LoD: Loss-difference OOD Detection by Intentionally Label-Noisifying Unlabeled Wild Data

Chuanxing Geng^{1,2,3}, Qifei Li¹, Xinrui Wang¹, Dong Liang^{1,3}, Songcan Chen^{1,3} and Pong C. Yuen^{2*}

¹College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics ²Department of Computer Science, Hong Kong Baptist University

Department of Computer Science, Hong Kong Daptist Oniversit

³MIIT Key Laboratory of Pattern Analysis and Machine Intelligence {gengchuanxing, liqifei, wangxinrui, liangdong, s.chen}@nuaa.edu.cn, pcyuen@comp.hkbu.edu.hk

Abstract

Using unlabeled wild data containing both indistribution (ID) and out-of-distribution (OOD) data to improve the safety and reliability of models has recently received increasing attention. Existing methods either design customized losses for labeled ID and unlabeled wild data then perform joint optimization, or first filter out OOD data from the latter then learn an OOD detector. While achieving varying degrees of success, two potential issues remain: (i) Labeled ID data typically dominates the learning of models, inevitably making models tend to fit OOD data as IDs; (ii) The selection of thresholds for identifying OOD data in unlabeled wild data usually faces dilemma due to the unavailability of pure OOD samples. To address these issues, we propose a novel loss-difference OOD detection framework (LoD) by intentionally label-noisifying unlabeled wild data. Such operations not only enable labeled ID data and OOD data in unlabeled wild data to jointly dominate the models' learning but also ensure the distinguishability of the losses between ID and OOD samples in unlabeled wild data, allowing the classic clustering technique (e.g., K-means) to filter these OOD samples without requiring thresholds any longer. We also provide theoretical foundation for LoD's viability, and extensive experiments verify its superiority.

1 Introduction

The safety and reliability of traditional machine learning models often face challenges when deployed in realworld environments due to unexpected occurrence of outof-distribution (OOD) data [Nguyen *et al.*, 2015]. To meet this challenge, the OOD detection problem has been studied [Hendrycks and Gimpel, 2016; Yang *et al.*, 2024], which requires the models not only predict the true class of indistribution (ID) data but also effectively reject the OOD data. To date, numerous OOD detection methods have been developed [Liu *et al.*, 2020b; Abati *et al.*, 2019; Wang *et al.*, 2022; Hendrycks *et al.*, 2018; Katz-Samuels *et al.*, 2022], and among them, the methods leveraging unlabeled wild data containing ID and OOD samples to improve the performance of OOD detection has recently received increasing attention [Katz-Samuels *et al.*, 2022]. This mainly attributed to the fact that such data can be freely collected during the deployment of any machine learning model in its operational environment, while also allowing for the capture of the true test-time OOD distribution.

Despite the promise, harnessing the power of unlabeled wild data is non-trivial due to the heterogeneous mixture of ID and OOD samples. Existing methods either adopt a joint optimization strategy [Katz-Samuels *et al.*, 2022] or a twostep strategy (i.e., filtering and learning) [Du *et al.*, 2024]. The former aims to design customized losses for labeled ID and unlabeled wild data to jointly optimize the models in a semi-supervised learning manner. The latter first filters out OOD samples from the unlabeled wild data using customized OOD score (usually based on labeled ID data), then uses them along with labeled ID data to learn an OOD detector. While achieving varying degrees of success, two potential issues of these methods remain:

- ✓ The model-bias issue. Labeled ID data typically dominates the model learning in both two strategies, especially for the two-step strategy, thus inevitably making the model tend to fit OOD data as IDs.
- ✓ Threshold selection dilemma. The selection of thresholds for determining OOD samples in unlabeled wild data usually faces challenges due to the unavailability of pure OOD samples.

To address these issues, this work proposes a novel lossdifference OOD detection framework (abbreviated as LoD) by *intentionally label-noisifying* unlabeled wild data. LoD adopts the filtering and learning strategy and its key lies in the loss-difference filtering module with *intentional label-noises*. In this module, the whole unlabeled wild data is intentionally labeled as a single K + 1-th class (assuming that ID data contain K classes), and then trained together with the labeled ID data through the *fully-supervised* manner of K + 1 classification. We would like to emphasize that such operations ingeniously transform the OOD filtering problem in unlabeled wild data into a label-noise learning problem, allowing us to solve the aforementioned issues by leveraging the inherent properties in label-noise learning. In this way, the OOD

^{*}Corresponding author



Figure 1: The cross-entropy loss changes of ID (*label-noise*) and OOD (*label-clean*) samples in unlabeled wild data when they are intentionally labeled as K + 1-th class. These two types of samples typically exhibit different loss curves due to the differences in how learning progresses for each.

samples in unlabeled wild data is intentionally transformed into *label-clean* samples, while the ID counterparts become *label-noise* ones. The former naturally and seamlessly enables OOD samples in unlabeled wild data to jointly dominate the model learning with the labeled ID data, effectively addressing the model-bias issue. Meanwhile, the latter provides the key clues for differentiating ID and OOD samples in unlabeled wild data due to the significant differences in the loss curves between ID (*label-noise*) and OOD (*label-clean*) samples during training.

As shown in Figure 1, as the OOD samples in the unlabeled wild data are correctly labeled (*label-clean*), the model fits them well as learning progresses, leading to a gradual decrease and convergence of the loss curve. In contrast, the corresponding ID part, due to being incorrectly labeled (*labelnoise*) and conflicting with the originally labeled ID data, exhibits not only higher loss values but also larger fluctuations during training. Such significant and natural differences allow us to employ classic clustering models, like K-means, to filter these OOD samples without requiring thresholds any longer. In particular, we also provide theoretical foundation to support the viability of such a module. Overall, our contributions can be highlighted as follows:

- Two potential issues (i.e., the model-bias issue and threshold selection dilemma) in this OOD research line are identified, providing some new insights for the sub-sequent modeling of OOD detection.
- The OOD filtering problem in unlabeled wild data is elegantly reformulated as a label-noise learning problem, leading to a novel LoD OOD detection framework, which not only effectively addresses the modelbias issue but also circumvents the threshold selection dilemma.
- Theoretical foundation is provided to support the viability of LoD. Meanwhile, extensive experiments are also conducted to demonstrate its superiority.

2 Related Works

2.1 Out-of-Distribution Detection

To improve the safety and reliability of models in detecting OOD data, various OOD methods have been developed [Zhu et al., 2023; Zheng et al., 2023; Wang et al., 2023b; Yang et al., 2024; Li et al., 2024b; Behpour et al., 2024; Fang et al., 2024; Sharifi et al., 2025], including adopting the classification confidence or entropy, modeling the ID density, leveraging auxiliary OOD data, and more. Among these, methods using auxiliary OOD data have demonstrated encouraging OOD detection performance over the counterpart without auxiliary data [Lee et al., 2017; Bevandić et al., 2018; Malinin and Gales, 2018; Liu et al., 2020b; Chen et al., 2021; Wei et al., 2022; Du et al., 2022; Wang et al., 2023a; Sharifi et al., 2025]. Despite the promise, there are two primary limitations: First, such data may not match the true distribution of OOD data in the wild; Second, collecting such data can be labor-intensive and inflexible. To address these limitations, recent works [Katz-Samuels et al., 2022; Du et al., 2024] proposed to leverage the unlabeled "in-thewild" data due to they are freely collected during the deployment of any machine learning model in its operational environment, while also allowing for the capture of the true testtime OOD distribution.

Our work falls into this research line, and as mentioned earlier, though the methods in this research line have achieved varying degrees of success, they still face two potential weaknesses, i.e., the model-bias issue and the threshold selection dilemma. These motivate us to seek new methods to address these issues.

2.2 Training Neural Networks with Label Noises

In many applications [Guan et al., 2018], due to the cost or difficulty of manual labeling, datasets are often annotated through online queries [Yuan et al., 2024a] or crowdsourcing [Li et al., 2024a]. Such annotations inevitably contain numerous mistakes, i.e., label-noises. When trained on the data mixed clean labels and noise labels, deep neural networks have been observed to first fit label-clean data during an early learning phase, and then start memorizing the label-noise data after sufficient epochs of training [Liu et al., 2020a]. This phenomenon is independent of the optimizations used during training or the architectures of neural networks employed [Arpit et al., 2017]. In particular, during the early learning phase, label-clean and label-noise data will have different loss curves due to the difference in how learning progresses for each type. This has been exploited in many label-noise learning works [Forouzesh et al., 2022; Li et al., 2023; Yuan et al., 2024b; Lin et al., 2024; Lienen and Hüllermeier, 2024; Yue and Jha, 2024]. For more information, please refer to the recent review work [Song et al., 2022].

In this paper, we propose a novel loss-difference OOD detection framework by *intentionally label-noisifying* unlabeled wild data, which interestingly transforms the OOD filtering problem in unlabeled wild data into a label-noise learning problem. This enables us to leverage the aforementioned inherent phenomenon of label-noise learning to effectively filter OOD data from the unlabeled wild data.



Figure 2: Overview of the loss-difference OOD detection framework by intentionally label-noisifying unlabeled wild data.

3 Methodology

3.1 **Problem Formulation**

Labeled ID Data Let \mathcal{X} denote the input space and $\mathcal{Y} = \{1, ..., K\}$ represent the label space. Let $\mathcal{D}_{in}^{\text{train}} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$ denote the labeled training set drawn independently and identically from $\mathbb{P}_{\mathcal{XY}}$. \mathbb{P}_{in} is the marginal distribution of $\mathbb{P}_{\mathcal{XY}}$ on \mathcal{X} , which is also referred to as the ID distribution.

Unlabeled Wild Data The main challenge in OOD detection is the lack of labeled OOD data. In particular, the sample space for potential OOD data can be prohibitively large, making it expensive to collect labeled OOD data. To model the realistic environment, recent works [Katz-Samuels *et al.*, 2022; Du *et al.*, 2024] incorporated unlabeled wild data $\mathcal{D}_{wild} = \{\tilde{x}_1, ..., \tilde{x}_m\}$ into OOD detection. Unlabeled wild data consists of potentially both ID and OOD data, and can be freely collected upon deploying an existing model in its natural habitats. Following [Katz-Samuels *et al.*, 2022], the Huber contamination model is employed to characterize the marginal distribution of the unlabeled wild data:

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

where $\pi \in (0, 1]$, and \mathbb{P}_{out} is the OOD distribution defined over \mathcal{X} .

Learning Goal The learning framework aims to build the OOD detector g_{θ} and the multi-class classifier f_{θ} by leveraging data from both \mathcal{D}_{in}^{train} and \mathcal{D}_{wild} . Following [Du *et al.*, 2024], we here are interested in the following measurements for model evaluation:

$$\downarrow \operatorname{FPR}(g_{\theta}) := \mathbb{E}_{\boldsymbol{x} \sim \mathbb{P}_{\operatorname{out}}}(\mathbb{1}\{g_{\theta}(\boldsymbol{x}) = \operatorname{in}\}),$$

$$\uparrow \operatorname{TPR}(q_{\theta}) := \mathbb{E}_{\boldsymbol{x} \sim \mathbb{P}_{\operatorname{out}}}(\mathbb{1}\{q_{\theta}(\boldsymbol{x}) = \operatorname{in}\})$$

3.2 Loss-Difference OOD Detection Framework

To effectively address the two aforementioned potential issues, i.e., the model-bias issue and the threshold-selection dilemma, we innovatively propose a novel loss-difference OOD detection framework (abbreviated as LoD) by *intentionally label-noisifying* unlabeled wild data. As shown in Figure 2, LoD follows the two-step strategy and contains two main modules, i.e., loss-difference OOD filtering module and OOD detector learning module. Next, we will elaborate on the specific details of each module.

Loss-difference OOD Filtering Module

In this part, a loss-difference filtering mechanism with *intentional label-noises* is developed, which ingeniously reformulates the OOD filtering problem in unlabeled wild data as a label-noise learning problem with *controllable label-noise ratio*. This allows us to leverage the inherent properties of label-noise learning demonstrated in Section 2.2 to effectively filter OOD data from the unlabeled wild data.

In specific, we first intentionally label the whole unlabeled wild data as a single K + 1-th class (assuming that ID data contains K classes) and then train them together with labeled ID data in a *fully-supervised manner of* K + 1 *classification*, as follows:

$$\mathcal{L} = \frac{1}{|\mathcal{B}_{\text{in}}^{\text{train}}|} \sum_{(\boldsymbol{x}_{i}, y_{i}) \sim \mathcal{B}_{\text{in}}^{\text{train}}} \ell(\hat{y}_{i}, y_{i}) + \frac{1}{|\mathcal{B}_{\text{wild}}|} \sum_{(\boldsymbol{x}_{i}, y_{i}) \sim \mathcal{B}_{\text{wild}}} \ell(\hat{y}_{i}, y_{K+1}), \quad \hat{y}_{i} = f(\boldsymbol{x}_{i}, \boldsymbol{\theta}),$$

$$(2)$$

where $f(\cdot, \theta) \in \mathcal{F}$ denotes the K + 1 classifier, $l(\cdot, \cdot)$ represents the vanilla cross-entropy (CE) loss. Each training batch consists of two parts: \mathcal{B}_{in}^{train} and \mathcal{B}_{wild} , respectively sampled from labeled ID data and unlabeled wild data. Note that the ratio of $|\mathcal{B}_{in}^{train}| : |\mathcal{B}_{wild}| \ge 1$ is controllable. In fact, we indirectly control the label-noise ratio of the learning task by controlling this ratio (for more details, please refer to Section 4).

By labeling the entire unlabeled wild data as a single K+1th class, the ID samples in \mathcal{D}_{wild} are intentionally converted to *label-noise* samples while the OOD samples in \mathcal{D}_{wild} become label-clean ones. According to the inherent phenomenon of early learning stage in label-noise learning, the loss curves of these two types of labeled samples will exhibit significant difference during the early learning stage, as shown in Figure 1. This discrepancy provides us a critical clue for effectively distinguishing between them. Therefore, after training the K+1classifier, we conduct clustering operations on the loss values of unlabeled wild data gained during training so as to filter the OOD samples from \mathcal{D}_{wild} , which does not need the filtering thresholds any more. Particularly, clustering in our case has a well-defined number of clusters - Two - corresponding to the ID and OOD clusters represented by their distinct loss behaviors throughout the training process, e.g., higher lossvalues for ID data while lower counterparts for OOD data.

Considering the efficiency issue, we here utilize the mean of loss-values during training as the new features for each



Figure 3: The mean cross-entropy loss curves respectively for all ID (*label-noise*) and OOD samples (*label-clean*) in unlabeled wild data when they are intentionally labeled as K + 1-th class.

sample in $\mathcal{D}_{\text{wild}}$. Then the classic K-means clustering technique is employed to achieve the OOD samples filtering from unlabeled wild data. Let μ_1 , μ_2 ($\mu_1 > \mu_2$) respectively denote the ID and OOD cluster centers, while d_1 , d_2 respectively denote the distances between the corresponding sample and the two cluster centers. We filter the OOD samples from the unlabeled wild data by the following rule:

$$\hat{y} = \begin{cases} \text{ID data, if } d_1 < d_2, \\ \text{OOD data, otherwise.} \end{cases}$$
(3)

Remark. At first glance, labeling the entire set of unlabeled wild data as a single K + 1-th class seems potentially to undermine the model learning. Intriguingly, however, once we switch to consider filtering the OOD samples from the label-noise perspective, such operations, just on the contrary, bring at least the following three-fold advantages:

- **First**, OOD samples in unlabeled wild data are correctly labeled (*label-clean*), naturally and seamlessly enabling them to jointly dominate the model learning with labeled ID data, thus effectively circumventing the model-bias issue.
- Second, as mentioned earlier, the ID samples in \mathcal{D}_{wild} being erroneously labeled (*label-noise*) as K + 1-th class contradict the label-correct ones in labeled ID data \mathcal{D}_{in}^{train} , thereby resulting in their loss curves exhibiting both higher values and greater fluctuations compared to those of OOD samples, such as their mean-loss curves shown in Figure 3. This discrepancy provides us a fairly clear signal to distinguish ID and OOD samples in the unlabeled wild data.
- Third, our LoD is data-centric in nature, wherein we just relabel the unlabeled wild data as the intentionally K + 1-th class without any modifications to the network architectures we employed. This endows our LoD stronger applicability (for related experiments, please refer to Appendix E).

To solidly demonstrate the viability of this OOD filtering module, we also provide the theoretical analyses to support our first two claims, which will be detailed in Section 4.

OOD Detector Learning Module

After obtaining the candidate OOD samples \mathcal{D}_{out} from the unlabeled wild data, we training an OOD detector g_{θ} using

them together with labeled ID data \mathcal{D}_{in}^{train} . Similar to [Du *et al.*, 2024], we adopt the following optimization objective:

$$\mathcal{L}(g_{\theta}) = \mathbb{E}_{\boldsymbol{x} \in \mathcal{D}_{\text{in}}^{\text{train}}} \mathbb{1}\{g_{\theta}(\boldsymbol{x}) \le 0\} + \mathbb{E}_{\widetilde{\boldsymbol{x}} \in \mathcal{D}_{\text{out}}} \mathbb{1}\{g_{\theta}(\widetilde{\boldsymbol{x}}) > 0\},$$
(4)

where the binary sigmoid loss is employed as the smooth approximation of the 0/1 loss to make it tractable. In addition, a *K*-class classifier f_{θ} is also trained using CE loss on labeled ID data along with g_{θ} to ensure the ID accuracy. Algorithm 1 denotes the entire workflow of our LoD.

Algorithm 1 LoD OOD Detection Framework

Input: In-distribution data \mathcal{D}_{in}^{train} , unlabeled wild data \mathcal{D}_{wild} , Max Epoch *T*, Batch size $|\mathcal{B}|$.

Output: OOD detector g_{θ} and classifier f_{θ} .

- 1: # Loss-difference OOD detection module
- 2: Initializing: Model parameters θ , \mathcal{D}_{wild} labeled as K + 1-th class, loss record matrix $\mathcal{V} = \{\} \in \mathbb{R}^{|\mathcal{D}_{wild}| \times T}$.
- 3: for epoch = 1 to T do
- 4: Batch $\mathcal{B} = \mathcal{B}_{in}^{\text{train}} \cup \mathcal{B}_{\text{wild}}$, where $\mathcal{B}_{in}^{\text{train}}$ samples from $\mathcal{D}_{in}^{\text{train}}$ and $\mathcal{B}_{\text{wild}}$ samples from $\mathcal{D}_{\text{wild}}$.
- 5: Update K + 1 classifier $f(\cdot, \theta)$ based on Eq.(2).
- 6: Record losses of unlabeled wild data. $\mathcal{V} \leftarrow \mathcal{V} \cup \{l_i \mid i \in (1, |\mathcal{B}_{wild}|)\}$
- 7: end for
- 8: Calculate the mean-loss set of wild data $\{u_i\} = \text{mean}(\mathcal{V})$, where $\{u_i\} \in \mathbb{R}^{|\mathcal{D}_{\text{wild}}| \times 1}$.
- 9: Cluster and detect candidate OOD samples set \mathcal{D}_{out} based on Eq.(3).
- 10: # OOD detector learning module
- 11: for epoch = 1 to T do
- 12: Batch $\mathcal{B} = \mathcal{B}_{in}^{train} \cup \mathcal{B}_{out}$, where \mathcal{B}_{in}^{train} samples from \mathcal{D}_{in}^{train} and \mathcal{B}_{out} samples from \mathcal{D}_{out} .
- 13: Update f_{θ} and g_{θ} based on Eq.(4).

4 Theoretical Analysis

4.1 Mitigation of The Model Bias

For the first claim in Subsection 3.2, we here provide a theoretical analysis at the gradient level to demonstrate that in our LoD framework, labeled ID data and OOD data in \mathcal{D}_{wild} can jointly dominate the model learning. Let N_1 denote the number of samples in \mathcal{B}_{in} , while N_2 and N_3 denote the number of samples respectively from IDs and OODs in \mathcal{B}_{wild} . Then Eq.(2) can be rewritten in the following form:

$$\mathcal{L} = \underbrace{\frac{1}{N_1} \sum_{(\boldsymbol{x}_i, y_i) \sim \mathcal{B}_{\text{in}}} \ell(\hat{y}_i, y_i) + \frac{1}{N_2} \sum_{(\boldsymbol{x}_i, y_i) \sim \mathcal{B}_{\text{wild}}} \ell(\hat{y}_i, y_{K+1})}_{\text{ID data}} + \underbrace{\frac{1}{N_3} \sum_{(\boldsymbol{x}_i, y_i) \sim \mathcal{B}_{\text{wild}}} \ell(\hat{y}_i, y_{K+1})}_{\text{OOD data}}.$$
(5)

Let $\nabla \mathcal{L}_{N_k} = \frac{1}{N_k} \sum_{i=1}^{N_k} \nabla l(\hat{y}_i, y_i), k = 1, 2, 3$, denote the gradient of the corresponding part with respect to the model parameters $\boldsymbol{\theta}$. For the OOD samples in \mathcal{B}_{wild} , evidently,

^{14:} end for

they are correctly labeled (label-clean), the model parameters therefore will be updated in the correct gradient direction.

As for ID samples, they consist of two parts: one part sampled from \mathcal{D}_{in}^{train} (*label-clean*), and the other part sampled from \mathcal{D}_{wild} (*label-noise*). Then the update of the model parameters $\boldsymbol{\theta}$ is as follows:

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t - \eta (\nabla \mathcal{L}_{N_1} + \nabla \mathcal{L}_{N_2}), \tag{6}$$

where t denotes the number of steps for model update, and η is the learning rate. According to Eq.(6), the update of θ is determined by $(\nabla \mathcal{L}_{N_1} + \nabla \mathcal{L}_{N_2})$. Since $|\mathcal{B}_{in}^{train}| > |\mathcal{B}_{wild}|$ and $|\mathcal{B}_{wild}| \ge N_2$, we have

$$|\mathcal{B}_{in}^{train}| > |\mathcal{B}_{wild}| \ge N_2.$$

This indicates that correctly labeled ID samples dominate the updating of model parameters, especially when $|\mathcal{B}_{in}^{train}| \gg N_2$. In summary, we have the labeled ID data \mathcal{D}_{in}^{train} and the OOD data in \mathcal{D}_{wild} that can jointly dominate the model learning, thus effectively addressing the model-bias issue.

4.2 Discriminability between ID and OOD CE Mean-Losses

As mentioned earlier, the key to our LoD lies in ingeniously transforming the OOD filtering problem into a label-noise learning problem with controllable label-noise ratio, which allows us to leverage the established theoretical foundation of label-noise learning [Liu *et al.*, 2020a; Yue and Jha, 2024] to ensure the feasibility of our LoD. The work [Liu *et al.*, 2020a] has shown that the phenomenon in early learning stage, when training with noisy labels, is intrinsic to high-dimensional classification tasks, even in the simplest setting, far from being a peculiar feature of deep neural networks. Therefore, for the second claim in Subsection 3.2, a theoretical analysis of loss gap between ID (label-noise) and OOD (label-clean) data in \mathcal{D}_{wild} is provided here using a similar setting in [Liu *et al.*, 2020a].

Considering a two class dataset that consists of n independent samples (x_i, y_i) drawn from a mixture of two Gaussians in \mathbb{R}^d as follows.

$$\begin{aligned} \boldsymbol{x} &\sim \mathcal{N}(+\boldsymbol{v}, \sigma^2 \boldsymbol{I}_{d \times d}), & \text{if } y = +1 \\ \boldsymbol{x} &\sim \mathcal{N}(-\boldsymbol{v}, \sigma^2 \boldsymbol{I}_{d \times d}), & \text{if } y = -1, \end{aligned}$$

where v is an arbitrary unit vector in \mathbb{R}^d and σ^2 is a small constant. Denote y as the true hidden label and \tilde{y} as the observed label. Assume that for any sample x_i ,

$$\widetilde{y} = \begin{cases} y_i, & \text{with probability } 1 - \Delta, \\ -y_i, & \text{with probability } \Delta, \end{cases}$$
(7)

where $\Delta \in (0, 1/2)$ is the label-noise ratio. Let us consider a linear classifier $f(\cdot, \theta)$ trained by gradient descent on CE loss:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^{2 \times d}} \mathcal{L}_{CE}(\boldsymbol{\theta}) := -\frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{2} y_i \log(f(\boldsymbol{x}_i, \boldsymbol{\theta})).$$
(8)

In order to correctly classify the true classes well (and not overfit to the noisy labels), the rows of θ should be correlated with the vector v. Let $\nabla \mathcal{L}_{CE}(\theta)$ denote the gradient of Eq.(8). According to [Liu *et al.*, 2020a], we have the following lemma.

Lemma 1 (Early-learning succeeds [Liu *et al.*, 2020a]). *De*note by $\{\boldsymbol{\theta}_t\}$ the iterates of gradient descent with step size η . For any $\Delta \in (0, 1/2)$, there exists a constant δ_{Δ} , depending only on Δ , such that if $\delta \leq \delta_{\Delta}$, then with high probability 1 - o(1), there exists a $T = \Omega(1/\eta)$ such that: for all t < T, we have $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_0\| \leq 1$ and

$$-\nabla \mathcal{L}_{CE}(\boldsymbol{\theta}_t)^T \boldsymbol{v} / \|\nabla \mathcal{L}_{CE}(\boldsymbol{\theta}_t)\| \ge 1/6.$$

Lemma 1 indicates that under the condition of label-noise ratio Δ , the model parameters θ update along the proper gradient direction during the early learning stage. This means, during this period, the loss curves of ID (*label-noise*) and OOD (*label-clean*) samples in \mathcal{D}_{wild} will have significantly different characteristics, with larger loss values and greater fluctuations for ID samples versus smaller loss values and smaller fluctuations for OOD ones. To theoretically analyze this, we have the following proposition.

Proposition 1. Let l_i denote the loss value of each sample in \mathcal{D}_{wild} , which is bounded by R. $\overline{l_{in}} = \frac{1}{|\mathcal{D}_{im}^{wild}|} \sum_{i \in \mathcal{D}_{in}^{wild}} l_i$ and $\overline{l_{out}} = \frac{1}{|\mathcal{D}_{out}^{wild}|} \sum_{i \in \mathcal{D}_{out}^{wild}} l_i$ respectively denote the mean losses of ID and OOD sets from unlabeled wild data \mathcal{D}_{wild} , and $n = |\mathcal{D}_{in}^{wild}| + |\mathcal{D}_{out}^{wild}|$. Under the Lemma 1, with high probability, we have

$$\overline{l_{in}} - \overline{l_{out}} \ge 1 - 2e^{-\boldsymbol{\theta}^T \boldsymbol{v} + \frac{1}{2} \|\boldsymbol{\theta}\|^2 \delta^2} - \mathcal{O}(\frac{R}{\sqrt{n}})$$

Proposition 1 demonstrates that the cross-entropy mean losses of ID and OOD samples in \mathcal{D}_{wild} are distinguishable, just as the two curves shown in Figure 3. The proof is provided in Appendix A of supplementary materials (https://github.com/ChuanxingGeng/LoD).

5 Experiments

5.1 Implementation Details

Our LoD (https://github.com/ChuanxingGeng/LoD) framework contains two main modules, i.e., loss-difference OOD filtering module and OOD detector learning module. For these two modules, we follow [Du et al., 2024; Katz-Samuels et al., 2022] and employ Wide ResNet [Zagoruyko, 2016] with 40 layers and widen factor of 2 as the backbone. Moreover, for the loss-difference OOD filtering module, we use stochastic gradient descent with a momentum of 0.9 as the optimizer, and set the initial learning rate to 0.01. We train for 100 epochs using cosine learning rate decay, a batch size of 128 in which $|\mathcal{B}_{in}^{train}|$: $|\mathcal{B}_{wild}|$ = 3 : 1 , and a dropout rate of 0.3. For the OOD detector learning module, similar to [Du et al., 2024], we load a pre-trained ID classifier and add an additional linear layer which utilize the penultimatelayer features of ID classifier for binary classification. The initial learning rate is set to 0.001, and the remaining training configurations are consistent with those of the former module. All experiments are conducted on a single NVIDIA RTX 3090 GPU.

Evaluation Metrics. Similar to [Du *et al.*, 2024; Katz-Samuels *et al.*, 2022], we adopt the following evaluation metrics: (1) the false positive rate (FPR95) of OOD examples when true positive rate of ID examples is at 95%, (2)

						OOD	Dataset						
Methods	SV	/HN	Pl	aces	LSU	N-Crop	LSUN	-Resize	Tex	tures	Ave	erage	ACC
	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	
					$\pi =$	0.1							
OE (ICLR'19)	1.57	99.63	60.24	83.43	3.83	99.26	0.93	99.79	27.89	93.35	18.89	95.09	71.65
Energy(w/OE) (NeurIPS'20)	1.47	99.68	54.67	86.09	2.52	99.44	2.68	99.50	37.26	91.26	19.72	95.19	73.46
WOODS (ICML'22)	0.12	99.96	29.58	90.60	0.11	99.96	0.07	99.96	9.12	96.65	7.80	97.43	75.22
SAL (ICLR'24)	0.07	99.95	3.53	99.06	0.06	99.94	0.02	99.95	5.73	98.65	1.88	99.51	73.71
LoD (Ours)	0	100	3.34	99.16	0	100	0	100	4.79	98.87	1.63	99.61	73.85
					$\pi =$	0.5							
OE (ICLR'19)	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) (NeurIPS'20)	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 (ICML'22)	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
SAL (ICLR'24)	0.02	99.98	1.27	99.62	0.04	99.96	0.01	99.99	5.64	99.16	1.40	99.74	73.77
LoD (Ours)	0	100	1.53	99.66	0	100	0	100	3.72	99.19	1.05	99. 77	74.32
					$\pi =$	0.9							
OE (ICLR'19)	0.84	99.36	19.78	96.29	1.64	99.57	0.51	99.75	12.74	94.95	7.10	97.98	72.02
Energy(w/OE) (NeurIPS'20)	0.97	99.64	17.52	96.53	1.36	99.73	0.94	99.59	14.01	95.73	6.96	98.24	73.62
WOODS (ICML'22)	0.05	99.98	11.34	95.83	0.07	99.99	0.03	99.99	6.72	98.73	3.64	98.90	73.86
SAL (ICLR'24)	0.03	99.99	2.79	99.89	0.05	99.99	0.01	99.99	5.88	99.53	1.75	99.88	74.01
LoD (Ours)	0	100	0.48	99.90	0	100	0	100	2.78	99.41	0.65	99.86	74.34

Table 1: Evaluation results of FPR95 \downarrow (%), AUROC \uparrow (%) and ACC \uparrow (%) on standard benchmarks. CIFAR100 is ID, and bold numbers highlight the best results.

	Dataset										
Methods	CIF	AR10	CIFA	AR+10	CIFA	AR+50	TinyIr	nageNet	Av	erage	ACC
	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	ĂUROC	FPR95	AUROC	
				$\pi =$	0.1						
OE (ICLR'19)	30.83	94.9	11.40	97.98	22.21	95.98	82.3	75.34	36.69	91.05	91.45
Energy(w/OE) (NeurIPS'20)	38.36	89.85	16.40	96.51	36.18	90.49	88.48	74.30	44.86	87.79	86.98
WOODS (ICML'22)	32.33	93.70	22.39	95.95	22.12	95.76	74.60	78.62	37.86	91.01	92.43
SAL (ICLR'24)	12.95	97.35	4.76	98.88	10.66	97.63	48.35	86.71	19.18	95.14	91.50
LoD (Ours)	2.56	99.40	1.50	99.62	1.96	99.39	47.61	91.55	13.41	97.49	91.44
				$\pi =$	0.5						
OE (ICLR'19)	13.77	97.68	4.08	99.09	9.80	98.27	76.13	80.62	25.95	93.92	91.82
Energy(w/OE) (NeurIPS'20)	9.16	97.70	3.70	98.98	10.01	97.43	75.93	83.58	24.70	94.42	87.91
WOODS (ICML'22)	17.89	96.64	12.50	97.69	12.68	97.68	70.60	81.42	28.42	93.36	92.53
SAL (ICLR'24)	12.76	97.38	4.84	98.87	10.86	97.60	48.17	86.77	19.16	95.16	91.39
LoD (Ours)	2.32	99.47	1.04	99.71	1.96	99.46	46.44	91.52	12.94	97.54	91.33
				$\pi =$	0.9						
OE (ICLR'19)	6.40	98.71	1.56	99.50	4.94	98.97	67.45	84.98	20.09	95.54	92.10
Energy(w/OE) (NeurIPS'20)	2.95	98.63	1.30	99.41	2.18	98.52	58.84	88.92	16.32	96.37	89.58
WOODS (ICML'22)	12.82	97.50	10.98	98.03	10.51	98.07	68.01	82.82	25.58	94.11	92.17
SAL (ICLR'24)	12.95	97.34	4.30	98.91	11.11	97.56	49.19	86.66	19.39	95.12	91.41
LoD (Ours)	2.19	99.45	1.04	99.77	1.90	99.45	45.24	91.80	12.59	97.62	91.50

Table 2: Evaluation results of FPR95 \downarrow (%), AUROC \uparrow (%) and ACC \uparrow (%) on hard benchmarks, and bold numbers highlight the best results...

Area Under the Receiver Operating Characteristic curve (AU-ROC), and (3) ID classification Accuracy (ACC).

To comprehensively evaluate our LoD framework, we conduct extensive experiments on both standard benchmarks and hard benchmarks (newly curated in this paper) detailed in the following subsections. Moreover, limited by space, we defer additional experiments in the supplementary materials, including results on CIFAR10 (Appendix C), results on unseen OOD datasets (Appendix D), and results on different network structures (Appendix E).

5.2 Experiments on Standard Benchmarks

Datasets. For standard benchmarks, we here follow [Du *et al.*, 2024; Katz-Samuels *et al.*, 2022], and choose CI-FAR100 as in-distribution (ID) datasets (\mathbb{P}_{in}). For the out-of-distribution (OOD) test datasets (\mathbb{P}_{out}), we use a diverse collection of natural image datasets including SVHN [Netzer *et al.*, 2011], Textures [Cimpoi *et al.*, 2014], Places [Zhou *et*

al., 2017], LSUN-Crop [Yu *et al.*, 2015] and LSUN-Resize [Yu *et al.*, 2015]. For the unlabeled wild data (\mathbb{P}_{wild}), we follow [Du *et al.*, 2024] and mix datasets by combining a subset of ID data with OOD data under different mixture proportions $\pi \in \{0.1, 0.5, 0.9\}$. Specifically, the ID dataset is split into two equal halves (25,000 images per half), with one half used to mix with an OOD dataset (e.g., SVHN) to create the unlabeled wild data (\mathbb{P}_{wild}).

Main Results. We mainly compare our LoD with 4 latest methods using unlabeled wild data including Outlier Exposure (OE) [Hendrycks *et al.*, 2018], energy-regularization learning (Energy) [Liu *et al.*, 2020b], WOODS [Katz-Samuels *et al.*, 2022], and SAL [Du *et al.*, 2024]. Table 1 presents a comprehensive comparison of different methods on standard benchmarks, highlighting the substantial advantages of our proposed LoD. Across all datasets and π values, our approach consistently delivers superior performance, achieving an FPR95 close to 0%, which is significantly lower

than the current SOTA baseline, SAL. Notably, on the most challenging Textures, our method outperforms SAL with substantial reductions in FPR95 by 0.94%, 1.92%, and 3.10% for $\pi = 0.1, 0.5, 0.9$, respectively. Moreover, while existing SOTA methods demonstrate strong performance in AUROC, our LoD achieves notable improvements even in this aspect. Importantly, our LoD maintains competitive in-distribution accuracy, matching or surpassing the performance of SOTA methods such as SAL and WOODS across various π values.

Experiments on Hard Benchmarks 5.3

Datasets. In the settings of standard benchmarks, the ID and OOD samples are sourced from different datasets with inherently distinct distributions, which actually indirectly reduces the difficulty of OOD detection. As shown in Table 1, many methods, including ours, have achieved exceptionally high performance. To further demonstrate the advantages of our LoD, we here curate more challenging benchmarks, called hard benchmarks. Different from standard benchmarks, the ID and OOD samples on hard benchmarks come from the same dataset with different classes.

In specific, taking CIFAR10 as an example, we first randomly select 6 classes as ID data and the remaining 4 classes as OOD data. Then, similar to the splitting protocol of standard benchmarks, the training set of 6 ID classes is divided into two halves (15,000 images per half). One half is used as labeled ID data, while the other half is mixed with the data from 4 OOD classes to create the unlabeled wild data. We here select CIFAR10, CIFAR+10, CIFAR+50, and TinyImageNet [Vaze *et al.*, 2022] to curate the hard OOD benchmarks, and more details can be found in Appendix B of supplementary materials.

Main Results. Since the four methods we compared do not conduct the experiments on these benchmarks, we reproduce the results according to the source codes provided by them. Table 2 reports the detailed results on hard benchmarks. Across all datasets and under various π values, our LoD achieves better FPR95 and AUROC performance compared to existing methods, indicating that its OOD detection has stronger generalization. Notably, compared to the SOTA baseline SAL [Du et al., 2024], our method reduces FPR95 by substantial margins of 5.77%, 6.22%, and 6.80% on average when $\pi = 0.1, 0.5, 0.9$, respectively. Especially on CIFAR10, where LoD outperforms SAL more than 10% in case of FPR95. In particular, on the most challenging Tiny-ImageNet, LoD consistently surpasses SAL by a large margin of 4.84%, 4.75%, and 5.14% in terms of AUROC when $\pi = 0.1, 0.5, 0.9$, respectively. Besides, our LoD also maintains competitive ID classification accuracy compared to the SOTA baseline, comprehensively demonstrating the effectiveness of our LoD.

5.4 Experiments on Different Ratios and Epochs

Results on Different Ratios of $|\mathcal{B}_{in}^{train}|/|\mathcal{B}_{wild}|$ According to Section 4.1, the larger the ratio $|\mathcal{B}_{in}^{train}|/|\mathcal{B}_{wild}|$, the more dominant the labeled ID data in \mathcal{D}_{in}^{train} and the OOD data in \mathcal{D}_{wild} are in model learning, thus leading to better model performance. To verify this, we conduct experiments in different ratios of $|\mathcal{B}_{in}^{train}|/|\mathcal{B}_{wild}|$. Figure 4 shows



Figure 4: Experiments in different rations $(|\mathcal{B}_{in}^{train}|/|\mathcal{B}_{wild}|)$ on standard benchmarks (dashed lines) and hard benchmarks (solid lines).

the results. As the ratio increases, the model performance consistently improves across all benchmarks, strongly supporting our claim. Considering computational efficiency, $|\mathcal{B}_{in}^{\text{train}}|/|\mathcal{B}_{\text{wild}}|$ is set to 3 : 1 in all of our experiments.



Figure 5: The impacts of training epochs on results respectively in standard and hard benchmarks.

Impact of Epoch in Early-learning Succeeds

As shown in Proposition 1, the early-learning succeeds is the key to our LoD. To clearly demonstrate the appropriate number of training epochs, we conduct the epoch experiments on standard benchmark (take Textures as an example) and hard benchmark (take CIFAR10 as an example) respectively. Figure 5 shows the results, and we can observe a steady performance improvement in our LoD from 100 to 500 training epochs. At first glance, this phenomenon seems inconsistent with the early-learning succeeds in the traditional labelnoise learning field, which is usually shorter. However, please note that in our work setting, the label-noise ratio is controlled within an appropriate range by controlling the ratio of $|\mathcal{B}_{in}^{train}|/|\mathcal{B}_{wild}|$, meaning that correctly labeled samples all along dominate the network's learning. This further verifies the operability of LoD due to the long period early-learning succeeds. Considering efficiency issues, the training epochs of all experiments in this paper are set to 100 epochs.

6 Conclusion

In this paper, we innovatively propose a loss-difference OOD detection framework by intentionally label-noisifying unlabeled wild data, which ingeniously transforms the OOD filtering problem in unlabeled wild data into a label-noise learning problem with controllable label-noise ratio. Importantly, LoD not only effectively addresses the model-bias issue commonly associated with existing methods, but also circumvents the threshold selection dilemma inherent in these approaches.

Acknowledgments

This research was supported in part by the National Natural Science Foundation of China (62106102, 62376126, 62272229), in part by the NSFC-Hong Kong joint collaboration research fund CRS_HKU703/24, in part by the Hong Kong Scholars Program under Grant XJ2023035, in part by the Fundamental Research Funds for the Central Universities under Grant NS2024058.

References

- [Abati et al., 2019] Davide Abati, Angelo Porrello, Simone Calderara, and Rita Cucchiara. Latent space autoregression for novelty detection. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 481–490, 2019.
- [Arpit et al., 2017] Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at memorization in deep networks. In *International conference on machine learning*, pages 233–242. PMLR, 2017.
- [Behpour *et al.*, 2024] Sima Behpour, Thang Long Doan, Xin Li, Wenbin He, Liang Gou, and Liu Ren. Gradorth: a simple yet efficient out-of-distribution detection with orthogonal projection of gradients. *Advances in Neural Information Processing Systems*, 36, 2024.
- [Bevandić *et al.*, 2018] Petra Bevandić, Ivan Krešo, Marin Oršić, and Siniša Šegvić. Discriminative out-ofdistribution detection for semantic segmentation. *arXiv preprint arXiv:1808.07703*, 2018.
- [Chen et al., 2021] Jiefeng Chen, Yixuan Li, Xi Wu, Yingyu Liang, and Somesh Jha. Atom: Robustifying out-ofdistribution detection using outlier mining. In Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part III 21, pages 430–445. Springer, 2021.
- [Cimpoi et al., 2014] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 3606–3613, 2014.
- [Deng *et al.*, 2009] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee, 2009.
- [Djurisic *et al.*, 2022] Andrija Djurisic, Nebojsa Bozanic, Arjun Ashok, and Rosanne Liu. Extremely simple activation shaping for out-of-distribution detection. *arXiv preprint arXiv*:2209.09858, 2022.
- [Du et al., 2022] Xuefeng Du, Xin Wang, Gabriel Gozum, and Yixuan Li. Unknown-aware object detection: Learning what you don't know from videos in the wild. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 13678–13688, 2022.

- [Du *et al.*, 2024] Xuefeng Du, Zhen Fang, Ilias Diakonikolas, and Yixuan Li. How does unlabeled data provably help out-of-distribution detection? *arXiv preprint arXiv:2402.03502*, 2024.
- [Fang et al., 2024] Zhen Fang, Yixuan Li, Feng Liu, Bo Han, and Jie Lu. On the learnability of out-of-distribution detection. Journal of Machine Learning Research, 25, 2024.
- [Forouzesh *et al.*, 2022] Mahsa Forouzesh, Hanie Sedghi, and Patrick Thiran. Leveraging unlabeled data to track memorization. *arXiv preprint arXiv:2212.04461*, 2022.
- [Guan *et al.*, 2018] Melody Guan, Varun Gulshan, Andrew Dai, and Geoffrey Hinton. Who said what: Modeling individual labelers improves classification. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- [Hendrycks and Gimpel, 2016] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-ofdistribution examples in neural networks. arXiv preprint arXiv:1610.02136, 2016.
- [Hendrycks *et al.*, 2018] Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. *arXiv preprint arXiv:1812.04606*, 2018.
- [Katz-Samuels et al., 2022] Julian Katz-Samuels, Julia B Nakhleh, Robert Nowak, and Yixuan Li. Training ood detectors in their natural habitats. In *International Confer*ence on Machine Learning, pages 10848–10865. PMLR, 2022.
- [Krizhevsky, 2009] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- [Lee *et al.*, 2017] Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. Training confidence-calibrated classifiers for detecting out-of-distribution samples. *arXiv preprint arXiv:1711.09325*, 2017.
- [Lee *et al.*, 2018] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. *Advances in neural information processing systems*, 31, 2018.
- [Li et al., 2023] Yifan Li, Hu Han, Shiguang Shan, and Xilin Chen. Disc: Learning from noisy labels via dynamic instance-specific selection and correction. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 24070–24079, 2023.
- [Li et al., 2024a] Huiru Li, Liangxiao Jiang, and Chaoqun Li. Certainty weighted voting-based noise correction for crowdsourcing. *Pattern Recognition*, 150:110325, 2024.
- [Li et al., 2024b] Tianqi Li, Guansong Pang, Xiao Bai, Wenjun Miao, and Jin Zheng. Learning transferable negative prompts for out-of-distribution detection. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 17584–17594, 2024.
- [Liang *et al.*, 2017] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017.

- [Lienen and Hüllermeier, 2024] Julian Lienen and Eyke Hüllermeier. Mitigating label noise through data ambiguation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 13799–13807, 2024.
- [Lin *et al.*, 2024] Yexiong Lin, Yu Yao, and Tongliang Liu. Learning the latent causal structure for modeling label noise. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [Liu et al., 2020a] Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda. Early-learning regularization prevents memorization of noisy labels. Advances in neural information processing systems, 33:20331–20342, 2020.
- [Liu *et al.*, 2020b] Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. *Advances in neural information processing systems*, 33:21464–21475, 2020.
- [Malinin and Gales, 2018] Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. *Advances in neural information processing systems*, 31, 2018.
- [Neal et al., 2018] Lawrence Neal, Matthew Olson, Xiaoli Fern, Weng-Keen Wong, and Fuxin Li. Open set learning with counterfactual images. In Proceedings of the European Conference on Computer Vision (ECCV), pages 613– 628, 2018.
- [Netzer *et al.*, 2011] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, page 4. Granada, 2011.
- [Nguyen et al., 2015] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 427–436, 2015.
- [Sharifi *et al.*, 2025] Sina Sharifi, Taha Entesari, Bardia Safaei, Vishal M Patel, and Mahyar Fazlyab. Gradient-regularized out-of-distribution detection. In *European Conference on Computer Vision*, pages 459–478. Springer, 2025.
- [Song et al., 2022] Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy labels with deep neural networks: A survey. *IEEE* transactions on neural networks and learning systems, 34(11):8135–8153, 2022.
- [Sun and Li, 2022] Yiyou Sun and Yixuan Li. Dice: Leveraging sparsification for out-of-distribution detection. In *European Conference on Computer Vision*, pages 691– 708. Springer, 2022.
- [Sun et al., 2021] Yiyou Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations. Advances in Neural Information Processing Systems, 34:144–157, 2021.

- [Sun *et al.*, 2022] Yiyou Sun, Yifei Ming, Xiaojin Zhu, and Yixuan Li. Out-of-distribution detection with deep nearest neighbors. In *International Conference on Machine Learning*, pages 20827–20840. PMLR, 2022.
- [Tack *et al.*, 2020] Jihoon Tack, Sangwoo Mo, Jongheon Jeong, and Jinwoo Shin. Csi: Novelty detection via contrastive learning on distributionally shifted instances. *Advances in neural information processing systems*, 33:11839–11852, 2020.
- [Vaze *et al.*, 2022] Sagar Vaze, Kai Han, Andrea Vedaldi, and Andrew Zisserman. Open-set recognition: A good closed-set classifier is all you need. *the International Conference on Learning Representations*, abs/2110.06207, 2022.
- [Wang et al., 2022] Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtuallogit matching. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 4921–4930, 2022.
- [Wang *et al.*, 2023a] Qizhou Wang, Zhen Fang, Yonggang Zhang, Feng Liu, Yixuan Li, and Bo Han. Learning to augment distributions for out-of-distribution detection. *Advances in neural information processing systems*, 36:73274–73286, 2023.
- [Wang *et al.*, 2023b] Qizhou Wang, Junjie Ye, Feng Liu, Quanyu Dai, Marcus Kalander, Tongliang Liu, Jianye Hao, and Bo Han. Out-of-distribution detection with implicit outlier transformation. *In International Conference on Learning Representations*, 2023.
- [Wei et al., 2022] Hongxin Wei, Renchunzi Xie, Hao Cheng, Lei Feng, Bo An, and Yixuan Li. Mitigating neural network overconfidence with logit normalization. In *International conference on machine learning*, pages 23631– 23644. PMLR, 2022.
- [Yang et al., 2024] Jingkang Yang, Kaiyang Zhou, Yixuan Li, and Ziwei Liu. Generalized out-of-distribution detection: A survey. *International Journal of Computer Vision*, 132(12):5635–5662, 2024.
- [Yu *et al.*, 2015] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.
- [Yuan *et al.*, 2024a] Shunjie Yuan, Xinghua Li, Yinbin Miao, Haiyan Zhang, Ximeng Liu, and Robert H Deng. Combating noisy labels by alleviating the memorization of dnns to noisy labels. *IEEE Transactions on Multime-dia*, 2024.
- [Yuan *et al.*, 2024b] Suqin Yuan, Lei Feng, and Tongliang Liu. Early stopping against label noise without validation data. In *The Twelfth International Conference on Learning Representations*, 2024.
- [Yue and Jha, 2024] Chang Yue and Niraj K Jha. Ctrl: Clustering training losses for label error detection. *IEEE Transactions on Artificial Intelligence*, 2024.

- [Zagoruyko, 2016] Sergey Zagoruyko. Wide residual networks. arXiv preprint arXiv:1605.07146, 2016.
- [Zheng et al., 2023] Haotian Zheng, Qizhou Wang, Zhen Fang, Xiaobo Xia, Feng Liu, Tongliang Liu, and Bo Han. Out-of-distribution detection learning with unreliable outof-distribution sources. Advances in Neural Information Processing Systems, 36:72110–72123, 2023.
- [Zhou et al., 2017] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *IEEE* transactions on pattern analysis and machine intelligence, 40(6):1452–1464, 2017.
- [Zhu *et al.*, 2023] Jianing Zhu, Yu Geng, Jiangchao Yao, Tongliang Liu, Gang Niu, Masashi Sugiyama, and Bo Han. Diversified outlier exposure for outof-distribution detection via informative extrapolation. *Advances in Neural Information Processing Systems*, 36:22702–22734, 2023.

A Detailed Proof

To prove Proposition 1, we first reintroduce Lemma 1 from [Liu *et al.*, 2020a] and Proposition 1 as follows:

Lemma 1 (Early-learning succeeds). Denote by $\{\theta_t\}$ the iterates of gradient descent with step size η . For any $\Delta \in (0, 1/2)$, there exists a constant δ_{Δ} , depending only on Δ , such that if $\delta \leq \delta_{\Delta}$, then with high probability 1 - o(1), there exists a $T = \Omega(1/\eta)$ such that: for all t < T, we have $\|\theta_t - \theta_0\| \leq 1$ and

$$-\nabla \mathcal{L}_{CE}(\boldsymbol{\theta}_t)^T \boldsymbol{v} / \| \nabla \mathcal{L}_{CE}(\boldsymbol{\theta}_t) \| \ge 1/6$$

Proposition 1. Let l_i denote the loss value of each sample in \mathcal{D}_{wild} , which is bounded by R. $\overline{l_{in}} = \frac{1}{|\mathcal{D}_{in}^{wild}|} \sum_{i \in \mathcal{D}_{in}^{wild}} l_i$ and $\overline{l_{out}} = \frac{1}{|\mathcal{D}_{out}^{wild}|} \sum_{i \in \mathcal{D}_{out}^{wild}} l_i$ respectively denote the mean losses of ID and OOD sets from unlabeled wild data \mathcal{D}_{wild} , and $n = |\mathcal{D}_{in}^{wild}| + |\mathcal{D}_{out}^{wild}|$. Under the Lemma 1, with high probability, we have

$$\overline{l_{in}} - \overline{l_{out}} \ge 1 - 2e^{-\theta^T v + \frac{1}{2} \|\theta\|^2 \delta^2} - \mathcal{O}(\frac{R}{\sqrt{n}}).$$

Proof. Lemma 1 indicates that under the condition of noise level Δ , the model parameters θ update along the proper gradient direction during the early learning stage. This means, during this period, the loss curves of ID (*label-noise*) and OOD (*label-clean*) samples in test-set will have significantly different characteristics, with larger loss values and greater fluctuations for ID samples versus smaller loss values and smaller fluctuations for OOD ones. Next, we analyze the mean loss gap between ID (label-noise) samples in \mathcal{D}_{in}^{wild} and OOD (label-clean) samples in \mathcal{D}_{out}^{wild} during this stage. Following [Yue and Jha, 2024], we adopt sigmoid function as the activation function for the network outputs. For each sample $(\boldsymbol{x_i}, y_i)$, we have

$$p(y_i = 1) = \operatorname{sig}(\theta^T x_i) = \frac{1}{1 + e^{-\theta^T x_i}},$$

$$p(y_i = -1) = 1 - p(y_i = 1).$$

Let $\boldsymbol{x} = \boldsymbol{v} + \boldsymbol{z}_i$, where $\boldsymbol{z}_i \sim \mathcal{N}(\boldsymbol{0}, \sigma^2 \boldsymbol{I}_{d \times d})$. For each sample $\boldsymbol{x}_i \in \mathcal{D}_{\text{out}}^{\text{wild}}$ (label-clean), we use log for its loss, and have

$$l_i(\boldsymbol{\theta}) = \log(1 + e^{-\boldsymbol{\theta}^T(\boldsymbol{v} + \boldsymbol{z}_i)}) \le e^{-\boldsymbol{\theta}^T(\boldsymbol{v} + \boldsymbol{z}_i)}.$$

Similarly, for each sample $x_j \in \mathcal{D}_{\text{in}}^{\text{wild}}$ (label-noise), we have

$$l_j(\boldsymbol{\theta}) = \log(1 + e^{\boldsymbol{\theta}^T(\boldsymbol{v} + \boldsymbol{z}_i)}) \ge 1 - e^{-\boldsymbol{\theta}^T(\boldsymbol{v} + \boldsymbol{z}_i)}.$$

Taking the expectation on the difference between OOD (label-clean) and ID (label-noise), we have

$$\mathbb{E}[l_i(\boldsymbol{\theta}) - l_i(\boldsymbol{\theta})] = \mathbb{E}[l_i(\boldsymbol{\theta})] - \mathbb{E}[l_j(\boldsymbol{\theta})] \ge 1 - 2 \cdot \mathbb{E}[e^{-\boldsymbol{\theta}^T(\boldsymbol{v} + \boldsymbol{z})}].$$

Note that the term $1-2 \cdot \mathbb{E}[e^{-\theta^T(\boldsymbol{v}+\boldsymbol{z})}]$ bounds the loss gap between OOD (label-clean) and know-class (label-noise) samples, and it is independent of the label type. Since

$$\mathbb{E}[e^{-\boldsymbol{\theta}^T(\boldsymbol{v}+\boldsymbol{z})}] = e^{-\boldsymbol{\theta}^T\boldsymbol{v}} \cdot \mathbb{E}[e^{-\boldsymbol{\theta}^T\boldsymbol{z}}] = e^{-\boldsymbol{\theta}^T\boldsymbol{v}} \cdot e^{\frac{1}{2}\|\boldsymbol{\theta}\|^2 \sigma^2}.$$
 (9)

Eq.(1) indicates that the smaller the σ or the projection θ has on v, the larger the expected loss gap. Interestingly, Lemma 1 ensures that we can obtain a good θ at least within T epochs. Define the mean losses of ID (label-noise) samples and OOD (label-clean) samples as follows:

$$\overline{l_{\text{in}}} = \frac{1}{|\mathcal{D}_{\text{in}}^{\text{wild}}|} \sum_{i \in \mathcal{D}_{\text{in}}^{\text{wild}}} l_i, \quad \overline{l_{\text{out}}} = \frac{1}{|\mathcal{D}_{\text{out}}^{\text{wild}}|} \sum_{i \in \mathcal{D}_{\text{out}}^{\text{wild}}} l_i.$$

By Hoeffding's Inequality on bounded variables and the Union Bound, with probability $\geq 1 - \delta$, we have

$$\overline{l_{\text{in}}} \ge \mathbb{E}[\overline{l_{\text{in}}}] - \mathcal{O}(\frac{R}{\sqrt{|\mathcal{D}_{\text{in}}^{\text{wild}}|}}\sqrt{\log\frac{1}{\delta}}).$$
(10)

and

$$\overline{l_{\text{out}}} \le \mathbb{E}[\overline{l_{\text{out}}}] + \mathcal{O}(\frac{R}{\sqrt{|\mathcal{D}_{\text{out}}^{\text{wild}}|}}\sqrt{\log\frac{1}{\delta}}).$$
(11)

According to Eq.(2) and Eq.(3), we have

$$\overline{l_{\mathrm{in}}} - \overline{l_{\mathrm{out}}} \ge 1 - 2e^{-\theta^T v + \frac{1}{2} \|\theta\|^2 \delta^2} - \mathcal{O}(\frac{R}{\sqrt{n}}).$$

B Details in Hard Benchmarks

To further demonstrate the advantages of our LoD, we conduct experiments on curated hard OOD benchmarks including CIFAR10, CIFAR+10, CIFAR+50, and TinyImageNet. The details of these benchmarks are as follows:

- **CIFAR10.** CIFAR10 [Krizhevsky, 2009] contains 10 classes, where 6 classes are randomly selected as indistribution (ID) classes, and the remaining 4 classes are used as out-of-distribution (OOD) classes.
- CIFAR+10 & CIFAR+50. In this set of experiments, 4 classes from CIFAR10 are randomly selected as ID classes, and 10/50 non-overlapping classes randomly selected from CIFAR100 [Krizhevsky, 2009] are OOD classes.

						OOD	Dataset						
Methods	SV	/HN	Pl	aces	LSUI	N-Crop	LSUN	I-Resize	Tex	tures	Av	erage	ACC
	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	
					With \mathbb{P}	in only							
MSP (ICLR'17)	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 (ICLR'18)	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 (NeurIPS'18)	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 (NeurIPS'20)	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
CSI (NeurIPS'20)	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
ReAct (NeurIPS'21)	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
KNN (ICML'22)	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
KNN+ (ICML'22)	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
DICE (ECCV'22)	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 (ICLR'23)	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
					With \mathbb{P}_{in}	and \mathbb{P}_{wild}							
OE (ICLR'19)	0.85	99.82	23.47	94.62	1.84	99.65	0.33	99.93	10.42	98.01	7.38	98.41	94.07
Energy(w/OE) (NeurIPS'20)	4.95	98.92	17.26	95.84	1.93	99.49	5.04	98.83	13.43	96.69	8.52	97.95	94.81
WOODS (ICML'22)	0.15	99.97	12.49	97.00	0.22	99.94	0.03	99.99	5.95	98.79	3.77	99.14	94.84
SAL (ICLR'24)	0.02	99.98	2.57	99.24	0.07	99.99	0.01	99.99	0.90	99.74	0.71	99.78	93.65
LoD (Ours)	0	100	1.72	99.52	0	100	0	100	0.66	99.90	0.48	99.88	94.06

Table 3: Evaluation results of FPR95 \downarrow (%), AUROC \uparrow (%) and ACC \uparrow (%) on standard benchmarks. CIFAR10 is ID dataset, and bold numbers highlight the best results.

							OOD	Dataset							
Methods	SVHN		Pl	aces	LSU	N-Crop	LSUN	-Resize	Tex	tures	25K RA	ND.IMG.	Average		ACC
	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	FPR95	AUROC	
						π=().1								
OE (ICLR'19)	77.74	82.84	60.70	84.02	31.06	93.99	55.74	88.45	57.39	88.27	50.95	87.44	56.53	87.51	83.04
Energy(w/OE) (NeurIPS'20)	55.89	90.19	49.08	88.03	22.74	94.94	34.10	93.42	39.33	90.63	48.91	88.12	40.23	91.44	90.02
WOODS (ICML'22)	4.90	98.70	18.53	96.27	1.94	99.53	5.73	98.78	17.71	96.17	10.37	96.92	9.76	97.89	94.50
SAL (ICLR'24)	5.83	97.63	17.96	96.23	2.50	98.77	5.67	98.56	8.44	97.93	8.95	97.40	8.08	97.82	93.65
LoD (Ours)	4.41	98.96	11.82	97.50	1.84	99.56	5.61	98.63	4.66	99.10	8.68	97.59	5.67	98.75	93.99

Table 4: Evaluation results of FPR95 \downarrow (%), AUROC \uparrow (%) and ACC \uparrow (%) on unseen datasets. We use CIFAR10 as ID and a subset (25K images) of 300K Random Images as wild OOD data. Bold numbers highlight the best results

• **TinyImageNet.** TinyImageNet is a subset derived from ImageNet [Deng *et al.*, 2009] with a total of 200 classes, of which 20 classes are randomly selected as ID classes and the rest 180 classes are treated as OOD classes.

Please note that, since ID and OOD are randomly divided, to mitigate the effects of randomness, each dataset is evaluated across five distinct "ID/OOD" splits following [Neal *et al.*, 2018; Vaze *et al.*, 2022], and the results are averaged. Moreover, similar to standard benchmarks [Katz-Samuels *et al.*, 2022; Du *et al.*, 2024], we use 70% of data from the OOD classes as the OOD part of the unlabeled wild data.

C Additional Results on CIFAR10

In this part, we utilize CIFAR10 as the ID dataset to evaluate our LoD under $\pi = 0.1$. In addition to the four methods utilizing wild data compared in the main paper, we here also evaluate methods that rely solely on labeled ID data (\mathbb{P}_{in} only) including MSP [Hendrycks and Gimpel, 2016], ODIN [Liang *et al.*, 2017], Mahalanobis [Lee *et al.*, 2018], Energy [Liu *et al.*, 2020b], CSI [Tack *et al.*, 2020], ReAct [Sun *et al.*, 2021], KNN and KNN+ [Sun *et al.*, 2022], DICE [Sun and Li, 2022] and ASH [Djurisic *et al.*, 2022]. The detailed results are presented in Table 3, which demonstrate that methods trained using both ID and wild data exhibit significantly better performance compared to those trained solely with ID data. Additionally, compared with methods utilizing \mathbb{P}_{wild} , LoD continues to exhibit superior performance, outperforming other SOTA methods in terms of FPR95 and AUROC metrics. Furthermore, LoD achieves competitive ID classification accuracy, either matching or exceeding the performance of leading SOTA methods such as SAL and WOODS.

D Additional Results on Unseen OOD Datasets

In this part, we follow [Du *et al.*, 2024] and evaluate our LoD on unseen OOD datasets, which are different from the OOD data we use in the wild. Table 2 and Table 3 report the results.

In Table 2, we employ CIFAR10 as the ID dataset. As for the wild OOD data, [Du *et al.*, 2024] utilizes the full 300Kimage dataset. However, we argue that this setting seems inappropriate due to a significant imbalance: the ID data in the wild data contains only 25K images, while the OOD counterpart comprises 300K images–12 times larger than the ID data. Therefore, we randomly sample a subset of 25K images from the 300K as the wild OOD data. The detailed results presented in Table 4 demonstrate that our LoD consistently outperforms SOTA baselines such as SAL and WOODS on the unseen OOD datasets, highlighting the effectiveness of our method.

In Table 3, we employ CIFAR100 as ID data. As for the wild OOD data, we follow [Du *et al.*, 2024] and utilize TinyImageNet-crop (TINc)/TinyImageNet-resize (TINr) dataset as the wild OOD data using during training and TINr/TINc as the unseen OOD data during testing. The results in Table 5 demonstrate the advantages of our LoD.

		OOD Dataset								
Methods	Т	INr	TINc							
	FPR95	AUROC	FPR95	AUROC						
STEP (NeurIPS'21)	72.31	74.59	48.68	91.14						
TSL (MM'23)	57.52	82.29	29.48	94.62						
SAL (ICLR'24)	43.11	89.17	19.30	96.29						
LoD (Ours)	23.54	92.81	9.67	98.10						

Table 5: Evaluation results of FPR95 \downarrow (%), AUROC \uparrow (%) on unseen datasets. CIFAR100 is ID, and bold numbers highlight the best results.

E Additional Results on Different Networks

To verify the applicability of LoD, the data-centric method, this part conducts experiments on different network structures on CIFAR10 and CIFAR+10. Table 4 reports the results, and we can find these networks mentioned here are all suitable for our LoD. In particular, LoD seems to follow scaling laws: the larger the network, the better it performs.

Naturalizad		CIFAR10	CIFAR+10				
INCLWOIKS(# params)	FPR95	AUROC	ACC	FPR95	AUROC	ACC	
WideResNet-40-2 (2.2M)	2.56	99.40	96.34	1.50	99.62	97.29	
ResNet18 (11.2M)	2.47	99.44	96.46	0.96	99.72	97.32	
ResNet34 (21.3M)	2.29	99.51	96.48	0.90	99.73	97.39	

Table 6: Evaluation results of FPR95 \downarrow (%), AUROC \uparrow (%) and ACC \uparrow (%) on different networks on hard benchmarks.

		Places		Textures				
Ratios	FPR95	AUROC	ACC	FPR95	AUROC	ACC		
1:6	10.50	98.07	72.9	8.88	97.64	74.10		
1:3	8.21	98.36	73.14	8.09	98.00	73.58		
1:1	4.08	99.12	73.06	6.17	98.53	74.06		
3:1	3.34	99.16	72.21	4.79	98.87	73.30		
6:1	3.91	98.95	71.38	3.63	99.14	73.22		

Table 7: Detailed results of FPR95 \downarrow (%), AUROC \uparrow (%) and ACC \uparrow (%) across different ratios on standard benchmarks.

F Detailed Results on Different Ratios $|\mathcal{B}_{in}^{train}|/|\mathcal{B}_{wild}|$

Table 7 reports the detailed results in different ratios on representative standard benchmark Places and Textures, while Table 8 reports the detailed results on CIFAR10 and Tiny-ImageNet. As the ratio $|\mathcal{B}_{in}^{train}|/|\mathcal{B}_{wild}|$ increasing, the performance of our method consistently improves.

Derter		CIFAR10		CIFAR+10					
Ratios	FPR95	AUROC	ACC	FPR95	AUROC	ACC			
1:6	3.07	99.30	96.18	8.30	97.66	97.35			
1:3	2.78	99.34	96.27	7.04	98.22	97.25			
1:1	2.57	99.38	96.28	2.52	99.31	97.25			
3:1	2.56	99.40	96.34	1.50	99.62	97.29			
6:1	2.33	99.44	96.21	1.13	99.73	97.34			

Table 8: Detailed results of FPR95 \downarrow (%), AUROC \uparrow (%) and ACC \uparrow (%) across different ratios on hard benchmarks.

Dation		Textures		CIFAR10				
Katios	FPR95	AUROC	ACC	FPR95	AUROC	ACC		
100	4.79	98.87	73.30	2.56	99.40	96.34		
200	4.17	99.07	72.86	2.37	99.45	96.28		
300	4.02	99.05	72.14	2.36	99.45	96.28		
400	3.82	99.05	72.19	2.33	99.46	96.30		
500	3.60	99.09	72.28	2.31	99.56	96.30		

Table 9: Detailed results of FPR95 \downarrow (%), AUROC \uparrow (%) and ACC \uparrow (%) in different training epochs.

G Detailed Results on the Impact of Epoch in Early-learning Succeeds

Table 9 reports the detailed results on the impact of training epochs in early-learning success. The results demonstrate a consistent performance improvement in our LoD model as the number of training epochs increases from 100 to 500.