inside  $B_\varepsilon$ ) are currently excluded from  $\mathcal{S}^*$  and  
 1256 available for swapping in.

1257 Construct a new set  $\hat{\mathcal{S}}$  by removing  $m'_\varepsilon$  points from  $\mathcal{S}^*$  that lie outside  $B_\varepsilon$  (OOD points or bad inliers)  
 1258 and inserting  $m'_\varepsilon$  good inliers from  $B_\varepsilon$  in their place. This maintains cardinality:

$$1259 |\hat{\mathcal{S}}| = |\mathcal{S}^*|.$$

1260 Recall that

$$1261 \text{ERR}_{\text{in}} := \frac{m_\varepsilon - m_{\text{in}}^*}{m_{\text{in}}}, \quad (27)$$

1262 which measures the fraction of good inliers excluded from  $\mathcal{S}^*$ . Rearranging, we obtain:

$$1263 m_\varepsilon - m_{\text{in}}^* = \text{ERR}_{\text{in}} \cdot m_{\text{in}}. \quad (28)$$

1264 Now we know, for contradiction, that

$$1265 \text{ERR}_{\text{in}} > 2\eta + \frac{\pi}{2(1-\pi)}, \quad (29)$$

1266 Subtracting  $\eta$  from both sides and multiplying by  $m_{\text{in}}$ , we obtain:

$$1267 (\text{ERR}_{\text{in}} - \eta)m_{\text{in}} > \left( \eta + \frac{\pi}{2(1-\pi)} \right) m_{\text{in}}. \quad (30)$$

1268 We know that at most a  $\eta$  fraction of inliers are outside  $B_\varepsilon$ , and the fraction of OOD points among all  
 1269 samples is  $\pi$ . Therefore, the total number of samples in the dataset that lie outside  $B_\varepsilon$  satisfies:

$$1270 m_\varepsilon \leq \eta m_{\text{in}} + \frac{\pi}{1-\pi} m_{\text{in}} = \left( \eta + \frac{\pi}{1-\pi} \right) m_{\text{in}}. \quad (31)$$

1271 It follows that

$$1272 \left( \eta + \frac{\pi}{2(1-\pi)} \right) m_{\text{in}} > \frac{1}{2} \cdot \left( \eta + \frac{\pi}{1-\pi} \right) m_{\text{in}} \geq \frac{m_\varepsilon}{2}. \quad (32)$$

1296 Combining with the previous inequality gives:

$$1297 \quad (\text{ERR}_{\text{in}} - \eta)m_{\text{in}} > \frac{m_\varepsilon}{2}, \quad (33)$$

1300 So  $m'_\varepsilon = m_\varepsilon/2 + \eta m_{\text{in}}$ , which implies

$$1303 \quad m_\varepsilon - m'_\varepsilon = \frac{m_\varepsilon}{2} - \eta m_{\text{in}}, \quad (34)$$

$$1307 \quad m_{\text{in}}^* + m'_\varepsilon = m_{\text{in}}^* + \frac{m_\varepsilon}{2} + \eta m_{\text{in}}. \quad (35)$$

1310 Because  $\eta m_{\text{in}} > 0$ , it follows that

$$1311 \quad m_{\text{in}}^* + m'_\varepsilon > m_\varepsilon - m'_\varepsilon, \quad (36)$$

1313 so the points lying inside  $B_\varepsilon$  form a strict majority in  $\widehat{\mathcal{S}}$ .

1315 Fix any coordinate  $j \in \{1, \dots, d\}$ . Let  $X^{(j)} = \{x_1^{(j)}, \dots, x_m^{(j)}\}$  be the  $j$ -th coordinate values of all  
1316 gradients in  $\mathcal{S}^*$  and  $\widehat{X}^{(j)} = \{\widehat{x}_1^{(j)}, \dots, \widehat{x}_m^{(j)}\}$  the same in  $\widehat{\mathcal{S}}$ . Let  $\mu_j = \bar{\nabla}_{\text{in},j}$ .

1317 Each inserted gradient lies inside the interval  $I_j = [\mu_j - \varepsilon, \mu_j + \varepsilon]$ , and each removed gradient lies  
1318 outside  $I_j$ . So for some coordinate  $j^*$ , the number of values inside  $I_{j^*}$  increases from below  $m/2$  to  
1319 above  $m/2$  due to the swaps. Therefore:

$$1321 \quad \text{med}(X^{(j^*)}) \notin I_{j^*}, \quad \text{med}(\widehat{X}^{(j^*)}) \in I_{j^*}, \quad (37)$$

1323 which implies

$$1325 \quad |\text{med}(\widehat{X}^{(j^*)}) - \mu_{j^*}| \leq \varepsilon \quad \text{and} \quad |\text{med}(\widehat{X}^{(j^*)}) - \mu_{j^*}| < |\text{med}(X^{(j^*)}) - \mu_{j^*}|. \quad (38)$$

1327 For the remaining coordinates  $j \neq j^*$ , the swapped points either leave the majority unchanged or  
1328 increase it, so the median distance from  $\mu_j$  does not increase and possibly decreases.

1329 Thus, component-wise:

$$1331 \quad |\text{EWM}_j(G_{\widehat{\mathcal{S}}}) - \mu_j| \leq |\text{EWM}_j(G_{\mathcal{S}^*}) - \mu_j|, \quad (39)$$

1333 with strict inequality in at least one coordinate  $j^*$ , implying

$$1335 \quad \|\text{EWM}(G_{\widehat{\mathcal{S}}}) - \bar{\nabla}_{\text{in}}\|_2 < \|\text{EWM}(G_{\mathcal{S}^*}) - \bar{\nabla}_{\text{in}}\|_2. \quad (40)$$

1337 This contradicts the assumption that  $\mathcal{S}^*$  is a minimizer of the optimization problem, which concludes  
1338 the proof.

1339

1340

$$1341 \quad \Rightarrow \quad \text{ERR}_{\text{in}} \leq 2\eta + \frac{\pi}{2(1-\pi)}. \quad (41)$$

1342

1343

1344

## 1345 C.2 BOUND ON THE OUTLIER RETENTION RATE

1346

1347 We now establish a non-asymptotic upper bound on the fraction of OOD points that remain in  
1348 the subset returned by the MEDIX filter. This result mirrors the inlier misclassification bound  
1349 (Appendix C.1), but instead of relying on inlier concentration, it exploits a separation assumption  
between the OOD and InD gradient means.

**Theorem C.2** (Outlier Retention Bound under Vector Separation). *Assume that for every OOD point  $x \sim P_{\text{out}}$ , the gradient  $g(x) := \nabla \ell(f_{\varphi_{\text{in}}}(x)) \in \mathbb{R}^d$  satisfies:*

(I) **Sub-Gaussianity.** *Each coordinate  $g_j(x)$  is sub-Gaussian with variance proxy  $\sigma_{\text{out}}^2$ .*

(II) **Vector Mean Separation.** *The mean OOD gradient vector  $\mu_{\text{out}}$  satisfies*

$$\|\mu_{\text{out}} - \bar{\nabla}_{\text{in}}\|_2 \geq \Delta\sqrt{d}, \quad \text{for some } \Delta > 0.$$

Fix any tolerance  $\varepsilon \in (0, \Delta)$  and confidence parameter  $\delta \in (0, 1)$ . Define:

$$p_{\text{out}}(\varepsilon) := 2d \exp\left(-\frac{(\Delta - \varepsilon)^2}{2\sigma_{\text{out}}^2}\right), \quad \eta_{\text{out}} := p_{\text{out}}(\varepsilon) + \sqrt{\frac{\log(1/\delta)}{2m_{\text{out}}}}.$$

Then with probability at least  $1 - \delta$ , the outlier retention rate satisfies:

$$\text{ERR}_{\text{out}} := \frac{|\mathcal{S}_{\text{in}}^* \cap \text{OOD}|}{m_{\text{out}}} \leq 2d \exp\left(-\frac{(\Delta - \varepsilon)^2}{2\sigma_{\text{out}}^2}\right) + \sqrt{\frac{\log(1/\delta)}{2m_{\text{out}}}} + \frac{1 - \pi}{2\pi}. \quad (42)$$

*Proof. Step 1: Concentration of OOD Gradients.*

By assumption, each coordinate of  $g(x)$  is sub-Gaussian, and the mean vector is separated from  $\bar{\nabla}_{\text{in}}$  by at least  $\Delta\sqrt{d}$ . Therefore, for any OOD point  $x$ , the event

$$\|g(x) - \bar{\nabla}_{\text{in}}\|_2 \leq \varepsilon\sqrt{d} \quad (43)$$

has probability at most  $p_{\text{out}}(\varepsilon)$ .

**Counting the ‘‘ambiguous’’ OOD points.** For every OOD example  $x_i$  define the indicator

$$J_i := \mathbf{1}\left\{\|g(x_i) - \bar{\nabla}_{\text{in}}\|_2 \leq \varepsilon\sqrt{d}\right\}, \quad i = 1, \dots, m_{\text{out}}.$$

Each  $J_i \sim \text{Bernoulli}(p_{\text{out}}(\varepsilon))$  and the  $J_i$ ’s are independent. Let

$$Y := \sum_{i=1}^{m_{\text{out}}} J_i, \quad \mu := \mathbb{E}[Y] = m_{\text{out}} p_{\text{out}}(\varepsilon).$$

Because the indicators  $J_i$  take values in  $\{0, 1\}$ , Hoeffding’s inequality states that for any  $s > 0$

$$\Pr\{Y - \mu \geq s\} \leq \exp\left(-\frac{2s^2}{m_{\text{out}}}\right).$$

Choose  $s = \sqrt{\frac{m_{\text{out}} \log(1/\delta)}{2}}$ ; then equation C.2 equals  $\delta$ , and with probability at least  $1 - \delta$

$$Y \leq m_{\text{out}} \left(p_{\text{out}}(\varepsilon) + \sqrt{\frac{\log(1/\delta)}{2m_{\text{out}}}}\right) = m_{\text{out}} \eta_{\text{out}}.$$

Define the high-probability event

$$\mathcal{E}_1 := \{Y \leq m_{\text{out}} \eta_{\text{out}}\}, \quad (35)$$

which holds with probability at least  $1 - \delta$ .

**Step 2: Contradiction via a Swapping Argument** Suppose, toward contradiction, that

$$\text{ERR}_{\text{out}} > \eta_{\text{out}} + \frac{1 - \pi}{2\pi}. \quad (44)$$

Let  $m_{\text{out}}^* := |\mathcal{S}_{\text{in}}^* \cap \text{OOD}|$ , and define  $m_\varepsilon^* := \#\{\text{ambiguous OODs in } \mathcal{S}_{\text{in}}^*\}$ . Then, under  $\mathcal{E}_1$ ,

$$m_{\text{sep}}^* := m_{\text{out}}^* - m_\varepsilon^* > \frac{1 - \pi}{2\pi} m_{\text{out}} = \frac{1}{2} m_{\text{in}}. \quad (45)$$

Thus, more than half the points in  $\mathcal{S}_{\text{in}}^*$  are separated OODs whose gradients lie outside the ball

$$B_\varepsilon := \left\{ g \in \mathbb{R}^d : \|g - \bar{\nabla}_{\text{in}}\|_2 \leq \varepsilon\sqrt{d} \right\}. \quad (46)$$

By the definition of the coordinate-wise median, this implies that in every coordinate, the majority of entries in  $G_{\mathcal{S}_{\text{in}}^*}^{(j)}$  lie outside the interval  $[\bar{\nabla}_{\text{in},j} - \varepsilon, \bar{\nabla}_{\text{in},j} + \varepsilon]$ , hence:

$$\|\text{EWM}(G_{\mathcal{S}_{\text{in}}^*}) - \bar{\nabla}_{\text{in}}\|_2^2 \geq \varepsilon^2 d. \quad (47)$$

Now consider a new subset  $\hat{\mathcal{S}}$  formed by swapping the  $m_{\text{sep}}^*$  separated OODs in  $\mathcal{S}_{\text{in}}^*$  with an equal number of InD points whose gradients lie inside  $B_\varepsilon$ . Such inliers exist because fewer than  $m_{\text{out}}\eta_{\text{out}}$  are ambiguous, and the total number of InD points is  $m_{\text{in}}$ .

Then at least half the entries in every coordinate of  $G_{\hat{\mathcal{S}}}^{(j)}$  lie inside the interval  $[\bar{\nabla}_{\text{in},j} - \varepsilon, \bar{\nabla}_{\text{in},j} + \varepsilon]$ , implying:

$$\|\text{EWM}(G_{\hat{\mathcal{S}}}) - \bar{\nabla}_{\text{in}}\|_2^2 \leq \varepsilon^2 d. \quad (48)$$

Combining with equation 47 yields a contradiction:

$$\|\bar{\nabla}_{\text{in}} - \text{EWM}(G_{\hat{\mathcal{S}}})\|_2^2 < \|\bar{\nabla}_{\text{in}} - \text{EWM}(G_{\mathcal{S}_{\text{in}}^*})\|_2^2. \quad (49)$$

This contradicts the optimality of  $\mathcal{S}_{\text{in}}^*$ , and so the assumption must be false. Therefore, our assumption must be false. With probability at least  $1 - \delta$ , we conclude that the outlier retention rate is bounded by

$$\text{ERR}_{\text{out}} \leq 2d \exp\left(-\frac{(\Delta - \varepsilon)^2}{2\sigma_{\text{out}}^2}\right) + \sqrt{\frac{\log(1/\delta)}{2m_{\text{out}}}} + \frac{1 - \pi}{2\pi}, \quad (50)$$

□

### C.3 THEOREM C.3 WITHOUT SUB-GAUSSIANITY

**Theorem C.3** (Inlier Misclassification Bound without Sub-Gaussianity). *Assume that the gradients of InD points in  $\mathcal{S}_{\text{wild}}$  are i.i.d., and that each coordinate has variance  $\sigma^2$  and finite fourth moment bounded by  $\mu_4$ , i.e.,*

$$\mathbb{E}[(\nabla\ell(f_\phi(x)))_j - \bar{\nabla}_{\text{in},j}]^2 = \sigma^2, \quad \mathbb{E}[(\nabla\ell(f_\phi(x)))_j - \bar{\nabla}_{\text{in},j}]^4 \leq \mu_4.$$

Let  $m_{\text{in}} = (1 - \pi)m$  denote the number of InD samples in  $\mathcal{S}_{\text{wild}}$ , and fix any tolerance  $\varepsilon > \sigma$  and confidence level  $\delta \in (0, 1)$ . Then, with probability at least  $1 - \delta$ , the inlier misclassification rate of the element-wise median (EWM) filtering rule satisfies:

$$\text{ERR}_{\text{in}} \leq \underbrace{2 \left( \frac{\mu_4 - \sigma^4}{d(\varepsilon^2 - \sigma^2)^2} + \sqrt{\frac{\log(1/\delta)}{2m_{\text{in}}}} \right)}_{\text{Concentration term}} + \underbrace{\frac{\pi}{2(1 - \pi)}}_{\text{Contamination term}}.$$

1458 **Differences with the original bound.** The original bound, which assumes sub-Gaussianity, features  
 1459 a decay term

$$1460 \quad 2 \exp\left(-\frac{\epsilon^2 d}{2\sigma_{\text{sub}}^2}\right)$$

1461  
 1462 that decreases exponentially with dimension  $d$ , making the bound extremely tight for large  $d$ . In  
 1463 contrast, the new bound, relying solely on the assumption of finite fourth moments, exhibits a decay  
 1464 term

$$1465 \quad \frac{\mu_4 - \sigma^4}{d(\epsilon^2 - \sigma^2)^2}$$

1466  
 1467 that decays as  $O(1/d)$ , yielding a looser bound for larger  $d$ . The original bound requires sub-Gaussian  
 1468 tails, while the new one only assumes finite fourth moments, allowing for heavier tails. In particular,  
 1469 both bounds share the finite-sample term

$$1470 \quad \sqrt{\frac{\log(1/\delta)}{2m_{\text{in}}}},$$

1471  
 1472 that shrinks as  $m_{\text{in}}$  increases, independent of the tail assumptions.  
 1473

1474  
 1475 *Proof.* The proof follows the structure of the original Theorem 4.1, consisting of two main steps: (1)  
 1476 establishing a concentration bound on the fraction of “bad” InD samples (those far from  $\nabla_{\text{in}}$ ) without  
 1477 relying on sub-Gaussianity, and (2) using a swapping argument to bound the inlier misclassification  
 1478 rate. The key difference is the replacement of the sub-Gaussian tail bound with a moment-based  
 1479 bound derived from the finite fourth moment assumption.  
 1480

1481 **Step 1: Concentration of InD Gradients** Define the gradient deviation for an InD sample  $x$  as:

$$1482 \quad v = \nabla \ell(f_\phi(x)) - \nabla_{\text{in}} \in \mathbb{R}^d, \quad (51)$$

1483  
 1484 where  $\nabla_{\text{in}} = \frac{1}{n} \sum_{(x_i, y_i) \in S_{\text{in}}} \nabla \ell(f_\phi(x_i))$  is the mean gradient over the labeled InD training set  $S_{\text{in}}$ ,  
 1485 and we assume  $\nabla_{\text{in}}$  is the population mean. An InD sample is classified as “bad” if:

$$1486 \quad \|v\|_2 > \epsilon\sqrt{d}, \quad (52)$$

1487  
 1488 and “good” otherwise. Our goal is to bound the fraction of bad inliers among the  $m_{\text{in}}$  InD samples in  
 1489  $S_{\text{wild}}$ .  
 1490

1491  
 1492 Since  $\|v\|_2^2 = \sum_{j=1}^d v_j^2$ , we analyze the  $\ell_2$  norm via the sum of squared coordinates. Assume that  
 1493 the coordinates  $v_j = \nabla \ell(f_\phi(x))_j - \nabla_{\text{in},j}$  are independent across  $j = 1, \dots, d$ , each with:

$$1494 \quad \mathbb{E}[v_j] = 0 \quad (\text{assuming } \nabla_{\text{in},j} = \mathbb{E}[\nabla \ell(f_\phi(x))_j]), \quad (53)$$

$$1495 \quad \mathbb{E}[v_j^2] = \sigma^2, \quad \mathbb{E}[v_j^4] \leq \mu_4. \quad (54)$$

1500  
 1501 Thus:

$$1502 \quad \mathbb{E}[\|v\|_2^2] = \mathbb{E}\left[\sum_{j=1}^d v_j^2\right] = d\sigma^2, \quad (55)$$

$$1503 \quad \text{Var}(\|v\|_2^2) = \text{Var}\left(\sum_{j=1}^d v_j^2\right) = \sum_{j=1}^d \text{Var}(v_j^2). \quad (56)$$

1512 Compute the variance of  $v_j^2$ :

1513

1514

$$1515 \quad \text{Var}(v_j^2) = \mathbb{E}[v_j^4] - (\mathbb{E}[v_j^2])^2 \leq \mu_4 - \sigma^4, \quad (57)$$

1516

1517 so:

1518

1519

$$1520 \quad \text{Var}(\|v\|_2^2) \leq d(\mu_4 - \sigma^4). \quad (58)$$

1521

1522 Using Chebyshev's inequality for  $\|v\|_2^2$ :

1523

1524

$$1525 \quad \mathbb{P}(|\|v\|_2^2 - d\sigma^2| \geq t) \leq \frac{d(\mu_4 - \sigma^4)}{t^2}. \quad (59)$$

1526

1527 We need  $\mathbb{P}(\|v\|_2 > \epsilon\sqrt{d})$ , which corresponds to  $\|v\|_2^2 > \epsilon^2 d$ . Set the threshold:

1528

1529

$$1530 \quad \|v\|_2^2 > \epsilon^2 d \implies \|v\|_2^2 - d\sigma^2 > \epsilon^2 d - d\sigma^2 = d(\epsilon^2 - \sigma^2). \quad (60)$$

1531

1532 Thus:

1533

1534

$$1535 \quad \mathbb{P}(\|v\|_2^2 > \epsilon^2 d) = \mathbb{P}(\|v\|_2^2 - d\sigma^2 > d(\epsilon^2 - \sigma^2)) \leq \frac{d(\mu_4 - \sigma^4)}{[d(\epsilon^2 - \sigma^2)]^2} = \frac{\mu_4 - \sigma^4}{d(\epsilon^2 - \sigma^2)^2}, \quad (61)$$

1536

1537 provided  $\epsilon > \sigma$  so that  $\epsilon^2 - \sigma^2 > 0$ . Define:

1538

1539

$$1540 \quad p = \mathbb{P}(\|\nabla \ell(f_\phi(x)) - \nabla_{\text{in}}\|_2 > \epsilon\sqrt{d}) \leq \frac{\mu_4 - \sigma^4}{d(\epsilon^2 - \sigma^2)^2}. \quad (62)$$

1541

1542 Let  $X$  be the number of bad inliers among the  $m_{\text{in}}$  InD samples in  $S_{\text{wild}}$ , where each InD sample is bad with probability at most  $p$ . Then:

1543

1544

$$1545 \quad \mathbb{E}[X] \leq m_{\text{in}} p \leq m_{\text{in}} \cdot \frac{\mu_4 - \sigma^4}{d(\epsilon^2 - \sigma^2)^2}. \quad (63)$$

1546

1547 Since  $X = \sum_{i=1}^{m_{\text{in}}} I_i$ , where  $I_i = 1$  if the  $i$ -th InD sample is bad and 0 otherwise, and  $I_i$  are i.i.d. Bernoulli( $p$ ) variables bounded in  $[0, 1]$ , apply Hoeffding's inequality:

1548

1549

$$1550 \quad \mathbb{P}(X \geq \mathbb{E}[X] + t) \leq \exp\left(-\frac{2t^2}{m_{\text{in}}}\right). \quad (64)$$

1551

1552 Set  $t = \sqrt{\frac{m_{\text{in}} \log(1/\delta)}{2}}$ , so:

1553

1554

$$1555 \quad \mathbb{P}\left(X \geq \mathbb{E}[X] + \sqrt{\frac{m_{\text{in}} \log(1/\delta)}{2}}\right) \leq \delta. \quad (65)$$

1556

1557 With probability at least  $1 - \delta$ :

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

$$X \leq m_{\text{in}} \cdot \frac{\mu_4 - \sigma^4}{d(\epsilon^2 - \sigma^2)^2} + \sqrt{\frac{m_{\text{in}} \log(1/\delta)}{2}}. \quad (66)$$

The fraction of bad inliers is:

$$\eta = \frac{X}{m_{\text{in}}} \leq \frac{\mu_4 - \sigma^4}{d(\epsilon^2 - \sigma^2)^2} + \sqrt{\frac{\log(1/\delta)}{2m_{\text{in}}}}. \quad (67)$$

This  $\eta$  represents the upper bound on the fraction of InD samples whose gradients deviate from  $\bar{\nabla}_{\text{in}}$  by more than  $\epsilon\sqrt{d}$ , relying only on finite fourth moments and coordinate independence.

**Step 2: Contradiction via a Swapping Argument** We demonstrate that if  $\mathcal{S}^*$  excludes an excessive number of InD, specifically when

$$\text{ERR}_{\text{in}} > 2\eta + \frac{\pi}{2(1-\pi)}, \quad (68)$$

then it is possible to construct a new set  $\hat{\mathcal{S}}$  whose EWM is strictly closer to  $\bar{\nabla}_{\text{in}}$  than that of  $\mathcal{S}^*$ , thereby contradicting the optimality of  $\mathcal{S}^*$ .

Let  $m_\epsilon$  be the number of indices  $i$  such that either  $x_i$  is OOD or  $x_i$  is an InD point whose gradient lies outside the closed  $\ell_2$  ball

$$B_\epsilon = \left\{ g \in \mathbb{R}^d : \|g - \bar{\nabla}_{\text{in}}\|_2 \leq \epsilon\sqrt{d} \right\}. \quad (69)$$

Let  $m_{\text{in}}^*$  be the number of InD points contained in the candidate minimizer  $\mathcal{S}^*$ .

Assume for contradiction that

$$\text{ERR}_{\text{in}} > 2\eta + \frac{\pi}{2(1-\pi)}, \quad (70)$$

Define the swap size as

$$m'_\epsilon := \min \{ m_\epsilon, (\text{ERR}_{\text{in}} - \eta)m_{\text{in}} \}. \quad (71)$$

This quantity counts how many good inliers (those inside  $B_\epsilon$ ) are currently excluded from  $\mathcal{S}^*$  and available for swapping in.

Construct a new set  $\hat{\mathcal{S}}$  by removing  $m'_\epsilon$  points from  $\mathcal{S}^*$  that lie outside  $B_\epsilon$  (OOD points or bad inliers) and inserting  $m'_\epsilon$  good inliers from  $B_\epsilon$  in their place. This maintains cardinality:

$$|\hat{\mathcal{S}}| = |\mathcal{S}^*|.$$

Recall that

$$\text{ERR}_{\text{in}} := \frac{m_\epsilon - m_{\text{in}}^*}{m_{\text{in}}}, \quad (72)$$

which measures the fraction of good inliers excluded from  $\mathcal{S}^*$ . Rearranging, we obtain:

$$m_\epsilon - m_{\text{in}}^* = \text{ERR}_{\text{in}} \cdot m_{\text{in}}. \quad (73)$$

Now we know, for contradiction, that

$$\text{ERR}_{\text{in}} > 2\eta + \frac{\pi}{2(1-\pi)}, \quad (74)$$

1620 Subtracting  $\eta$  from both sides and multiplying by  $m_{\text{in}}$ , we obtain:

$$1621 \quad (\text{ERR}_{\text{in}} - \eta)m_{\text{in}} > \left( \eta + \frac{\pi}{2(1-\pi)} \right) m_{\text{in}}. \quad (75)$$

1622 We know that at most a  $\eta$  fraction of inliers are outside  $B_\varepsilon$ , and the fraction of OOD points among all  
1623 samples is  $\pi$ . Therefore, the total number of samples in the dataset that lie outside  $B_\varepsilon$  satisfies:

$$1624 \quad m_\varepsilon \leq \eta m_{\text{in}} + \frac{\pi}{1-\pi} m_{\text{in}} = \left( \eta + \frac{\pi}{1-\pi} \right) m_{\text{in}}. \quad (76)$$

1625 It follows that

$$1626 \quad \left( \eta + \frac{\pi}{2(1-\pi)} \right) m_{\text{in}} > \frac{1}{2} \cdot \left( \eta + \frac{\pi}{1-\pi} \right) m_{\text{in}} \geq \frac{m_\varepsilon}{2}. \quad (77)$$

1627 Combining with the previous inequality gives:

$$1628 \quad (\text{ERR}_{\text{in}} - \eta)m_{\text{in}} > \frac{m_\varepsilon}{2}, \quad (78)$$

1629 So  $m'_\varepsilon = m_\varepsilon/2 + \eta m_{\text{in}}$ , which implies

$$1630 \quad m_\varepsilon - m'_\varepsilon = \frac{m_\varepsilon}{2} - \eta m_{\text{in}}, \quad (79)$$

$$1631 \quad m_{\text{in}}^* + m'_\varepsilon = m_{\text{in}}^* + \frac{m_\varepsilon}{2} + \eta m_{\text{in}}. \quad (80)$$

1632 Because  $\eta m_{\text{in}} > 0$ , it follows that

$$1633 \quad m_{\text{in}}^* + m'_\varepsilon > m_\varepsilon - m'_\varepsilon, \quad (81)$$

1634 so the points lying inside  $B_\varepsilon$  form a strict majority in  $\widehat{\mathcal{S}}$ .

1635 Fix any coordinate  $j \in \{1, \dots, d\}$ . Let  $X^{(j)} = \{x_1^{(j)}, \dots, x_m^{(j)}\}$  be the  $j$ -th coordinate values of all  
1636 gradients in  $\mathcal{S}^*$  and  $\widehat{X}^{(j)} = \{\widehat{x}_1^{(j)}, \dots, \widehat{x}_m^{(j)}\}$  the same in  $\widehat{\mathcal{S}}$ . Let  $\mu_j = \bar{\nabla}_{\text{in},j}$ .

1637 Each inserted gradient lies inside the interval  $I_j = [\mu_j - \varepsilon, \mu_j + \varepsilon]$ , and each removed gradient lies  
1638 outside  $I_j$ . So for some coordinate  $j^*$ , the number of values inside  $I_{j^*}$  increases from below  $m/2$  to  
1639 above  $m/2$  due to the swaps. Therefore:

$$1640 \quad \text{med}(X^{(j^*)}) \notin I_{j^*}, \quad \text{med}(\widehat{X}^{(j^*)}) \in I_{j^*}, \quad (82)$$

1641 which implies

$$1642 \quad |\text{med}(\widehat{X}^{(j^*)}) - \mu_{j^*}| \leq \varepsilon \quad \text{and} \quad |\text{med}(\widehat{X}^{(j^*)}) - \mu_{j^*}| < |\text{med}(X^{(j^*)}) - \mu_{j^*}|. \quad (83)$$

1643 For the remaining coordinates  $j \neq j^*$ , the swapped points either leave the majority unchanged or  
1644 increase it, so the median distance from  $\mu_j$  does not increase and possibly decreases.

1645 Thus, component-wise:

$$1646 \quad |\text{EWM}_j(G_{\widehat{\mathcal{S}}}) - \mu_j| \leq |\text{EWM}_j(G_{\mathcal{S}^*}) - \mu_j|, \quad (84)$$

1647 with strict inequality in at least one coordinate  $j^*$ , implying

$$1648 \quad \|\text{EWM}(G_{\widehat{\mathcal{S}}}) - \bar{\nabla}_{\text{in}}\|_2 < \|\text{EWM}(G_{\mathcal{S}^*}) - \bar{\nabla}_{\text{in}}\|_2. \quad (85)$$

1649 This contradicts the assumption that  $\mathcal{S}^*$  is a minimizer of the optimization problem, which concludes  
1650 the proof.

$$1651 \quad \Rightarrow \quad \text{ERR}_{\text{in}} \leq 2\eta + \frac{\pi}{2(1-\pi)}. \quad (86)$$

1652  $\square$