

---

# LEMON: Label Error Detection using Multimodal Neighbors

---

Haoran Zhang<sup>\*1</sup> Aparna Balagopalan<sup>\*1</sup> Nassim Oufattole<sup>1</sup> Hyewon Jeong<sup>1</sup> Yan Wu<sup>1</sup> Jiacheng Zhu<sup>1</sup>  
Marzyeh Ghassemi<sup>1</sup>

## Abstract

Large repositories of image-caption pairs are essential for the development of vision-language models. However, these datasets are often extracted from noisy data scraped from the web, and contain many mislabeled instances. In order to improve the reliability of downstream models, it is important to identify and filter images with incorrect captions. However, beyond filtering based on image-caption embedding similarity, no prior works have proposed other methods to filter noisy multimodal data, or concretely assessed the impact of noisy captioning data on downstream training. In this work, we propose, theoretically justify, and empirically validate LEMON, a method to identify label errors in image-caption datasets. Our method leverages the multimodal neighborhood of image-caption pairs in the latent space of contrastively pretrained multimodal models to automatically identify label errors. Through empirical evaluations across eight datasets and twelve baselines, we find that LEMON outperforms the baselines by over 3% in label error detection, and that training on datasets filtered using our method improves downstream captioning performance by more than 2 BLEU points over noisy training.

## 1. Introduction

Machine learning datasets used to train and finetune large vision, language, and vision-language models frequently contain millions of labeled instances (Schuhmann et al., 2021; Li et al., 2022; Wang et al., 2022a; Changpinyo et al., 2021). Prior work highlights that some instances in such datasets may be mislabeled (Northcutt et al., 2021b; Lucioni & Rolnick, 2023; Liao et al., 2021; Beyer et al., 2020; Plummer et al., 2015), as seen in Figure 1. This is especially

---

<sup>\*</sup>Equal contribution <sup>1</sup>Massachusetts Institute of Technology. Correspondence to: Haoran Zhang <haoranz@mit.edu>, Aparna Balagopalan <aparnab@mit.edu>.



Figure 1: Samples from classification and captioning datasets discovered to be mislabeled by our method.

problematic in settings such as healthcare, where the reliability of downstream models may depend on the quality of data used for pretraining (Chen et al., 2024; Liu et al., 2023; Longpre et al., 2023).

Identifying and correcting label errors in existing datasets at scale would lead to more reliable and accurate models in the real world (Zhu et al., 2022; Vasudevan et al., 2022; Liao et al., 2021; Beyer et al., 2020). However, given the large size of such datasets, manual detection of errors is practically infeasible. This is evidenced by the growth of models trained on noisy data with the web (Li et al., 2022; Wang et al., 2022a; Liu et al., 2024), or with model generated pseudo-labels (Menghini et al., 2023; Lai et al., 2023).

Machine learning (ML) based approaches to automatically identifying label errors have also been proposed in prior work (Pleiss et al., 2020; Swayamdipta et al., 2020; Liang et al., 2023; Bahri et al., 2020; Zhu et al., 2022; Northcutt et al., 2021a). However, we identify two critical limitations: (1) the majority of such works are *unimodal*: i.e., they only utilize image-based representations in detection strategies, and (2) many of the best-performing approaches depend on having access to a model already trained on the downstream tasks of interest (Pleiss et al., 2020; Swayamdipta et al., 2020). We hypothesize that applying a neighborhood-based approach to multimodal representations in the form of image-text pairs can improve label error detection without requiring task-specific training, which may be costly and/or domain specific for some datasets.

Additionally, a common assumption made in prior works is that each label is one-of-k classes (Bahri et al., 2020; Zhu et al., 2022). The vast majority of label error detection meth-

ods proposed in prior works are hence for *classification* datasets. In contrast, datasets used to train large vision-language models contain natural language labels such as image captions (Li et al., 2022; 2023; Wang et al., 2022a). Methods to filter out instances with noisy labels — e.g., based on the similarity of image and caption representations — have been utilized in prior work with some success (Li et al., 2022; Kang et al., 2023) for such datasets. However, to the best of our knowledge, no prior works have proposed or rigorously compared methods to identify errors in datasets with natural language labels, or assessed the impact of detection on downstream tasks like image captioning.

In this work, we propose LEMON – Label Error detection using Multimodal Neighbors, a method for multimodal label error detection which can be applied to image-text pairs in datasets such as MSCOCO (Lin et al., 2014). While prior techniques utilize unimodal neighbors for label error detection, LEMON leverages multi-modal neighborhoods derived using contrastively pretrained vision-language models such as Contrastive Language-Image Pre-training (CLIP) (Radford et al., 2021). Specifically, in addition to considering pairwise image-text distances, we also retrieve nearest neighbors in the image and text space as illustrated in Figure 2. This is motivated by the rich neighborhood geometry in the joint embedding space of multimodal models (Liang et al., 2022; Schrodi et al., 2024). We then compute distance scores with neighbors in each modality and combine these into a single score measuring the likelihood of a label error, with the intuition that higher discordance (or higher distance) with neighbors indicates a higher chance of label error. We validate the utility of these scores across eight datasets, including one in a healthcare setting, and compare to over ten baselines.

Our key contributions and findings are as follows<sup>1</sup>:

- We propose LEMON, a novel, theoretically justified multimodal method capable of detecting label errors in large image-caption datasets (Section 3).
- We show that LEMON outperforms all training-free baselines for label error detection in three out of four classification datasets by up to 3.4% AUROC, and in three out of four captioning datasets by up to 6.3% AUROC (Section 6.1).
- We demonstrate that LEMON improves performance on downstream classification and captioning models by filtering out data predicted to be label errors. (Section 6.2).
- Finally, we verify that the predictions generated by LEMON are meaningful through a real world analysis of LEMON on existing datasets without known label errors (Section 6.5).

<sup>1</sup>Code: <https://github.com/MLforHealth/LEMON>

## 2. Related Works

**Label Noise Detection** Noisy and incorrect labels (Beyer et al., 2020) in training data may lead to decreased or “destabilized” (Northcutt et al., 2021a; Luccioni & Rolnick, 2023) performance on downstream tasks (Chen et al., 2023; Northcutt et al., 2021b). Two orthogonal approaches can be taken to reduce the adverse effects of such labels: developing methods to learn in the presence of label errors (Cui et al., 2020; Natarajan et al., 2013; Huang et al., 2023) referred to as “noise robust” training, and/or detecting and filtering out instances with label errors (Zhu et al., 2024). In this work, we focus on the latter direction. This approach may be preferable, as identifying mislabeled samples is more flexible, with applications beyond just removing these samples for downstream model training. For example, accurate label error detection can uncover systematic errors or biases in datasets (Rottmann & Reese, 2023), and these insights can then be used to guide higher-quality data collection practices (Bernhardt et al., 2022). This is especially important for practitioners releasing datasets intended for model evaluation (Northcutt et al., 2021b; Schubert et al., 2024).

Prior approaches (Swayamdipta et al., 2020; Bahri et al., 2020; Pleiss et al., 2020; Northcutt et al., 2021a; Liang et al., 2023; Wu et al., 2020; Kim et al., 2021) for automatic label error detection include relying on the training dynamics of task-specific downstream models (Swayamdipta et al., 2020) and neighborhood-based strategies (Bahri et al., 2020; Grivas et al., 2020). Some of these techniques are fully supervised (Northcutt et al., 2021a; Chen et al., 2023) or unsupervised (Pleiss et al., 2020; Swayamdipta et al., 2020; Grivas et al., 2020; Bahri et al., 2020), use pre-trained generative models (Gertz et al., 2024) or are fully training-free approaches (Zhu et al., 2022; Liang et al., 2023). Previous approaches for label error detection closest to this work includes deep k-nearest neighbor (deep k-NN) methods using k-NN entropy on vector space embeddings (Bahri et al., 2020; Grivas et al., 2020) and SimiFeat (Zhu et al., 2022) which employs a local neighborhood-based voting or ranking for noise identification. In contrast to these methods, our work enhances label noise detection by harnessing information across *multiple data modalities*, such as image and text.

**Contrastive Learning** Contrastive learning is a representation learning method, relying on positive and negative pairs of data instances (Chen et al., 2020; Misra & Maaten, 2020; Balestriero et al., 2023) to learn an embedding space. The core idea is to embed similar data points (positive pairs) closer together than dissimilar data points (negative pairs) (Schroff et al., 2015; Sohn, 2016; Oord et al., 2018). In this work, we primarily utilize pre-trained models that use the CLIP loss (where the pre-training objective is to predict which text caption is paired with which image)

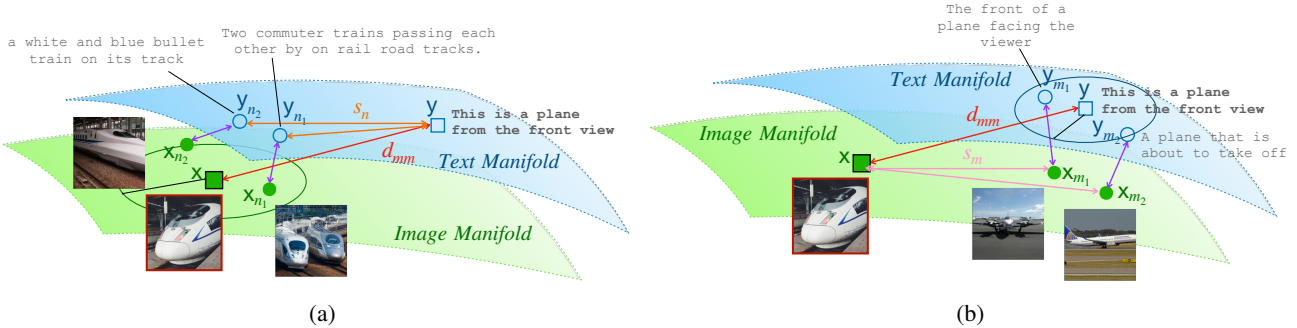


Figure 2: **Outline of LEMON**, our proposed method for multimodal label error detection. We demonstrate LEMON on a real sample from the MSCOCO dataset, where an image of a train ( $x$ ) is mislabeled as  $y$  = “This is a plane from the front view”. (a) We compute the simple CLIP similarity  $d_{mm}(x, y)$ . We then find the nearest neighbors of  $x$  in the image space ( $x_{n_j}$ ) and compute the distance between the corresponding texts and  $y$  to compute the score component  $s_n$ . (b) To compute the score component  $s_m$ , we find the nearest neighbors of  $y$  in the text space ( $y_{m_k}$ ), and compute the distance between the corresponding images and  $x$ .

for jointly embedding image and text data (Radford et al., 2021).

**Image Captioning** The goal of image captioning is to describe a given image (Fu et al., 2024) in natural language. Prior approaches for caption generation have included supervised training of end-to-end models from scratch (Wang et al., 2022b; Lin et al., 2022; Hu et al., 2023; Xu et al., 2015; Fu et al., 2024). More recently, vision-language models pretrained on large datasets of noisy image-caption pairs extracted from the web (Li et al., 2022; 2023; Wang et al., 2022a) – such as CC12M (Changpinyo et al., 2021) – have been utilized for captioning. Some of the pretraining tasks include image-text contrastive learning, image-text matching, and/or retrieval (Li et al., 2022), as well as general purpose text generation conditioned on an input image (Wang et al., 2022a). Given that datasets for training such large models are noisy (Kang et al., 2023), several steps have been utilized in prior work to filter out noisy captions during training. The most common strategy involves computing the similarity between representations of the image and caption text using another pretrained model (e.g., CLIP) prior to training (Kang et al., 2023). Another approach in training the BLIP (Li et al., 2022) model is to synthetically generate noisy captions and train a classifier to distinguish between high quality captions and noisy captions with a cross-entropy loss (Li et al., 2022). To the best of our knowledge, no previous work has conducted a comprehensive comparison of various strategies for label error detection in captioning datasets.

**Multimodal Neighborhood Methods** Previous studies (Li et al., 2021; Thomas & Kovashka, 2020; Huang et al., 2024; Liang et al., 2022; Cai et al., 2023) have examined the geometry of neighborhood spaces in mul-

timodal models, often with the goal of improving representation learning (Huang et al., 2024; Li et al., 2021) or retrieval (Thomas & Kovashka, 2020; 2022). The closest related work is Thomas & Kovashka (2022), where the authors use the semantic neighborhood of multimodal models to identify samples with high semantic diversity using text-based neighbors of neighbors. Importantly, the objective of their work is different from ours, which leads their proposed discrepancy and diversity scores to lack signal for label error in our setting. We further clarify this in Appendix B, and empirically compare against their discrepancy score as a baseline. We believe our work is the first to use multimodal neighbors for label error detection.

### 3. LEMON: Label Error Detection using Multimodal Neighbors

We are given a dataset  $\mathcal{D} = \{(x, y)_{i=1}^N\}$  consisting of two modalities  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ . For example,  $\mathcal{X}$  may represent the set of all natural images, and  $\mathcal{Y}$  may represent the set of all English text, or a restricted subset such as {cat, dog, ...}. We assume the existence of, but not access to, an oracle  $f^* : \mathcal{X} \times \mathcal{Y} \rightarrow \{0, 1\}$ , which is able to assign a binary mislabel indicator  $z_i = f^*(x_i, y_i)$  to each sample in  $\mathcal{D}$ . Here,  $z_i = 1$  indicates that the sample is mislabeled, and  $z_i = 0$  indicates that the sample is correctly labeled. Our goal is to output a score  $s \in \mathbb{R}$  with some model  $s := f(x, y)$  such that:

$$\text{AUROC} = \mathbb{E}_{(x, y) \sim \mathbb{P}(\cdot | z=1), (x', y') \sim \mathbb{P}(\cdot | z=0)} [\mathbf{1}_{f(x, y) \geq f(x', y')}]$$

is maximized. Prior works have alternatively aimed to maximize the F1 score, optimizing over a threshold  $t$ :

$$\text{F1} = \max_{t \in \mathbb{R}} \frac{2 \cdot \mathbb{P}(z = 1 | s \geq t) \cdot \mathbb{P}(s \geq t | z = 1)}{\mathbb{P}(z = 1 | s \geq t) + \mathbb{P}(s \geq t | z = 1)}$$

Here, we build on prior work for label error detection in unimodal data (Bahri et al., 2020; Zhu et al., 2022) and propose a method for  $f$  based on nearest neighbors, summarized in Figure 2. Suppose we have a query sample  $(\mathbf{x}, \mathbf{y})^2$ . Define  $B(\mathbf{x}, r) := \{\mathbf{x}' \in \mathcal{X} : d_{\mathcal{X}}(\mathbf{x}, \mathbf{x}') \leq r\}$ , the ball of radius  $r$  around  $\mathbf{x}$ , and  $B(\mathbf{y}, r)$  similarly. Let  $r_k(\mathbf{x}) := \inf\{r : |B(\mathbf{x}, r) \cap \mathcal{D}| \geq k\}$ , the minimum radius required to encompass at least  $k$  neighbors. Then, we define  $\{\mathbf{x}_{n_1}, \mathbf{x}_{n_2}, \dots, \mathbf{x}_{n_k}\} := B(\mathbf{x}, r_k(\mathbf{x})) \cap \mathcal{D}$  the top  $k$  nearest neighbors of  $\mathbf{x}$ , and  $\{\mathbf{y}_{m_1}, \mathbf{y}_{m_2}, \dots, \mathbf{y}_{m_k}\} := B(\mathbf{y}, r_k(\mathbf{y})) \cap \mathcal{D}$  the top  $k$  nearest neighbors of  $\mathbf{y}$ <sup>3</sup>. We assume that the neighbors are sorted in order of ascending distance, e.g.  $d_{\mathcal{X}}(\mathbf{x}, \mathbf{x}_{n_2}) \geq d_{\mathcal{X}}(\mathbf{x}, \mathbf{x}_{n_1})$ .

If  $\mathcal{Y}$  is a small discrete set, we could choose  $d(\mathbf{y}, \mathbf{y}') = \mathbf{1}_{\mathbf{y} \neq \mathbf{y}'}$ . If  $\mathcal{X}$  or  $\mathcal{Y}$  are unstructured or high dimensional, we assume access to multimodal encoders  $h_{\theta} = (h_{\theta}^{\mathcal{X}}, h_{\theta}^{\mathcal{Y}})$ , where  $h_{\theta}^{\mathcal{X}} : \mathcal{X} \rightarrow \mathbb{R}^d$  and  $h_{\theta}^{\mathcal{Y}} : \mathcal{Y} \rightarrow \mathbb{R}^d$ . Here,  $h_{\theta}$  may be a CLIP model (Radford et al., 2021) trained on a large internet corpus, or, as we show later, it may be sufficient to train  $h_{\theta}$  from scratch only on  $\mathcal{D}$ . Then, we could naturally use simple distance metrics in the embedding space, such as the cosine distance  $d_{\mathcal{X}}(\mathbf{x}, \mathbf{x}') = d_{\cos}(h_{\theta}^{\mathcal{X}}(\mathbf{x}), h_{\theta}^{\mathcal{X}}(\mathbf{x}')) = 1 - \frac{h_{\theta}^{\mathcal{X}}(\mathbf{x})^T h_{\theta}^{\mathcal{X}}(\mathbf{x}')}{\|h_{\theta}^{\mathcal{X}}(\mathbf{x})\|_2 \cdot \|h_{\theta}^{\mathcal{X}}(\mathbf{x}')\|_2}$ . Our proposed score is the linear combination of three terms:

$$s = f(\mathbf{x}, \mathbf{y}) = d_{mm}(\mathbf{x}, \mathbf{y}) + \beta s_n(\mathbf{x}, \mathbf{y}, \mathcal{D}) + \gamma s_m(\mathbf{x}, \mathbf{y}, \mathcal{D}), \quad (1)$$

where  $\beta, \gamma \geq 0$  are hyperparameters. Here,  $d_{mm}(\mathbf{x}, \mathbf{y}) := d_{\cos}(h_{\theta}^{\mathcal{X}}(\mathbf{x}), h_{\theta}^{\mathcal{Y}}(\mathbf{y}))$  is the multimodal distance, which has been shown empirically to provide a meaningful signal in prior label error detection work (Liang et al., 2023; Kang et al., 2023). We thus use this distance as the basis, and augment it with two additional terms based on nearest neighbors:

$$s_n(\mathbf{x}, \mathbf{y}, \mathcal{D}) = \frac{1}{k} \sum_{j=1}^k d_{\mathcal{Y}}(\mathbf{y}, \mathbf{y}_{n_j}) e^{-\tau_{1,n} d_{\mathcal{X}}(\mathbf{x}, \mathbf{x}_{n_j})} e^{-\tau_{2,n} d_{mm}(\mathbf{x}_{n_j}, \mathbf{y}_{n_j})}, \quad (2)$$

where  $(\mathbf{x}_{n_j}, \mathbf{y}_{n_j}) \in \mathcal{D}$ , and  $\tau_{1,n}, \tau_{2,n} \geq 0$  are hyperparameters. This corresponds to finding the nearest neighbors of  $\mathbf{x}$  in  $\mathcal{X}$  space, then averaging the distance between their *corresponding* modality in  $\mathcal{Y}$  and  $\mathbf{y}$ . We weight this average with two additional terms. The  $\tau_{1,n}$  term corresponds to downweighting neighbors which are far from  $\mathbf{x}$ . Intuitively, this is useful when  $k$  is too large for  $\mathbf{x}$  and not all neighbors are relevant, and can be thought of as an adaptive  $k$ . The  $\tau_{2,n}$  term corresponds to downweighting neighbors which are themselves likely to be mislabeled. If  $(\mathbf{x}_{n_j}, \mathbf{y}_{n_j})$  is itself mislabeled, then  $d_{\mathcal{Y}}(\mathbf{y}, \mathbf{y}_{n_j})$  would contribute an erroneous signal to whether  $(\mathbf{x}, \mathbf{y})$  is mislabeled, and we thus want to

<sup>2</sup>One could take, for any  $i$ ,  $(\mathbf{x}, \mathbf{y}) := (\mathbf{x}, \mathbf{y})_i, \mathcal{D}' := \mathcal{D} \setminus \{(\mathbf{x}, \mathbf{y})_i\}$

<sup>3</sup>We will use a subscript  $n_j$  to index nearest neighbors in  $\mathcal{X}$ , and subscript  $m_j$  for neighbors in  $\mathcal{Y}$ .

downweight those instances.

The third term is analogous to  $s_n$ , but uses neighbors of  $\mathbf{y}$ :

$$s_m(\mathbf{x}, \mathbf{y}, \mathcal{D}) = \frac{1}{k} \sum_{j=1}^k d_{\mathcal{X}}(\mathbf{x}, \mathbf{x}_{m_j}) e^{-\tau_{1,m} d_{\mathcal{Y}}(\mathbf{y}, \mathbf{y}_{m_j})} e^{-\tau_{2,m} d_{mm}(\mathbf{x}_{m_j}, \mathbf{y}_{m_j})}, \quad (3)$$

where  $(\mathbf{x}_{m_j}, \mathbf{y}_{m_j}) \in \mathcal{D}$ , and  $\tau_{1,m}, \tau_{2,m} \geq 0$  are hyperparameters. Crucially, note that notationally,  $\mathbf{x}_{n_j} \neq \mathbf{x}_{m_j}$ , and  $\mathbf{y}_{n_j} \neq \mathbf{y}_{m_j}$ . Specifically,  $\mathbf{y}_{n_j}$  corresponds to the  $\mathcal{Y}$  modality of nearest neighbors taken in  $\mathcal{X}$  space, and  $\mathbf{y}_{m_j}$  corresponds to the nearest neighbors of  $\mathbf{y}$  taken in  $\mathcal{Y}$  space.

We note that our method is a generalization of several prior methods. When  $\beta = \gamma = 0$ , the method is equivalent to CLIP similarity (Liang et al., 2023). When  $\beta$  is large,  $\tau_{1,n} = \tau_{2,n} = \gamma = 0$ , and  $d(\mathbf{y}, \mathbf{y}_{n_j}) = \mathbf{1}_{\mathbf{y} \neq \mathbf{y}_{n_j}}$ , the method is equivalent to Deep kNN (Bahri et al., 2020). An algorithm outline and high-level description of the method can be found in Appendix C.

Our method contains several hyperparameters:  $k, \beta, \gamma, \tau_{1,n}, \tau_{2,n}, \tau_{1,m}$ , and  $\tau_{2,m}$ . When there is a validation set with known mislabel flags, we perform a grid search over  $k$ , and use numerical optimization methods to search for an optimal value of the remaining hyperparameters which maximize label error detection performance on this set, which we describe further in Section 5.2. We refer to our method in this setting as LEMON<sub>OPT</sub>. We will empirically show that only a few hundred labeled validation samples may be sufficient to achieve optimal performance in this setting.

When there is no labeled validation set available, we will show that our method is fairly robust to these hyperparameter choices, and that choosing a set of reasonable fixed values for these hyperparameters yields nearly comparable results. We refer to our method in this setting as LEMON<sub>FIX</sub>.

## 4. Theoretical Analysis

We show that our multimodal kNN scores (Equations (2) and (3)) provide a signal for label error. Suppose there exists a ‘‘paraphrase function’’  $\mathcal{H} : \mathcal{Y} \rightarrow \mathcal{P}(\mathcal{Y})$ , where  $\mathcal{P}$  denotes the powerset, such that for a particular sample  $(x, y)$  with  $\mathcal{H}(y) = (\bar{y}_1, \bar{y}_2, \dots)$ ,  $(x, \bar{y}_i)$  is considered correctly labeled for all  $\bar{y}_i \in \mathcal{H}(y)$ . Informally,  $\mathcal{H}$  outputs the set of all possible captions which correctly describe  $x$ . Similarly define  $\mathcal{J}(x)$ , which outputs the set of images with identical semantics as  $x$ .

**Assumption 1** (Structure of  $\mathcal{H}, \mathcal{J}$ ):

- Let  $(x', y')$  be an arbitrary sample. If  $y' \notin \mathcal{H}(y)$ , then  $x' \notin \mathcal{J}(x)$ .
- Let  $(x', y')$  be an arbitrary mislabeled sample. Then,  $\forall y'' \in \mathcal{H}(y'), x'' \notin \mathcal{J}(x')$ .

**Assumption 2** (Distribution of Distances): *Let  $(X, Y)$  be a randomly drawn sample.*

- $\forall X' \notin \mathcal{J}(X) : d_{\mathcal{X}}(X, X') \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_1, \sigma_1^2)$ .
- $\forall \bar{X} \in \mathcal{J}(X) : d_{\mathcal{X}}(X, \bar{X}) \stackrel{\text{iid}}{\sim} \mathcal{N}(\mu_2, \sigma_2^2)$ .

We empirically validate this assumption in Appendix A.2.

Let  $N_k(Y) = \{Y_{m_1}, \dots, Y_{m_k}\}$  denote the nearest neighbors of  $Y$  in the text space. Let  $\frac{1}{k}|\mathcal{H}(Y) \cap N_k(Y)| = \zeta_Y$ , a random variable. Suppose that  $\frac{1}{k}|\{i : (X_{m_i}, Y_{m_i}) \text{ is mislabeled}\}| = p$  is constant for all samples in the support of  $(X, Y)$ .

Let  $S_m(X, Y) = \frac{1}{k} \sum_{Y_{m_i} \in N_k(Y)} d_{\mathcal{X}}(X, X_{m_i})$ , which is identical to the proposed Equation (3) with  $\tau_1 = \tau_2 = 0$ .

**Theorem 4.1** (AUROC of kNN Score). *Let  $(X, Y)$  be a randomly selected correctly labeled sample, and  $(X', Y')$  a randomly selected incorrectly labeled sample. Under Assumptions 1 and 2:*

$$\mathbb{P}(S_m(X', Y') > S_m(X, Y)) = 1 - \Phi\left(\frac{-\mu}{\sigma}\right)$$

where  $\mu = \mathbb{E}[\zeta_Y](1-p)(\mu_1 - \mu_2)$ ,  $\sigma = \left(\frac{\mathbb{E}[\zeta_Y](1-p)\sigma_2^2 + (2 - \mathbb{E}[\zeta_Y](1-p))\sigma_1^2}{k} + \text{Var}(\zeta_Y)(1-p)^2(\mu_2 - \mu_1)^2\right)^{1/2}$ , and  $\Phi$  is the Gaussian CDF.

This provides an expression for the detection AUROC of the score  $S_m$ . The same expression can be derived for  $S_n$  by symmetry.

**Lemma 4.2** (Non-random Signal of kNN Score). *If the following three conditions hold: (1)  $p < 1$ , (2)  $\mathbb{E}[\zeta_Y] > 0$ , (3)  $\mu_1 > \mu_2$ . Then,  $\mathbb{P}(S_m(X', Y') > S_m(X, Y)) > 0.5$ .*

Under these mild conditions,  $S_m$ , our proposed multimodal neighborhood score, provides a better than random signal at detecting mislabeled samples. The proof can be found in Appendix A.1. We additionally provide a theorem showing that embedding models trained via the contrastive multimodal objective are natural noisy label detectors in Appendix A.3.

## 5. Experiments

### 5.1. Datasets

We evaluate our method using eight datasets, as shown in Table 1. Four datasets (cifar10, cifar100, stanfordCars, miniImageNet) are label error detection datasets from the classification setting. The four remaining datasets are image captioning datasets. For msccoco and flickr30k, we use the Karpathy split (Karpathy & Fei-Fei, 2015). The remaining datasets were randomly split into: training or reference set for the label detection method (80%), validation set for hyperparameter selection (10%), and test set for performance evaluation (10%).

#### 5.1.1. NOISE TYPES

In cifar10 and cifar100, we utilize a dataset collected in prior work (Wei et al., 2021) with human mislabels (*human*). We also follow prior work (Zhu et al., 2022) in experimenting with class symmetric (*sym.*) and class asymmetric (*asym.*) synthetic noise. For stanfordCars and miniimagenet, we use datasets from Jiang et al. (2020), which contain noise from real-world (*real*) web annotators.

For the four captioning datasets, we devise several ways to inject synthetic noise of prevalence  $p$ . The simplest way is to randomly select  $p$  fraction (*random*) of the samples and assign their text modality to be that of another random caption. In datasets where additional metadata is available (msccoco: object category, mmimdb: genre of movie, mimiccxr: disease label), we can randomly swap the caption with that of another sample from the same category (*cat*). Finally, in all captioning datasets except mimiccxr, we tag each token of each caption with its part-of-speech using SpaCy (Honnibal & Montani, 2017), and then randomly assign a selected sample’s text modality to be from another sample with at least one noun in common (*noun*). Dataset processing details are in Appendix D.

These noise types are intended to simulate an array of realistic label corruptions in the real world. As such, the resulting synthetic dataset may not have an exact noise level  $p$ , as e.g. a randomly selected caption may actually be correct for the image, as well as due to noise in the base datasets, which we explore in Section 6.5. Unless otherwise stated, results shown in the main paper are for the bolded noise type in Table 1, with 40% synthetic noise. Additional results for other noise types can be found in Appendix I.

### 5.2. Model Selection and Evaluation

We run LEMON on each dataset, using the training split of each dataset to compute nearest neighbors. In classification datasets, we use the discrete metric  $d_Y(\mathbf{y}, \mathbf{y}') = \mathbf{1}_{\mathbf{y} \neq \mathbf{y}'}$ . In all other cases and for  $d_{\mathcal{X}}$ , we utilize cosine or euclidean distance computed in the embedding space of a pretrained CLIP model, selecting the best distance metric on the validation set for LEMON<sub>OPT</sub>, and keeping the distance as the cosine distance for LEMON<sub>FIX</sub>. In mimiccxr, we use BiomedCLIP (ViT-B/16) (Zhang et al., 2023b), and we use OpenAI CLIP ViT-B/32 (Radford et al., 2021) for all other datasets. A full list of hyperparameters for our method and the baselines are in Appendix G.

For LEMON<sub>OPT</sub>, we select the hyperparameter combination that maximizes F1 on a labeled validation set. We report the AUROC, macro-averaged AUPRC, and F1 for this model. For LEMON<sub>FIX</sub>, we fix the hyperparameters at the following reasonable values:  $k = 30$ ,  $\beta = \gamma = 5$ ,  $\tau_{1,n} = \tau_{1,m} = 0.1$ , and  $\tau_{2,n} = \tau_{2,m} = 5$ . We report AU-

Table 1: Classification and captioning datasets.  $n$  is the number of samples. In the main paper, results shown are for the bolded noise type with 40% noise level for synthetic noise. Results for remaining settings can be found in the appendices.

Dataset	$n$			Domain		Noise Types
	Train	Validation	Test	Image	Text	
cifar10	40,000	5,000	5,000	Natural images	Object labels	{ <b>human</b> (Wei et al., 2021), <i>sym.</i> , <i>asym.</i> }
cifar100	40,000	5,000	5,000	Natural images	Object labels	{ <b>human</b> (Wei et al., 2021), <i>sym.</i> , <i>asym.</i> }
miniImageNet (Jiang et al., 2020)	49,419	24,710	24,710	Natural images	Object labels	{ <b>real</b> }
stanfordCars (Jiang et al., 2020)	13,501	6,751	6,752	Car images	Car year and model	{ <b>real</b> }
mscoco (Lin et al., 2014)	82,783	5,000	5,000	Natural images	Captions	{ <b>cat.</b> , <i>noun</i> , <i>random</i> }
flickr30k (Young et al., 2014)	29,000	1,014	1,000	Natural images	Captions	{ <b>noun</b> , <i>random</i> }
mmimdb (Arevalo et al., 2017)	15,552	2,608	7,799	Movie Posters	Plot summaries	{ <b>cat.</b> , <i>noun</i> , <i>random</i> }
mimiccxr (Johnson et al., 2019)	368,909	2,991	5,159	Chest X-rays	Radiology reports	{ <b>cat.</b> , <i>random</i> }

ROC and AUPRC, as the F1 requires additional information to compute a threshold for the score. We recognize that access to such a validation set as in LEMON<sub>OPT</sub> may be unrealistic, but we will empirically show that (1) our method is fairly robust to selection of these hyperparameters, (2) only a few hundred labeled samples may be sufficient to select these hyperparameters, (3) using LEMON<sub>FIX</sub> with the fixed hyperparameter setting described above achieves nearly comparable results, and (4) hyperparameters optimized on a dataset with synthetic noise may transfer well to real datasets.

We repeat each experiment three times, using a different random seed for the noise sampling (for *human* and *real* noise, we use a different random data split). Performance metrics shown are test-set results averaged over these three runs, with error bounds corresponding to one standard deviation.

**Baselines** We compare our method versus previous state-of-the-art in both the classification and captioning settings. We additionally adapt several baselines from the classification setting to the captioning setting. We briefly list the baselines here, and a detailed description is in Appendix E.

In the classification setting, we experiment with the following baselines which require training a classifier on the particular dataset: **AUM** (Pleiss et al., 2020), **Datamap** (Swayamdipta et al., 2020), and **Confident Learning** (Northcutt et al., 2021a), and the following baselines which do not require classifier training: **Deep k-NN** (Bahri et al., 2020), **SimiFeat** (Zhu et al., 2022)-Voting and Ranking, discrepancy in the image space (**Discrepancy**) ( $\Upsilon_X^{DIS}$  from Thomas & Kovashka (2022)), **CLIP Similarity** (Kang et al., 2023), and **CLIP Logits** (Liang et al., 2023; Feng et al.).

In the captioning setting, we compare our method with **LLaVA** (Liu et al., 2024) prompting and **CapFilt** (Li et al., 2022). We note that the latter can be viewed as an oracle for natural image captioning, as it has been trained in a supervised manner on *clean* mscoco data. **CLIP Similarity** (Kang et al., 2023), **Discrepancy** (Thomas & Kovashka, 2022), and **Datamap** (Swayamdipta et al., 2020) can also

be used directly in this setting. Next, to adapt classification baselines to captioning, we embed the captions using the corresponding CLIP text encoder, and then use K-means clustering to assign the text caption into one of 100 clusters. We then apply **Deep k-NN** (Bahri et al., 2020) and **Confident Learning** (Northcutt et al., 2021a), using the cluster ID as the discretized class. Finally, we adapt **VDC** (Zhu et al., 2024) to the captioning setting using open source models. Further details can be found in Appendix E.

## 6. Results

### 6.1. LEMON Outperforms Baselines on Label Error Detection

**Classification** In Table 2, we show the performance of LEMON against the baselines for label error detection on four classification datasets. We find that our method performs on-par with, or outperforms, existing baselines which do not require classifier training on all classification datasets. Two downstream-task specific approaches (AUM and Datamap) outperform most training-free models (particularly on *cifar100*), but LEMON performs comparably and even outperforms them in two datasets. Similar results are also observed on the two synthetic error types (see Appendix Table I.2). We find that LEMON<sub>FIX</sub> performs almost comparably with LEMON<sub>OPT</sub>, and still beats almost all baselines.

**Captioning** In Table 3, we find that our method outperforms existing neighborhood and similarity-based baselines on three datasets. In two datasets, our model underperforms the oracle (CapFilt) which has been pre-trained on *clean* captions from the mscoco dataset. Results for synthetic error types show similar trends (see Appendix I.2).

### Label Error Detection Performance Across Noise Levels

In Figure I.1, we show the performance of LEMON versus the CLIP similarity baseline on mscoco and mmimdb, varying the level of the synthetic noise. We find that LEMON performs better uniformly across noise levels.

Table 2: Label error detection performance across classification datasets. We separate AUM, Datamap, and Confident learning, as they require training a classifier from scratch. Bold denotes best score within each training approach. A full version of this table with AUPRC can be found in Appendix I.1.

Method	Training-Free	cifar10		cifar100		miniImageNet		stanfordCars	
		AUROC	F1	AUROC	F1	AUROC	F1	AUROC	F1
AUM		<b>98.3</b> (0.1)	<b>92.9</b> (0.1)	<b>92.3</b> (0.3)	<b>81.1</b> (0.3)	83.1 (0.2)	68.3 (0.4)	70.4 (2.3)	47.2 (3.1)
Datamap	✗	98.2 (0.1)	92.2 (0.5)	91.8 (0.3)	80.8 (0.5)	<b>85.1</b> (0.3)	<b>70.6</b> (0.2)	<b>72.2</b> (1.7)	<b>50.4</b> (2.1)
Confident		89.6 (1.4)	88.2 (1.7)	78.6 (0.4)	73.7 (0.5)	59.5 (0.7)	37.7 (1.5)	60.7 (0.3)	39.9 (0.6)
CLIP Logits		95.5 (0.2)	86.2 (0.6)	84.9 (0.7)	72.0 (0.9)	<b>90.0</b> (0.2)	<b>77.1</b> (0.2)	68.8 (0.1)	47.3 (0.5)
CLIP Sim.		92.2 (0.2)	82.3 (0.3)	80.8 (0.9)	68.7 (1.1)	89.3 (0.2)	76.1 (0.4)	69.8 (0.4)	46.6 (0.5)
Simifeat-V		90.8 (0.3)	88.2 (0.5)	79.5 (0.1)	73.3 (0.2)	68.1 (0.2)	54.8 (0.3)	63.6 (1.3)	43.5 (1.6)
Simifeat-R		88.8 (0.1)	85.9 (0.1)	81.4 (0.1)	76.1 (0.1)	65.1 (0.3)	48.6 (0.5)	66.8 (0.3)	47.9 (0.5)
Discrepancy	✓	77.1 (1.9)	66.4 (2.2)	66.0 (1.5)	58.9 (0.8)	79.5 (0.2)	64.0 (0.1)	65.7 (0.3)	44.3 (0.7)
Deep k-NN		97.8 (0.1)	91.4 (0.6)	87.4 (0.3)	75.7 (0.3)	83.2 (0.2)	68.5 (0.2)	71.5 (0.6)	49.1 (0.6)
LEMON <sub>FIX</sub> (Ours)		97.7 (0.2)	90.9 (0.1)	88.9 (0.7)	75.4 (0.6)	89.5 (0.2)	74.7 (0.2)	72.6 (0.7)	47.7 (2.0)
LEMON <sub>OPT</sub> (Ours)		<b>98.1</b> (0.0)	<b>92.0</b> (0.2)	<b>90.8</b> (0.0)	<b>78.4</b> (0.0)	<b>90.0</b> (0.4)	76.9 (0.2)	<b>73.1</b> (0.5)	<b>51.3</b> (0.5)

Table 3: Label error detection performance on captioning datasets. We separate CapFilt as it has been trained on *clean* `mscoco` data. Bold denotes best (highest) score. A full version of this table with AUPRC can be found in Appendix I.2.

Method	flickr30k		mscoco		mmimdb		mimiccxr	
	AUROC	F1	AUROC	F1	AUROC	F1	AUROC	F1
LLaVA	79.3 (0.8)	65.0 (1.1)	80.3 (0.1)	74.9 (0.3)	58.4 (0.2)	58.5 (0.1)	53.9 (0.5)	57.0 (0.1)
Datamap	80.6 (1.3)	66.6 (0.6)	85.8 (0.4)	74.3 (0.3)	55.4 (0.5)	57.2 (0.0)	51.0 (0.5)	57.0 (0.1)
Discrepancy	73.0 (0.6)	59.0 (0.3)	72.7 (0.3)	62.5 (0.3)	57.8 (0.4)	57.4 (0.2)	60.0 (0.7)	57.2 (0.1)
VDC	92.9 (1.1)	81.1 (1.6)	94.1 (0.2)	86.3 (0.4)	80.5 (0.3)	69.3 (0.6)	50.8 (0.4)	57.0 (0.1)
Deep k-NN	71.1 (0.4)	59.2 (0.8)	76.6 (0.4)	67.5 (0.8)	61.2 (0.4)	58.3 (0.4)	62.9 (0.4)	59.2 (0.1)
Confident	63.1 (0.9)	54.0 (0.9)	71.5 (0.5)	66.5 (0.5)	52.8 (1.1)	51.8 (1.8)	61.8 (0.3)	58.1 (0.6)
CLIP Sim.	<b>94.8</b> (0.5)	<b>84.2</b> (0.9)	93.8 (0.2)	84.5 (0.4)	85.1 (0.3)	72.7 (0.6)	64.1 (0.4)	59.2 (0.0)
LEMON <sub>FIX</sub> (Ours)	93.6 (0.2)	-	92.0 (0.1)	-	84.3 (0.3)	-	66.3 (0.4)	-
LEMON <sub>OPT</sub> (Ours)	94.5 (0.2)	83.6 (1.4)	<b>95.6</b> (0.2)	<b>87.0</b> (0.2)	<b>86.0</b> (0.1)	<b>73.5</b> (0.3)	<b>70.4</b> (1.6)	<b>61.1</b> (0.8)
CapFilt (Supervised Training)	98.6 (0.1)	93.1 (0.7)	99.3 (0.0)	95.4 (0.4)	82.7 (0.7)	71.3 (0.3)	49.7 (0.3)	57.0 (0.0)

**Size of Labeled Validation Set** In Appendix Figure I.3, we examine how varying the size of the labeled validation set impacts the performance of LEMON<sub>OPT</sub>. We find that in all four captioning datasets, having about 100-500 labeled examples is sufficient to tune hyperparameters in LEMON<sub>OPT</sub> to outperform LEMON<sub>FIX</sub>. In the three datasets where LEMON<sub>FIX</sub> underperforms the CLIP similarity baseline, we find again that having 100-500 labeled validation samples is sufficient for tuning LEMON<sub>OPT</sub> to perform on par with this baseline.

**Robustness to Hyperparameters** Here, we test the robustness of our method when there is no labeled validation set available. First, in Appendix I.4, we visualize the F1 of the selected score when varying  $\beta$  and  $\gamma$ , keeping all other hyperparameters at their selected optimal values. We find that for most datasets and noise types, there is a reasonably large space of such hyperparameters, bounded away from the origin, which achieves close to optimal performance.

Next, we compare the performance of LEMON<sub>OPT</sub> and LEMON<sub>FIX</sub> with hyperparameters described in Section 5.2 across all datasets in Table I.8. We find that when there is no labeled validation set available, using these hyperparameters

results in an AUROC drop of only 1.6% on average (std = 1.3%), with a worst-case AUROC drop of 4.1% across all 18 dataset and noise type combinations. Thus, even when a labeled validation set is not available, LEMON<sub>FIX</sub> with reasonable hyperparameter settings is able to outperform most baselines which do use such information.

## 6.2. Filtering Misabeled Data Improves Downstream Performance

**Classification** To assess the impact of label error detection on the performance of the downstream classification tasks, we filter out samples from the training set with mislabel scores in the top  $q$  percentile. We vary  $q$ , train ViT (Dosovitskiy et al., 2020) models on the filtered dataset, and evaluate the downstream test accuracy using clean data. We compare the performance of LEMON<sub>OPT</sub> with all training-free baselines that produce a continuous score (i.e. all except Simifeat and Confident). In Figure 3, we find that training with LEMON<sub>OPT</sub> filtered samples leads to the highest accuracy on `cifar10` and `cifar100` after removing more than 30% of the data. Training with LEMON<sub>OPT</sub> filtered samples is also on par with baselines on the other datasets (either the best or within 1% of the best baseline) as shown in Ap-

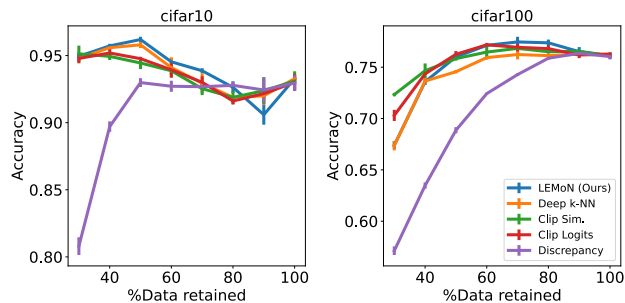


Figure 3: Downstream classification accuracy on `cifar10` (left) and `cifar100` (right) with `LEMONOPT` with *human* noise versus the baselines. Note that the noise prevalence is 40% in both datasets.

Table 4: Downstream captioning performance when removing 40% samples with highest mislabel scores. We find that filtering noisy data with `LEMONOPT` improves captioning.

Dataset	Method	B@4	CIDER	ROUGE
flickr30k	No Filtering	31.6 (0.6)	75.0 (1.6)	51.4 (0.4)
	CLIP Sim.	<b>32.7 (0.3)</b>	<b>79.6 (0.4)</b>	<b>52.8 (0.1)</b>
	<code>LEMON<sub>OPT</sub></code>	32.1 (0.6)	78.6 (1.3)	52.5 (0.2)
	Clean	33.3 (0.1)	79.9 (1.1)	53.2 (0.3)
mscoco	No Filtering	36.0 (0.1)	118.9 (0.3)	56.8 (0.1)
	CLIP Sim.	39.5 (0.1)	130.4 (0.3)	59.3 (0.0)
	<code>LEMON<sub>OPT</sub></code>	<b>39.6 (0.1)</b>	<b>130.8 (0.1)</b>	<b>59.4 (0.0)</b>
	Clean	39.6 (0.1)	131.2 (0.2)	59.5 (0.1)

pendix I.14. We also show that filtering data in this manner does not reduce classifier robustness (Appendix I.16).

**Captioning** We finetune a pre-trained GenerativeImage2Text (GIT) (Wang et al., 2022a) model to generate captions. Given the large size of the model, we use the parameter-efficient Low-Rank Adaptation (LoRA) (Hu et al., 2021) for all captioning models. We train models with clean data, noisy captions (*No Filtering*), and by filtering data detected as being mislabeled by a label detection method. In Table 4, we compare results of using either our model or a strong baseline (CLIP Sim.) for filtering data, as measured by the BLEU-4 (Papineni et al., 2002), CIDER (Vedantam et al., 2015), and ROUGE (Lin, 2004) scores. In all cases, we filtered out the top 40% percentile of data predicted to be mislabeled (i.e., equal to the expected prevalence of noisy data). We find that (1) filtering out data predicted to be mislabeled helps recover performance as compared to training on fully clean data along multiple metrics, and (2) our method performs comparably to the baseline in improving downstream results, with both filtering techniques perform within 1 score point of each other. However, as training with a fully clean dataset only outperforms training with a fully noisy (i.e. no filtering) dataset by 2-3 BLEU-4 points, the range of potential improvement for any filtering method is limited.

### 6.3. Ablations

In Table I.10, we show the performance of our method after ablating each component. We find that mislabel detection performance almost decreases monotonically as we remove additional components until we reach the CLIP Similarity baseline. We find that ablating the  $\tau_1$  and  $\tau_2$  terms results in a performance loss of about 1%. In Table I.11, we examine the performance of each of the three components of our score and their combinations. We find that  $d_{mm}$  is the most critical term. Of the two nearest neighbors terms, we find that  $s_n$  (nearest image neighbors) is more important in general, though this is highly dataset dependent, e.g. error detection in `mimicdb` relies much more on neighbors in the text space than the image space, while the opposite is true for `mscoco`.

### 6.4. External Pretraining May Not Be Required

**Medical Images** Thus far, all of the results for LEMON (and CLIP Similarity) have utilized CLIP models which have been pretrained on external datasets (e.g. PMC-15M in the case of BiomedCLIP). Here, we examine whether we can achieve comparable performance by pretraining CLIP from scratch *only on the noisy data*. We select `mimiccxr` as it has the most samples out of all captioning datasets. Similar to CheXzero (Tiu et al., 2022), we pretrain a CLIP ViT B/16 from scratch on the `mimiccxr` training set with 40% noise. We train this model for 10 epochs with a batch size of 64, and do not do any model selection or early stopping. We then apply LEMON and the CLIP similarity baseline using this model, for the same noise level and noise type. We present our results in Table 5. Surprisingly, we find that pretraining CLIP only on noisy data from MIMIC-CXR actually outperforms BiomedCLIP. This could be attributed to the pretraining domain (chest X-rays and radiology notes) matching the inference domain exactly (Nguyen et al., 2022). As an upper bound, we evaluate the same methods using CheXzero (Tiu et al., 2022), which has been pretrained on *clean* MIMIC-CXR data. We find that, as expected, it far outperforms this baseline. We conclude that, for large noisy datasets, pretraining a CLIP model from scratch could be a viable solution, though pretraining on clean data from the same domain is certainly superior.

**Web-Scale Corpus** Motivated by this result, we conduct a large scale experiment on the CC3M dataset (Changpinyo et al., 2021), which contains 2.9 million valid URLs to image-caption pairs. We pretrain CLIP from scratch on this dataset, then use this CLIP model to filter samples in the original dataset using `LEMONFIX` and the CLIP similarity baseline. We select the 1 million samples with the lowest mislabel scores from each method, and pretrain another CLIP from scratch on this clean subset. We evaluate the resulting model on zero-shot classification using the

Table 5: Performance of LEMON for label error detection versus the CLIP similarity baseline on `mimiccxr`, when external pretrained models may not be available. Biomed-CLIP (Zhang et al., 2023a) is trained on a large corpus of biomedical image-text pairs. We find that pretraining only on noisy data from MIMIC-CXR outperforms Biomed-CLIP, though pretraining on clean `mimiccxr` data (as in CheXzero (Tiu et al., 2022)) does perform better.

		Random Noise		Cat. Noise	
		AUROC	F1	AUROC	F1
BiomedCLIP	Clip Sim.	66.8 (0.8)	60.1 (0.4)	64.1 (0.4)	59.2 (0.0)
	LEMON <sub>FIX</sub> (Ours)	69.5 (0.7)	-	66.3 (0.4)	-
	LEMON <sub>OPT</sub> (Ours)	<b>73.7</b> (1.7)	<b>63.5</b> (0.8)	<b>70.4</b> (1.6)	<b>61.1</b> (0.8)
CLIP Pretrain On Noisy Data	Clip Sim.	78.8 (0.1)	66.8 (0.3)	76.5 (0.5)	64.4 (0.5)
	LEMON <sub>FIX</sub> (Ours)	<b>80.5</b> (0.1)	-	77.0 (0.5)	-
	LEMON <sub>OPT</sub> (Ours)	80.0 (0.9)	<b>67.7</b> (0.5)	<b>77.2</b> (0.8)	<b>64.6</b> (0.3)
CheXzero	Clip Sim.	90.8 (0.2)	79.5 (0.2)	88.4 (0.6)	76.1 (0.7)
	LEMON <sub>FIX</sub> (Ours)	91.4 (0.1)	-	88.4 (0.7)	-
	LEMON <sub>OPT</sub> (Ours)	<b>91.8</b> (0.2)	<b>80.9</b> (0.6)	<b>88.5</b> (0.4)	<b>76.2</b> (1.0)

VTAB benchmark (Zhai et al., 2019). We find filtering with LEMON marginally outperforms the baseline on average zero-shot accuracy, though both underperform pretraining on the full corpus. Full details are in Appendix I.9. We additionally conduct an experiment on Datacomp (Gadre et al., 2024) in Appendix I.10.

### 6.5. Real-World Analysis

We conduct a preliminary study of LEMON on real datasets without known label errors. We run LEMON<sub>FIX</sub> and the CLIP similarity baseline on `cifar10`, `cifar100`, `flickr30k`, and `mscoco`. As no labeled validation set is available, we use optimal hyperparameters from models previously run on each dataset with synthetic noise from Section 6.1 (Appendix I.11). For each dataset, we select the top 200 images from the validation and test splits with the highest mislabel scores. We then manually annotated each sample to determine whether it was mislabeled. Crucially, during labeling, images were randomly selected, so the labeler is unaware of whether the candidate image originated from the baseline or our method. We find that our method outperforms the baseline for every dataset (Table 6). Examples of real-world mislabels are also in Figures 1 and I.5. We present a further comparison of our identified error sets in `cifar10` and `cifar100` with crowd-sourced labels (Northcutt et al., 2021b) in Appendix I.13.

## 7. Conclusion

In this work, we proposed LEMON, a method that leverages the neighborhood structure of contrastively pretrained multimodal embeddings to automatically identify label errors in image datasets with natural language text labels.

Our approach is a promising step to automatically detecting

Table 6: We manually label 200 images from real-world datasets that each method identifies as the most likely to be mislabeled and show the percentage (%) of times where it is actually mislabeled. Numbers in parentheses are 95% confidence intervals from a binomial proportion.

	CLIP Sim.	Ours
<code>cifar10</code>	5.5 (3.2)	<b>10.0</b> (4.2)
<code>cifar100</code>	11.0 (4.3)	<b>20.5</b> (5.6)
<code>flickr30k</code>	32.5 (6.5)	<b>41.0</b> (6.8)
<code>mscoco</code>	19.5 (5.5)	<b>25.5</b> (6.0)

and filtering data mislabels at scale. Through experiments on multiple datasets with synthetic and real-world noise, we demonstrated LEMON’s effectiveness in detecting label errors and improving downstream model performance.

**Limitations** Our work has several limitations. As we primarily rely on existing open-sourced datasets, some parts of these datasets may have been used as training data in pretrained models. We specifically chose pretrained models that take care not to include the test sets of such datasets. Further, we run experiments on a real-world healthcare dataset with access controls (`mimiccxr`) to verify our results. Next, we assume that there exists an oracle binary indicator for whether a sample is mislabeled. As we saw in practice, real-world mislabels contain much more uncertainty and ambiguity, e.g. due to blurry images (Gao et al., 2017; Basile et al., 2021; Gordon et al., 2021; 2022). Evaluating the effectiveness of our score as a measure of this uncertainty is an area of future work.

## Acknowledgments

This work was supported in part by a National Science Foundation (NSF) 22-586 Faculty Early Career Development Award (#2339381), a Gordon & Betty Moore Foundation award, and a Google Research Scholar award. We would like to thank Walter Gerych and Olawale Salaudeen for their valuable feedback.

## Impact Statement

In this work, we have studied the identification and removal of samples deemed to be mislabels. We recognize that not all such samples are truly erroneous, and that these false positives may not be uniformly distributed across the space of images and text. For example, automated detectors such as LEMON may flag legitimate but atypical examples, such as rare concepts, under-represented languages or dialectic phrases, and images from minority classes or groups. Removal of these samples may then lead to downstream models with certain biases. To alleviate these issues, we encourage practitioners to use LEMON and similar tools

to surface candidates for expert review, rather than as an unsupervised pruning tool. We view the ability to manually examine such samples as a strength of the filtering approach, over loss-based approaches, for learning under label noise.

## References

- Arevalo, J., Solorio, T., Montes-y Gómez, M., and González, F. A. Gated multimodal units for information fusion. *arXiv preprint arXiv:1702.01992*, 2017.
- Bahri, D., Jiang, H., and Gupta, M. Deep k-nn for noisy labels. In *International Conference on Machine Learning*, pp. 540–550. PMLR, 2020.
- Balestriero, R., Ibrahim, M., Sobal, V., Morcos, A., Shekhar, S., Goldstein, T., Bordes, F., Bardes, A., Mialon, G., Tian, Y., et al. A cookbook of self-supervised learning. *arXiv preprint arXiv:2304.12210*, 2023.
- Basile, V., Cabitza, F., Campagner, A., and Fell, M. Toward a perspectivist turn in ground truthing for predictive computing. *arXiv preprint arXiv:2109.04270*, 2021.
- Bernhardt, M., Castro, D. C., Tanno, R., Schwaighofer, A., Tezcan, K. C., Monteiro, M., Bannur, S., Lungren, M. P., Nori, A., Glocker, B., et al. Active label cleaning for improved dataset quality under resource constraints. *Nature communications*, 13(1):1161, 2022.
- Beyer, L., Hénaff, O. J., Kolesnikov, A., Zhai, X., and Oord, A. v. d. Are we done with imagenet? *arXiv preprint arXiv:2006.07159*, 2020.
- Cai, S., Qiu, L., Chen, X., Zhang, Q., and Chen, L. Semantic-enhanced image clustering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 6869–6878, 2023.
- Changpinyo, S., Sharma, P., Ding, N., and Soricut, R. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3558–3568, 2021.
- Chen, H., Wang, J., Shah, A., Tao, R., Wei, H., Xie, X., Sugiyama, M., and Raj, B. Understanding and mitigating the label noise in pre-training on downstream tasks. *arXiv preprint arXiv:2309.17002*, 2023.
- Chen, S., Gallifant, J., Gao, M., Moreira, P., Munch, N., Muthukkumar, A., Rajan, A., Kolluri, J., Fiske, A., Hastings, J., et al. Cross-care: Assessing the healthcare implications of pre-training data on language model bias. *arXiv preprint arXiv:2405.05506*, 2024.
- Chen, T., Kornblith, S., Norouzi, M., and Hinton, G. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Cui, Z., Zhang, Y., and Ji, Q. Label error correction and generation through label relationships. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 3693–3700, 2020.
- Dai, W., Li, J., Li, D., Tiong, A., Zhao, J., Wang, W., Li, B., Fung, P., and Hoi, S. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv preprint arXiv:2305.06500*, 2, 2023.
- Diffenderfer, J., Bartoldson, B., Chaganti, S., Zhang, J., and Kailkhura, B. A winning hand: Compressing deep networks can improve out-of-distribution robustness. *Advances in neural information processing systems*, 34:664–676, 2021.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Feng, C., Tzimiropoulos, G., and Patras, I. Clipcleaner: Cleaning noisy labels with clip. In *ACM Multimedia 2024*.
- Fu, Z., Song, K., Zhou, L., and Yang, Y. Noise-aware image captioning with progressively exploring mismatched words. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 12091–12099, 2024.
- Gadre, S. Y., Ilharco, G., Fang, A., Hayase, J., Smyrnis, G., Nguyen, T., Marten, R., Wortsman, M., Ghosh, D., Zhang, J., et al. Datacomp: In search of the next generation of multimodal datasets. *Advances in Neural Information Processing Systems*, 36, 2024.
- Gao, B.-B., Xing, C., Xie, C.-W., Wu, J., and Geng, X. Deep label distribution learning with label ambiguity. *IEEE Transactions on Image Processing*, 26(6):2825–2838, 2017.
- Gertz, R. J., Dratsch, T., Bunck, A. C., Lennartz, S., Iuga, A.-I., Hellmich, M. G., Persigehl, T., Pennig, L., Gietzen, C. H., Fervers, P., et al. Potential of gpt-4 for detecting errors in radiology reports: Implications for reporting accuracy. *Radiology*, 311(1):e232714, 2024.

- Ghasemi, A. and Zahediasl, S. Normality tests for statistical analysis: a guide for non-statisticians. *International journal of endocrinology and metabolism*, 10(2):486, 2012.
- Gordon, M. L., Zhou, K., Patel, K., Hashimoto, T., and Bernstein, M. S. The disagreement deconvolution: Bringing machine learning performance metrics in line with reality. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pp. 1–14, 2021.
- Gordon, M. L., Lam, M. S., Park, J. S., Patel, K., Hancock, J., Hashimoto, T., and Bernstein, M. S. Jury learning: Integrating dissenting voices into machine learning models. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pp. 1–19, 2022.
- Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Grivas, A., Alex, B., Grover, C., Tobin, R., and Whiteley, W. Not a cute stroke: analysis of rule-and neural network-based information extraction systems for brain radiology reports. In *Proceedings of the 11th international workshop on health text mining and information analysis*, pp. 24–37, 2020.
- Hendrycks, D. and Dietterich, T. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations*, 2018.
- Honnibal, M. and Montani, I. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017.
- Hu, E. J., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W., et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.
- Hu, J. C., Cavicchioli, R., and Capotondi, A. Exploiting multiple sequence lengths in fast end to end training for image captioning. In *2023 IEEE International Conference on Big Data (BigData)*, pp. 2173–2182. IEEE, 2023.
- Huang, B., He, F., Wang, Q., Chen, H., Li, G., Feng, Z., Wang, X., and Zhu, W. Neighbor does matter: Global positive-negative sampling for vision-language pre-training. In *ACM Multimedia 2024*, 2024.
- Huang, R., Long, Y., Han, J., Xu, H., Liang, X., Xu, C., and Liang, X. Nlip: Noise-robust language-image pre-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 926–934, 2023.
- Jiang, L., Huang, D., Liu, M., and Yang, W. Beyond synthetic noise: Deep learning on controlled noisy labels. In *International conference on machine learning*, pp. 4804–4815. PMLR, 2020.
- Johnson, A. E., Pollard, T. J., Greenbaum, N. R., Lungren, M. P., Deng, C.-y., Peng, Y., Lu, Z., Mark, R. G., Berkowitz, S. J., and Horng, S. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs. *arXiv preprint arXiv:1901.07042*, 2019.
- Kang, W., Mun, J., Lee, S., and Roh, B. Noise-aware learning from web-crawled image-text data for image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2942–2952, 2023.
- Karpathy, A. and Fei-Fei, L. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3128–3137, 2015.
- Kim, T., Ko, J., Choi, J., Yun, S.-Y., et al. Fine samples for learning with noisy labels. *Advances in Neural Information Processing Systems*, 34:24137–24149, 2021.
- Lai, Z., Vesdapunt, N., Zhou, N., Wu, J., Huynh, C. P., Li, X., Fu, K. K., and Chuah, C.-N. Padclip: Pseudo-labeling with adaptive debiasing in clip for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 16155–16165, 2023.
- Li, J., Li, D., Xiong, C., and Hoi, S. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International conference on machine learning*, pp. 12888–12900. PMLR, 2022.
- Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pp. 19730–19742. PMLR, 2023.
- Li, Y., Liang, F., Zhao, L., Cui, Y., Ouyang, W., Shao, J., Yu, F., and Yan, J. Supervision exists everywhere: A data efficient contrastive language-image pre-training paradigm. In *International Conference on Learning Representations*, 2021.
- Liang, C., Zhu, L., Shi, H., and Yang, Y. Combating label noise with a general surrogate model for sample selection. *arXiv preprint arXiv:2310.10463*, 2023.
- Liang, V. W., Zhang, Y., Kwon, Y., Yeung, S., and Zou, J. Y. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. *Advances in Neural Information Processing Systems*, 35: 17612–17625, 2022.

- Liao, T., Taori, R., Raji, I. D., and Schmidt, L. Are we learning yet? a meta review of evaluation failures across machine learning. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.
- Lin, C.-Y. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81, 2004.
- Lin, K., Li, L., Lin, C.-C., Ahmed, F., Gan, Z., Liu, Z., Lu, Y., and Wang, L. Swinbert: End-to-end transformers with sparse attention for video captioning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 17949–17958, 2022.
- Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Liu, B., Bubeck, S., Eldan, R., Kulkarni, J., Li, Y., Nguyen, A., Ward, R., and Zhang, Y. Tinygsm: achieving 80% on gsm8k with small language models. *arXiv preprint arXiv:2312.09241*, 2023.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.
- Longpre, S., Yauney, G., Reif, E., Lee, K., Roberts, A., Zoph, B., Zhou, D., Wei, J., Robinson, K., Mimno, D., et al. A pretrainer’s guide to training data: Measuring the effects of data age, domain coverage, quality, & toxicity. *arXiv preprint arXiv:2305.13169*, 2023.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018.
- Luccioni, A. S. and Rolnick, D. Bugs in the data: How imagenet misrepresents biodiversity. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 14382–14390, 2023.
- Menghini, C., Delworth, A., and Bach, S. Enhancing clip with clip: Exploring pseudolabeling for limited-label prompt tuning. *Advances in Neural Information Processing Systems*, 36:60984–61007, 2023.
- Misra, I. and Maaten, L. v. d. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6707–6717, 2020.
- Natarajan, N., Dhillon, I. S., Ravikumar, P. K., and Tewari, A. Learning with noisy labels. *Advances in neural information processing systems*, 26, 2013.
- Nguyen, T., Ilharco, G., Wortsman, M., Oh, S., and Schmidt, L. Quality not quantity: On the interaction between dataset design and robustness of clip. *Advances in Neural Information Processing Systems*, 35:21455–21469, 2022.
- Northcutt, C., Jiang, L., and Chuang, I. Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021a.
- Northcutt, C. G., Athalye, A., and Mueller, J. Pervasive label errors in test sets destabilize machine learning benchmarks. *arXiv preprint arXiv:2103.14749*, 2021b.
- Oord, A. v. d., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- Pleiss, G., Zhang, T., Elenberg, E., and Weinberger, K. Q. Identifying mislabeled data using the area under the margin ranking. *Advances in Neural Information Processing Systems*, 33:17044–17056, 2020.
- Plummer, B. A., Wang, L., Cervantes, C. M., Caicedo, J. C., Hockenmaier, J., and Lazebnik, S. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pp. 2641–2649, 2015.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Rennie, S. J., Marcheret, E., Mroueh, Y., Ross, J., and Goel, V. Self-critical sequence training for image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7008–7024, 2017.
- Ridnik, T., Ben-Baruch, E., Noy, A., and Zelnik-Manor, L. Imagenet-21k pretraining for the masses. *Proceedings of 35th Conference on Neural Information Processing Systems, Track on Datasets and Benchmarks*, 2021.
- Rottmann, M. and Reese, M. Automated detection of label errors in semantic segmentation datasets via deep learning and uncertainty quantification. In *Proceedings*

- of the *IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3214–3223, 2023.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115: 211–252, 2015.
- Schrodi, S., Hoffmann, D. T., Argus, M., Fischer, V., and Brox, T. Two effects, one trigger: On the modality gap, object bias, and information imbalance in contrastive vision-language representation learning. *arXiv preprint arXiv:2404.07983*, 2024.
- Schroff, F., Kalenichenko, D., and Philbin, J. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 815–823, 2015.
- Schubert, M., Riedlinger, T., Kahl, K., Kröll, D., Schoenen, S., Šegvić, S., and Rottmann, M. Identifying label errors in object detection datasets by loss inspection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 4582–4591, 2024.
- Schuhmann, C., Vencu, R., Beaumont, R., Kaczmarczyk, R., Mullis, C., Katta, A., Coombes, T., Jitsev, J., and Komatsuzaki, A. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- Sohn, K. Improved deep metric learning with multi-class n-pair loss objective. *Advances in neural information processing systems*, 29, 2016.
- Swayamdipta, S., Schwartz, R., Lourie, N., Wang, Y., Hajishirzi, H., Smith, N. A., and Choi, Y. Dataset cartography: Mapping and diagnosing datasets with training dynamics. *arXiv preprint arXiv:2009.10795*, 2020.
- Thomas, C. and Kovashka, A. Preserving semantic neighborhoods for robust cross-modal retrieval. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVIII 16*, pp. 317–335. Springer, 2020.
- Thomas, C. and Kovashka, A. Emphasizing complementary samples for non-literal cross-modal retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4632–4641, 2022.
- Tiu, E., Talius, E., Patel, P., Langlotz, C. P., Ng, A. Y., and Rajpurkar, P. Expert-level detection of pathologies from unannotated chest x-ray images via self-supervised learning. *Nature Biomedical Engineering*, 6(12):1399–1406, 2022.
- Vasudevan, V., Caine, B., Gontijo Lopes, R., Fridovich-Keil, S., and Roelofs, R. When does dough become a bagel? analyzing the remaining mistakes on imagenet. *Advances in Neural Information Processing Systems*, 35: 6720–6734, 2022.
- Vedantam, R., Lawrence Zitnick, C., and Parikh, D. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4566–4575, 2015.
- Wang, J., Yang, Z., Hu, X., Li, L., Lin, K., Gan, Z., Liu, Z., Liu, C., and Wang, L. Git: A generative image-to-text transformer for vision and language. *Transactions on Machine Learning Research*, 2022a.
- Wang, Y., Xu, J., and Sun, Y. End-to-end transformer based model for image captioning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 2585–2594, 2022b.
- Wei, J., Zhu, Z., Cheng, H., Liu, T., Niu, G., and Liu, Y. Learning with noisy labels revisited: A study using real-world human annotations. *arXiv preprint arXiv:2110.12088*, 2021.
- Wu, P., Zheng, S., Goswami, M., Metaxas, D., and Chen, C. A topological filter for learning with label noise. *Advances in neural information processing systems*, 33: 21382–21393, 2020.
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., and Bengio, Y. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pp. 2048–2057. PMLR, 2015.
- Young, P., Lai, A., Hodosh, M., and Hockenmaier, J. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014.
- Zhai, X., Puigcerver, J., Kolesnikov, A., Ruysen, P., Riquelme, C., Lucic, M., Djolonga, J., Pinto, A. S., Neumann, M., Dosovitskiy, A., et al. The visual task adaptation benchmark. 2019.
- Zhang, S., Xu, Y., Usuyama, N., Bagga, J., Tinn, R., Preston, S., Rao, R., Wei, M., Valluri, N., Wong, C., et al. Large-scale domain-specific pretraining for biomedical vision-language processing. *arXiv preprint arXiv:2303.00915*, 2(3):6, 2023a.
- Zhang, S., Xu, Y., Usuyama, N., Xu, H., Bagga, J., Tinn, R., Preston, S., Rao, R., Wei, M., Valluri, N., et al. Biomed-clip: a multimodal biomedical foundation model pre-trained from fifteen million scientific image-text pairs. *arXiv preprint arXiv:2303.00915*, 2023b.

Zhu, Z., Dong, Z., and Liu, Y. Detecting corrupted labels without training a model to predict. In *International conference on machine learning*, pp. 27412–27427. PMLR, 2022.

Zhu, Z., Zhang, M., Wei, S., Wu, B., and Wu, B. Vdc: Versatile data cleanser based on visual-linguistic inconsistency by multimodal large language models. In *The Twelfth International Conference on Learning Representations*, 2024.

## A. Theoretical results

### A.1. Proof: Theorem 4.1

Suppose that  $\zeta_Y$  is distributed such that  $\text{supp}(k\zeta_Y(1-p)) \subseteq \{0, 1, \dots, k\}$ . For a correctly labeled sample  $(X, Y)$ , we have that  $k\zeta_Y(1-p)$  of the neighbors are relevant and have correct labels, and so each contribute  $d_{\mathcal{X}}(X, \bar{X})$  to  $S_m(X, Y)$ , and all remaining samples are either incorrectly labeled, or are not relevant to  $Y$ , and so each contribute  $d_{\mathcal{X}}(X, X')$ . Since  $S_m(X, Y)$  is the sum of iid Gaussians, it is also a Gaussian, with:

$$\begin{aligned} \mathbb{E}[S_m(X, Y)] &= \frac{1}{k} \left( \mathbb{E}[\mathbb{E}[d(X, \bar{X}_1) + \dots + d(X, \bar{X}_{k\zeta_Y(1-p)}) | \zeta] + \mathbb{E}[\mathbb{E}[d(X, X'_1) + \dots + d(X, X'_{k-k\zeta_Y(1-p)}) | \zeta]] \right) \\ &= \mathbb{E}[\zeta_Y](1-p)\mu_2 + (1 - \mathbb{E}[\zeta_Y](1-p))\mu_1 \\ &= \mathbb{E}[\zeta_Y](1-p)(\mu_2 - \mu_1) + \mu_1 \end{aligned}$$

$$\begin{aligned} \text{Var}[S_m(X, Y)] &= \mathbb{E}[\text{Var}(S_m(X, Y) | \zeta_Y)] + \text{Var}(\mathbb{E}[S_m(X, Y) | \zeta_Y]) \\ &= \mathbb{E}\left[\frac{1}{k^2} \text{Var}\left(d(X, \bar{X}_1) + \dots + d(X, \bar{X}_{k\zeta_Y(1-p)}) + d(X, X'_1) + \dots + d(X, X'_{k-k\zeta_Y(1-p)}) \mid \zeta_Y\right) \right. \\ &\quad \left. + \text{Var}(\mathbb{E}[S_m(X, Y) | \zeta_Y])\right] \\ &= \mathbb{E}\left[\frac{1}{k} (\zeta_Y(1-p)\sigma_2^2 + (1 - \zeta_Y(1-p))\sigma_1^2)\right] + \text{Var}(\zeta_Y(1-p)(\mu_2 - \mu_1) + \mu_1) \\ &= \frac{1}{k} (\mathbb{E}[\zeta_Y](1-p)\sigma_2^2 + (1 - \mathbb{E}[\zeta_Y](1-p))\sigma_1^2) + \text{Var}(\zeta_Y)(1-p)^2(\mu_2 - \mu_1)^2 \end{aligned}$$

Similarly,

$$S(X', Y') \sim \mathcal{N}(\mu_1, \frac{\sigma_1^2}{k})$$

Putting it all together:

$$\mathbb{P}(S_m(X', Y') - S_m(X, Y) > 0) = 1 - \Phi\left(\frac{-\mu}{\sigma}\right)$$

Where  $\mu = \mathbb{E}[\zeta_Y](1-p)(\mu_1 - \mu_2)$ ,  $\sigma = \sqrt{\frac{1}{k} (\mathbb{E}[\zeta_Y](1-p)\sigma_2^2 + (2 - \mathbb{E}[\zeta_Y](1-p))\sigma_1^2) + \text{Var}(\zeta_Y)(1-p)^2(\mu_2 - \mu_1)^2}$ , and  $\Phi$  is the Gaussian CDF. Note that  $\text{Var}(\zeta_Y)$  is finite as  $\zeta_Y$  is bounded by  $[0, 1]$ .

Setting  $\mu > 0$  gives Lemma 4.3.

### A.2. Empirically Validating Assumption 2

#### A.2.1. CLASSIFICATION AT DATASET LEVEL

To empirically validate Assumption 2, we first utilize the training sets from the original CIFAR-10 and CIFAR-100 datasets. As these are classification datasets, we naturally define  $\mathcal{J}$  as:  $x_2 \in \mathcal{J}(x_1) \iff y_1 = y_2$ , i.e. all images with the same label are paraphrases. We encode these images using the image encoder from OpenAI CLIP ViT-B/32 (Radford et al., 2021), and utilize the cosine distance as  $d_{\mathcal{X}}$ . We compute pairwise distance between all 40,000 samples, and categorize these distances into either  $x' \in \mathcal{J}(x)$  or  $x' \notin \mathcal{J}(x)$ . We plot a histogram of these distances in Figure A.1. Visually, both of these distributions appear to be normal, and we also observe that  $\mu_1 > \mu_2$  from Lemma 4.2. We then run a Shapiro–Wilk test on all four distributions to test for normality, randomly subsampling to 100 samples, as the Shapiro–Wilk test is not suitable for large sample sizes (Ghasemi & Zahediasl, 2012). We find that in all four cases, the null hypothesis cannot be rejected ( $p > 0.05$ ), and the test statistics are all greater than 0.97, indicating a high degree of normality.

#### A.2.2. CAPTIONING AT PER-SAMPLE LEVEL

We verify Assumption 2 at a sample-wise level, for the analogous  $\mathcal{H}(Y)$ , by conducting an experiment on captions from flickr30k. We select 20 random captions, then prompt Llama 3.1-8B-instruct (Grattafiori et al., 2024) to generate 50

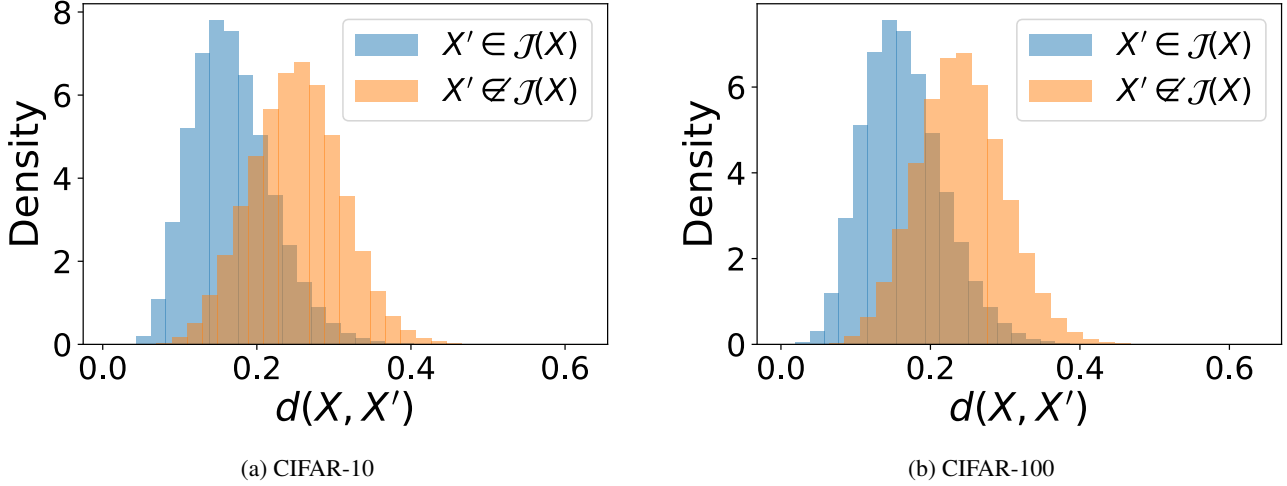


Figure A.1: Histogram of cosine distances in the CLIP image embedding space

paraphrasings of each caption (via sampling with temperature of 1), corresponding to 50 samples from  $\mathcal{H}(Y)$  for each caption. For the samples  $Y' \notin \mathcal{H}(Y)$ , we randomly select 50 other captions from the dataset. To match the support of the Gaussian, we take the distance function to be the log cosine distance (note that this does not change the ordering of the score across samples). We compute this distance using the text encoder from OpenAI CLIP ViT-B/32 (Radford et al., 2021), and plot histograms for each caption. We visualize our results in Figure A.2. Running the same Shapiro-Wilk test from Section A.2.1, we find that of the positive samples, 8/20 are Gaussian, and of the negative samples, 16/20 are Gaussian. Thus, there is some evidence the Gaussianity assumption holds for natural language and complex paraphrase functions.

### A.3. Second Theorem

We demonstrate that the embedding models trained via the contrastive multimodal objective are natural noisy label detectors.

**Theorem A.1** (Contrastive Multimodal Embedding Models Detect Noisy Labels). *Let  $\mathcal{Y} = \mathbb{R}$  and consider a training dataset  $\mathcal{D}$ . Suppose that  $\hat{h}_\theta^{\mathcal{X}} : \mathcal{X} \rightarrow \mathbb{R}^d$  is an embedding function, and  $\hat{h}_\theta^{\mathcal{Y}} : \mathcal{Y} \rightarrow \mathbb{R}^d$  is a Lipschitz continuous embedding function with constant  $L_{\mathcal{Y}} > 0$ , meaning that for all  $y, y' \in \mathcal{Y}$ ,*

$$\left\| \hat{h}_\theta^{\mathcal{Y}}(y) - \hat{h}_\theta^{\mathcal{Y}}(y') \right\|_2 \leq L_{\mathcal{Y}} |y - y'|.$$

*For an input  $x \in \mathcal{X}$  and its corresponding positive label  $y \in \mathcal{Y}$ , let  $\eta$  be a random variable drawn from a normal distribution:  $\eta \sim \mathcal{N}(0, \sigma^2)$ . Define a noisy label  $y' = y + \eta$ . Let  $d_{mm}(u, v) = \|u - v\|_2$ , which is proportional to  $\sqrt{d_{\cos}(u, v)}$  when  $\|u\|_2 = \|v\|_2 = 1$ . Then, with probability at least  $\delta(\epsilon) = 1 - 2\Phi\left(-\frac{\epsilon}{\sigma}\right)$ , where  $\epsilon > 0$  and  $\Phi$  is the cumulative distribution function of the standard normal distribution, the following inequality holds:*

$$d_{mm}\left(\hat{h}_\theta^{\mathcal{X}}(x), \hat{h}_\theta^{\mathcal{Y}}(y')\right) \geq d_{mm}\left(\hat{h}_\theta^{\mathcal{X}}(x), \hat{h}_\theta^{\mathcal{Y}}(y)\right) - L_{\mathcal{Y}} \epsilon.$$

When  $L_{\mathcal{Y}}$  is small, this means that the score for the mislabeled sample cannot be much lower than the score for the positive pair with high probability. Thus, we can see that multimodal embeddings are inherently capable of detecting mislabeled pairs, ensuring the distance between the embeddings of positive pairs is smaller than that of negative pairs. This motivates the use of  $d_{mm}$  in LEMON and in prior work (Kang et al., 2023; Liang et al., 2023).

**Proof:** Since  $\hat{h}_\theta^{\mathcal{Y}}$  is Lipschitz continuous with constant  $L_{\mathcal{Y}}$ , for any  $y, y' \in \mathcal{Y}$ , we have:

$$\left\| \hat{h}_\theta^{\mathcal{Y}}(y') - \hat{h}_\theta^{\mathcal{Y}}(y) \right\|_2 \leq L_{\mathcal{Y}} |y' - y| = L_{\mathcal{Y}} |\eta| \quad (4)$$

Let  $d_{mm}(u, v) = \|u - v\|_2$  be the Euclidean distance. Note that when  $\|u\|_2 = \|v\|_2 = 1$  (as in our experiments), we have that  $\|u - v\|_2 = \sqrt{2(1 - u^T v)} = \sqrt{2d_{\cos}(u, v)}$ , and so the two distances provide the same ordering of scores. Applying



Figure A.2: Histogram of log-cosine distances in the CLIP text embedding space

the triangle inequality, we get:

$$d_{mm}(\hat{h}_\theta^{\mathcal{X}}(x), \hat{h}_\theta^{\mathcal{Y}}(y')) \geq d_{mm}(\hat{h}_\theta^{\mathcal{X}}(x), \hat{h}_\theta^{\mathcal{Y}}(y)) - \|\hat{h}_\theta^{\mathcal{Y}}(y) - \hat{h}_\theta^{\mathcal{Y}}(y')\|_2.$$

When  $|\eta| \leq \epsilon$ , and substituting from Equation (4), it follows that:

$$d_{mm}(\hat{h}_\theta^{\mathcal{X}}(x), \hat{h}_\theta^{\mathcal{Y}}(y')) \geq d_{mm}(\hat{h}_\theta^{\mathcal{X}}(x), \hat{h}_\theta^{\mathcal{Y}}(y)) - L_{\mathcal{Y}}\epsilon$$

Since  $\eta \sim \mathcal{N}(0, \sigma^2)$ , the probability that  $|\eta| \leq \epsilon$  is:

$$P(|\eta| \leq \epsilon) = 1 - 2\Phi\left(-\frac{\epsilon}{\sigma}\right) = \delta(\epsilon),$$

where  $\Phi$  is the cumulative distribution function of the standard normal distribution.

Thus, with probability at least  $\delta(\epsilon)$ , we have:

$$d_{mm}(\hat{h}_\theta^{\mathcal{X}}(x), \hat{h}_\theta^{\mathcal{Y}}(y')) \geq d_{mm}(\hat{h}_\theta^{\mathcal{X}}(x), \hat{h}_\theta^{\mathcal{Y}}(y)) - L_{\mathcal{Y}}\epsilon$$

When  $L_{\mathcal{Y}}$  is small, this means that the score for the mislabeled sample cannot be much lower than the score for the positive pair with high probability.

## B. Comparison with Thomas & Kovashka (2022)

The goal of Thomas & Kovashka (2022) to identify samples with semantic diversity, which is different from our goal of identifying mislabeled examples. As such, their proposed scores (i.e.  $\Upsilon^{DIS}$  and  $\Upsilon^{DIV}$ ) may not be effective in identifying mislabeled samples. As an example, consider the score  $\Upsilon_Y^{DIS}$ , which computes the similarity between the original caption, and the captions of its second-degree neighbors in text-space. Given a particular caption, e.g. ‘‘This is a plane from the front view’’ in Figure 2, it could have second-degree neighbors in text-space that are semantically very similar to this caption (e.g. ‘‘A plane facing the viewer’’). However, only computing the distance of these captions in text space does not provide any signal for whether the *image* is correctly paired to the caption. Similarly, the  $\Upsilon^{DIV}$  scores also would not necessarily work, as the closeness of neighbors to each other in either modality do not provide a signal for whether the original sample is mislabeled.

However, the score from Thomas & Kovashka (2022) that would intuitively provide a signal for mislabeling is  $\Upsilon_X^{DIS}$ , which computes second-degree neighbors in text space, then examines similarity between images. This is essentially the sum over  $d_{\mathcal{X}}(\mathbf{x}, \mathbf{x}_{m_j})$  terms in our Equation (3), but using second-degree neighbors instead of nearest neighbors. In addition, our Equation (3) contains two additional weighting terms (which we show improve label error performance in our ablation experiments). Finally, our proposed score contains the sum of two additional terms, which are not explored in Thomas & Kovashka (2022).

We compare the performance of our method against the  $\Upsilon_X^{DIS}$  score in the main paper, and show performance of all four individual scores and two combined scores from Thomas & Kovashka (2022) in Appendix I.8.

## C. LEMON Algorithm

---

**Algorithm 1: LEMON: Label Error Detection Using Multimodal Neighbors**


---

**Input:** Dataset  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ , Multimodal encoders  $h_\theta^{\mathcal{X}}, h_\theta^{\mathcal{Y}}$ , Distance functions  $d_{\mathcal{X}}, d_{\mathcal{Y}}$

Hyperparameters:  $k, \beta, \gamma, \tau_{1,n}, \tau_{2,n}, \tau_{1,m}, \tau_{2,m}$

**Output:** Scores  $\{s_i\}_{i=1}^N$

- 1 Cache embeddings  $h_\theta^{\mathcal{X}}(\mathbf{x}_i)$  and  $h_\theta^{\mathcal{Y}}(\mathbf{y}_i)$  for  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}$ ;
  - 2 Cache  $d_{mm}(\mathbf{x}_i, \mathbf{y}_i) = 1 - \frac{h_\theta^{\mathcal{X}}(\mathbf{x}_i) \cdot h_\theta^{\mathcal{Y}}(\mathbf{y}_i)}{\|h_\theta^{\mathcal{X}}(\mathbf{x}_i)\|_2 \|h_\theta^{\mathcal{Y}}(\mathbf{y}_i)\|_2}$  for  $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}$ ;
  - 3 **for**  $i = 1$  **to**  $N$  **do**
  - 4     Find indices  $\{n_j\}_{j=1}^k$  of  $k$  nearest neighbors of  $\mathbf{x}_i$  from  $\mathcal{D} \setminus \{(\mathbf{x}_i, \mathbf{y}_i)\}$  using  $d_{\mathcal{X}}$ ; //  $d_{\mathcal{X}}$  can use cached  $h_\theta^{\mathcal{X}}$
  - 5     Find indices  $\{m_j\}_{j=1}^k$  of  $k$  nearest neighbors of  $\mathbf{y}_i$  from  $\mathcal{D} \setminus \{(\mathbf{x}_i, \mathbf{y}_i)\}$  using  $d_{\mathcal{Y}}$ ; //  $d_{\mathcal{Y}}$  can use cached  $h_\theta^{\mathcal{Y}}$
  - 6     Compute  $s_{n,i} := \frac{1}{k} \sum_{j=1}^k d_{\mathcal{Y}}(\mathbf{y}_i, \mathbf{y}_{n_j}) e^{-\tau_{1,n} d_{\mathcal{X}}(\mathbf{x}_i, \mathbf{x}_{n_j})} e^{-\tau_{2,n} d_{mm}(\mathbf{x}_{n_j}, \mathbf{y}_{n_j})}$ ;
  - 7     Compute  $s_{m,i} := \frac{1}{k} \sum_{j=1}^k d_{\mathcal{X}}(\mathbf{x}_i, \mathbf{x}_{m_j}) e^{-\tau_{1,m} d_{\mathcal{Y}}(\mathbf{y}_i, \mathbf{y}_{m_j})} e^{-\tau_{2,m} d_{mm}(\mathbf{x}_{m_j}, \mathbf{y}_{m_j})}$ ;
  - 8      $s_i := d_{mm}(\mathbf{x}_i, \mathbf{y}_i) + \beta s_{n,i} + \gamma s_{m,i}$
  - 9 **return**  $s$ ;
- 

For each image-caption pair in the dataset, we first compute how similar the image and caption are to each other using a pre-trained CLIP model ( $d_{mm}$ ), which gives a basic measure of how well they match. To compute  $s_m$ , we compute the nearest neighbors of the caption among other captions in the dataset. For each neighbor, we look at how similar their corresponding image is to the original image. The intuition is that if a sample is correctly labeled, the image should be similar to images of other samples with similar captions. We weight each neighbor based on how close it is to our original sample and how well-matched the neighboring pairs themselves are. Finally, we repeat this for nearest neighbors in the image space to get  $s_n$ . LEMON is then the weighted sum of these three scores.

Table C.1: Notation and definitions used in Section 3.

Symbol/Notation	Meaning
$\mathcal{D}$	Dataset consisting of samples $(\mathbf{x}, \mathbf{y})_{i=1}^N$
$\mathbf{x}, \mathcal{X}$	First modality and its corresponding space (e.g., images)
$\mathbf{y}, \mathcal{Y}$	Second modality and its corresponding space (e.g., text)
$f^*$	Oracle function that assigns a binary mislabel indicator $z_i$
$z_i$	Mislabel indicator for sample $i$ ( $z_i = 1$ if mislabeled, $z_i = 0$ otherwise)
$f(\mathbf{x}, \mathbf{y}) = s$	Model output score
$d_{\mathcal{X}}, d_{\mathcal{Y}}$	Distance functions in $\mathcal{X}$ and $\mathcal{Y}$ spaces
$B(\mathbf{x}, r)$	Ball of radius $r$ centered at $\mathbf{x}$ in $\mathcal{X}$ space
$B(\mathbf{y}, r)$	Ball of radius $r$ centered at $\mathbf{y}$ in $\mathcal{Y}$ space
$r_k(\mathbf{x})$	Radius such that the ball $B(\mathbf{x}, r)$ contains at least $k$ neighbors
$\mathbf{x}_{n_j}$	Nearest neighbor $j$ in $\mathcal{X}$ space
$\mathbf{y}_{m_j}$	Nearest neighbor $j$ in $\mathcal{Y}$ space
$h_\theta = (h_\theta^{\mathcal{X}}, h_\theta^{\mathcal{Y}})$	Multimodal encoder mapping $\mathcal{X}$ and $\mathcal{Y}$ to $\mathbb{R}^d$
$d_{mm}(\mathbf{x}, \mathbf{y})$	Multimodal distance between $\mathbf{x}$ and $\mathbf{y}$
$s_n(\mathbf{x}, \mathbf{y}, \mathcal{D})$	Score component based on $\mathbf{x}$ 's neighbors, see Equation (2).
$s_m(\mathbf{x}, \mathbf{y}, \mathcal{D})$	Score component based on $\mathbf{y}$ 's neighbors, see Equation (3).
$\beta, \gamma$	Hyperparameters weighting $s_n$ and $s_m$
$\tau_{1,n}, \tau_{2,n}, \tau_{1,m}, \tau_{2,m}$	Hyperparameters for weighting terms in $s_n$ and $s_m$
$k$	Number of nearest neighbors

## D. Data Processing

### D.1. Classification

We utilize CIFAR10N (`cifar10`) and CIFAR100N (`cifar100`) object detection (Zhu et al., 2022) datasets for all classification-based experiments. Each image is associated with a label indicating the primary object present in the image. These datasets contain 50,000 image-label pairs, with a clean and noisy label available per image. The noisy labels are examples of real human errors within the dataset. Further, we also generate synthetically noised labels as described in the main text. All images are resized to 224x224, center cropped, and normalized using mean and standard deviations corresponding to CLIP during the pre-processing stage. These two datasets are released under the Creative Commons Attribution-NonCommercial 4.0 license.

For `miniImageNet` and `stanfordCars`, we use the “red” datasets from Jiang et al. (2020), which contain noise from real-world web annotators. We split the full dataset (containing all annotations) into 75%/12.5%/12.5% train/val/test sets, stratifying by the mislabel flag. The annotations are licensed by Google under CC BY 4.0 license, and the images are under CC BY 2.0 license.

To generate the “text” modality for these classification datasets, we utilize the label name correspond to each class. For example, class 0 in `cifar10` is “airplane”, and this is the caption associated we associate with all images of that class. In contrast to the caption-based datasets, there will be multiple k-nearest neighbors in the text modality with zero distance (i.e., with the same class label).

### D.2. Captioning

We preprocess MSCOCO (Lin et al., 2014) and Flickr30k (Young et al., 2014) by using the Karpathy split (Karpathy & Fei-Fei, 2015), and then selecting one random annotation from the ones available. For the MMIMDB dataset (Arevalo et al., 2017), we utilize the plot outline as the text, and use the dataset splits provided. For MIMIC-CXR (Johnson et al., 2019), we use all images in the database and the provided data splits, and extract the findings and impression sections from the radiology note for the text modality. Images were normalized and transformed using the same procedure described above.

For downstream captioning, we use the pre-trained tokenizer and image processor corresponding to the pre-trained model (GIT (Wang et al., 2022a)) to pre-process image and captions.

Note that `flickr30k` is available under Flickr terms of use for non-commercial research and/or educational purposes<sup>4</sup>. `mscoco` is available under Creative Commons Attribution 4.0 License. `mmimdb` is available for personal and non-commercial use<sup>5</sup>. Finally, `mimiccxr` is available under the PhysioNet Credentialed Health Data License 1.5.0<sup>6</sup>.

## E. Baseline Methods

### E.1. Classification

TRAINING-DEPENDENT

**AUM** (Pleiss et al., 2020): This model assumes access to a classifier that can predict the class that an image likely belongs to. Then, the margin of difference between the prediction probability from the trained classifier for the assigned class and the class with the (next) highest probability is computed and averaged over training epochs. This score is thresholded to identify potential label errors.

**Datamap** (Swayamdipta et al., 2020): Similar to AUM, this method requires training a classifier on the dataset. In this baseline, it is assumed that instances with label errors are ‘hard to learn’, and thus low confidence in prediction throughout training epochs. To produce a single score, we combine the mean and standard deviation of the probability associated with the assigned class into a single score<sup>7</sup>.

**Confident Learning** (Northcutt et al., 2021a) is designed to identify labeling errors in classification datasets by modeling

<sup>4</sup><https://shannon.cs.illinois.edu/DenotationGraph/>

<sup>5</sup><https://developer.imdb.com/non-commercial-datasets/>

<sup>6</sup><https://physionet.org/content/mimic-cxr/view-license/2.0.0/>

<sup>7</sup>We experimented with different strategies, and the square root of the product of the mean and (1-standard deviation) and (1-mean) and standard deviation led to comparable, high validation F1 scores.

the relationship between true class labels and noisy ones. It sets thresholds for each true-noisy label pair. Using these thresholds, the model employs predicted class probabilities to rank predictions for each class, filtering out the noisy data.

#### TRAINING-FREE

**CLIP Logits** (Liang et al., 2023): CLIP is used as a zero-shot classifier to obtain the softmax-based probability for the assigned class. This value is then thresholded to identify label errors. Recently, (Feng et al.) used a similar zero-shot prediction jointly with a semi-supervised training approach for learning in the presence of label noise.

**CLIP Similarity** (Kang et al., 2023): The distance (either euclidean or cosine) between image and text embeddings from CLIP are computed and thresholded.

**Deep k-NN** (Bahri et al., 2020): The proportion of  $k$  nearest neighbors<sup>8</sup> with the same label is computed for each image of interest. Prior works have utilized different representations for obtaining neighbors, including logits and representations from pre-trained (Zhu et al., 2022) vision models. Here, we use pre-trained representations from CLIP.

**SimiFeat** (Zhu et al., 2022) uses nearby features to detect noisy labels under the assumption that local groups of features share clean or noisy labels. **SimiFeat-V** uses local voting and **SimiFeat-R** leverages ranking to detect noisy labels based on HOC estimator. The binary outputs produced are used for all score computations. Note that the difference between Simifeat-V and deep k-NN is in the data processing and augmentation.

**Discrepancy** (Thomas & Kovashka, 2022) finds second-degree nearest neighbors in the text space, then computes the average distance of these neighbors to the original sample in image space. We utilize the same CLIP model to compute semantic distance here as in LEMON.

Note that AUM and Datamap use the entire dataset as the reference set for label error detection, whereas SimiFeat uses both the train and test sets as reference sets, following the implementation in their original papers (Pleiss et al., 2020; Swayamdipta et al., 2020; Zhu et al., 2022).

## E.2. Captioning

#### PRE-TRAINED OR SUPERVISED

**LLaVA** (Liu et al., 2024): We prompt LLaVA (v1.6-vicuna-13b) with the following prompt: The proposed caption for this image is "{ }". Is this caption correct? Only answer with "Yes" or "No". We examine the probability distribution over the first non-special token, and find the likelihood of the token with the highest probability. If the corresponding token in lower case starts with "yes", we return 1— this probability as the mislabel score. Otherwise, we return the probability.

**VDC** (Zhu et al., 2024): As the original VDC paper does not explore the captioning setting, we make the following modifications to adapt it to our setup:

- As we only utilize open-source models in our method, to preserve fairness, we also implement VDC with open-source models. Specifically, we use Llama-3.1-8B-Instruct (Dubey et al., 2024) for the LLM in the Visual Question Generation (VQG) and Visual Answer Evaluation (VAE) stages (note that the VDC paper uses the OpenAI API), and InstructBLIP-Vicuna-7b (Dai et al., 2023) as the VLLM in the Visual Question Answering (VQA) stage (as in the VDC paper).
- In the VQG stage, instead of generating specific questions for each class, we generate six specific questions for each caption. We slightly modify the VQG prompt (Table 8 in the VDC paper) to adapt it to captioning, and omit providing the label set, as the set of all possible captions is infinite. We keep the two general questions used in the VDC paper.

**CapFilt (oracle-like)**: We generate predictions using a pre-trained model trained to distinguish between high-quality MSCOCO and noisy synthetic captions (Li et al., 2022). This forms an oracle-like, supervised baseline. We use a pre-trained checkpoint (pre-trained on the train split of MSCOCO)<sup>9</sup> for producing prediction scores.

<sup>8</sup>Note that this score can only take value in  $\{0, 1/k, 2/k, \dots, 1\}$ .

<sup>9</sup><https://huggingface.co/Salesforce/blip-itm-base-coco>

UNSUPERVISED

**Datamap:** We compute the causal language modeling loss across training epochs and compute the product of the mean and variance in loss across epochs. That is, we expect captioning loss for instances with label errors to be consistently high. We train captioning models for 3 epochs, with LoRA rank set to 4, and a maximum length of corresponding to maximum model length<sup>10</sup> for the finetuning task.

**Confident Learning:** We adapt this approach for dual-modality datasets, such as image-text pairs, by clustering text embeddings to serve as class labels for noise detection.

DOWNSTREAM-TASK UNAWARE

**Deep KNN:** We cluster captions similar to confident learning, adapting classification baseline.

**CLIP Similarity:** This is the same setup as classification.

**Discrepancy:** This is the same setup as classification.

## F. Compute Setup

We run our experiments on a shared Slurm cluster. Most experiments used one RTX A6000 with 48 GB VRAM, 10 CPU cores of Intel Xeon Ice Lake Platinum 8368, and 50 GB RAM.

## G. Hyperparameters in Label Error Detection

The hyperparameters in each case were selected based on the validation set F1-score. Note that LEMON<sub>FIX</sub> does not require hyperparameter tuning. Baseline code is included in the supplementary material. For SimiFeat-V and -R, we use the official open-sourced implementation directly.

### G.1. Classification

The search space for each method:

1. AUM, Datamap: learning rate  $\in \{5e - 5, 5e - 6\}$ , training for epochs  $\in \{5, 10\}$ <sup>11</sup>
2. Confident learning: learning rate  $\in \{5e - 6, 5e - 5\}$ , upto 30 epochs with early stopping with a patience of 10.
3. CLIP Sim.: cosine distance metric, no other hyperparameters
4. CLIP Zero shot: distance metric
5. Discrepancy:  $k \in \{1, 2, 5, 10, 15, 20, 30, 50\}$
6. deep k-NN:  $k$ , cosine distance metric
7. Simifeat: we set  $k = 10$  following the original paper (Zhu et al., 2022).

### G.2. Captioning

For most baselines requiring a class index—obtained by clustering captions—we set the number of clusters to be 100.

1. LLaVA: Small amount of prompt tuning. The optimal prompt selected was `The proposed caption for this image is ``{}``. Is this caption correct? Only answer with ``Yes`` or ``No``.`
2. Confident learning: learning rate  $\in \{5e - 6, 5e - 5\}$ , upto 30 epochs with early stopping with a patience of 10, number of clusters for captions<sup>12</sup>

<sup>10</sup>This is longer than captions in the train sets of all datasets except the medical dataset.

<sup>11</sup>Note that we experiment with training for fewer epochs to avoid memorization, following (Pleiss et al., 2020).

<sup>12</sup>For `mimiccxr`, we use 10 clusters.

3. Discrepancy:  $k \in \{1, 2, 5, 10, 15, 20, 30, 50\}$
4. deep k-NN:  $k \in \{1, 2, 5, 10, 15, 20, 30, 50\}$ , number of text clusters
5. VDC: no tunable hyperparameters

### G.3. Our Method

We search the following hyperparameters for our LEMON<sub>OPT</sub>:

1.  $k \in \{1, 2, 5, 10, 15, 20, 30, 50\}$
2. Distance metric (either cosine or euclidean)
3.  $\beta, \gamma, \tau_{1,n}, \tau_{2,n}, \tau_{1,m}, \tau_{2,m}$ : We take the hyperparameter set which achieves the best validation set F1 from these two strategies: (1) Using Scipy’s `minimize` function, with initial guess  $(1, 1, \dots, 1)$ , and with no explicit bounds. (2) Using a grid search with the following grid:
  - $\beta \in \{0, 5, 10, 15, \dots, 100\}$
  - $\gamma \in \{0, 5, 10, 15, \dots, 100\}$
  - $\tau_{1,n}, \tau_{2,n}, \tau_{1,m}, \tau_{2,m} \in \{0, 1, 5, 10\}$

### G.4. Optimal Hyperparameters

Optimal hyperparameters for classification datasets can be found in Table G.1, and optimal hyperparameters for captioning datasets can be found in Table G.2.

Table G.1: Optimal hyperparameters for methods shown in Table I.1. Note that Simifeat, VDC, CLIP Sim., and LEMON<sub>FIX</sub> have no tunable hyperparameters.

	cifar10	cifar100	miniImageNet	stanfordCars
<b>AUM</b>	LR = 5E-6 Epochs = 5	LR = 5E-5 Epochs = 5	LR = 5E-6 Epochs = 10	LR = 5E-5 Epochs = 10
<b>Datamap</b>	LR = 5E-6 Epochs = 5	LR = 5E-5 Epochs = 5	LR = 5E-5 Epochs = 10	LR = 5E-5 Epochs = 10
<b>Confident</b>	LR=5e-06 Epochs=30 Batch size=128	LR=5e-06 Epochs=30 Batch size=128	LR=5e-05 Epochs=30 Batch size=128	LR=5e-05 Epochs=30 Batch size=128
<b>CLIP Logits</b>	Cosine distance	Cosine distance	Cosine distance	Cosine distance
<b>Discrepancy</b>	k=20	k=50	k=30	k=20
<b>Deep k-NN</b>	k=50 cosine distance	k=20 cosine distance	k=50 cosine distance	k=30 cosine distance
<b>LEMON<sub>OPT</sub></b>	k=50 cosine distance $\beta = 20$ $\gamma = 35$ $\tau_{1,n} = 0$ $\tau_{2,n} = 5$ $\tau_{1,m} = 0$ $\tau_{2,m} = 5$	k=20 cosine distance $\beta = 2.14$ $\gamma = -0.024$ $\tau_{1,n} = -1.71$ $\tau_{2,n} = 4.85$ $\tau_{1,m} = -0.068$ $\tau_{2,m} = -0.019$	k=50 Euclidean distance $\beta = 0.664$ $\gamma = 0.395$ $\tau_{1,n} = 1.91$ $\tau_{2,n} = 1.04$ $\tau_{1,m} = 1.00$ $\tau_{2,m} = 1.35$	k=15 Euclidean distance $\beta = 0.631$ $\gamma = 0.431$ $\tau_{1,n} = 0.898$ $\tau_{2,n} = -0.192$ $\tau_{1,m} = 0.0$ $\tau_{2,m} = -0.001$

Table G.2: Optimal hyperparameters for methods shown in Table 3. Note that LLaVA, VDC, CLIP Sim. and LEMOn<sub>FIX</sub> have no tunable hyperparameters.

	flickr30k	mscoco	mmimdb	mimiccxr
<b>Datamap</b>	Batch size = 16	Batch size = 16	Batch size = 16	Batch size = 16
	Epochs = 3	Epochs = 3	Epochs = 3	Epochs = 3
	LoRA rank = 4	LoRA rank = 4	LoRA rank = 4	LoRA rank = 4
<b>Discrepancy</b>	k=5	k=10	k=10	k=10
<b>Deep k-NN</b>	k=50 n_cluster=100	k=50 n_cluster=100	k=20 n_cluster=100	k=50 n_cluster=100
<b>Confident</b>	LR=5e-06	LR=5e-06	LR=5e-05	LR=5e-06
	Epochs=30	Epochs=30	Epochs=30	Epochs=30
	Batch size=128	Batch size=128,	Batch size=128	Batch size=16
	n_cluster=10	n_cluster=100	n_cluster=10	n_cluster=10
<b>LEMOn<sub>OPT</sub></b>	k=30	k=30	k=10	k=30
	cosine distance	cosine distance	Euclidean distance	cosine distance
	$\beta = 0.092$	$\beta = 5.324$	$\beta = 1.001$	$\beta = 5$
	$\gamma = 0.177$	$\gamma = 11.057$	$\gamma = 1.202$	$\gamma = 10$
	$\tau_{1,n} = 0.274$	$\tau_{1,n} = 5.143$	$\tau_{1,n} = 0.983$	$\tau_{1,n} = 5$
	$\tau_{2,n} = 0.074$	$\tau_{2,n} = 10.498$	$\tau_{2,n} = 1.000$	$\tau_{2,n} = 10$
	$\tau_{1,m} = 0.072$	$\tau_{1,m} = 7.233$	$\tau_{1,m} = 4.450$	$\tau_{1,m} = 5$
	$\tau_{2,m} = 0.0$	$\tau_{2,m} = 15.637$	$\tau_{2,m} = 1.080$	$\tau_{2,m} = 10$

## H. Hyperparameters in Downstream Models

### H.1. Classification

We train a Vision Transformer (ViT)-based image classification (Dosovitskiy et al., 2020)<sup>13</sup> model pre-trained on ImageNet-21k (Ridnik et al., 2021) and fine-tuned on ImageNet 2012 (Russakovsky et al., 2015) with an additional linear layer. We add a linear layer above the last hidden state. We train for up to 30 epochs with an initial learning rate of and early stopping with a patience of 4. We tune learning rates on the fully noisy set, for learning rate in  $\in \{5e - 4, 1e - 3, 1e - 5\}$  (with cosine annealing learning rate scheduling).

### H.2. Captioning

The hyperparameter tuning grid for the captioning model<sup>14</sup> are: learning rate of  $1e - 4$  or  $1e - 5$ , batch size: 16, maximum number of epochs: 10. The model checkpoint from the epoch with the lowest validation loss is used for caption generation at test time. For text generation, we use beam search with 4 beams, following Wang et al. (2022a). We use the AdamW optimizer (Loshchilov & Hutter, 2018), with cosine scheduling for learning rate with 1000 warmup steps. We tune LoRA rank in  $\{4, 16\}$  based on BLEU-4 scores on the validation set. We also experimented with a lower learning rate of  $1e - 5$  on the fully noisy set, and observed higher performance (in terms of validation BLEU-4) with a higher learning rate of  $1e - 4$  on flickr30k.

Here, for finding best epoch during training, we assume that validation loss is correlated with language model quality. We leave strategies such as self-critical training validation (Rennie et al., 2017) to future work. We verified that retaining the model from the last epoch leads to similar trends in results – see Table I.9 for models trained up to 10 epochs. We find that both LEMOn and the baseline perform within 1 point of each other on both datasets (similar to the results with loss-based checkpointing), and filtering out noisy data improves downstream captioning performance. Note that pre-conference, we benchmarked results without loss masking of padded tokens during training and observed similar result trends.

<sup>13</sup><https://huggingface.co/google/vit-base-patch16-224>

<sup>14</sup><https://huggingface.co/microsoft/git-base>

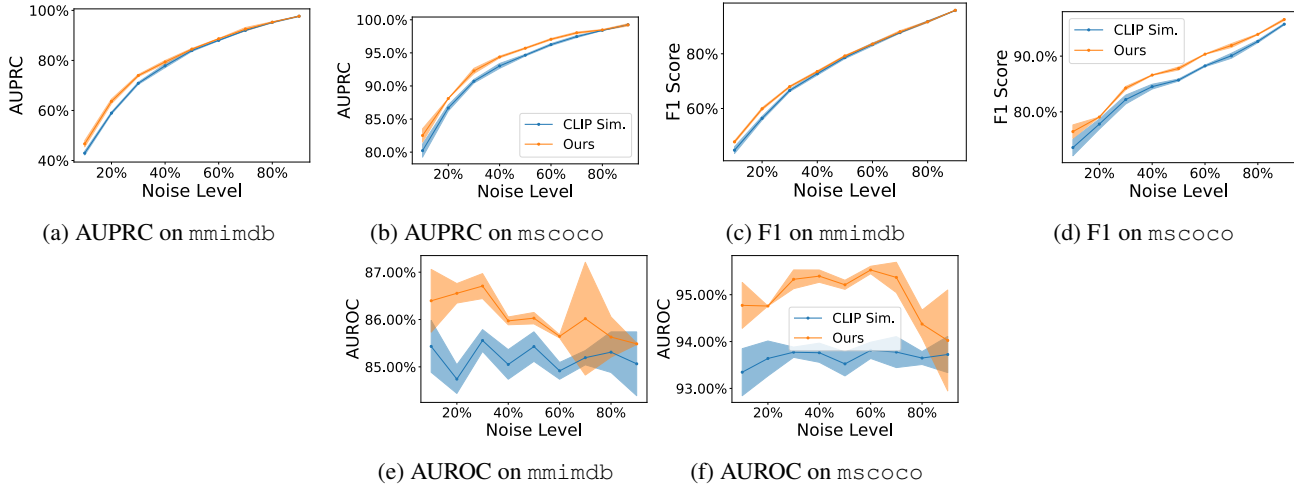


Figure I.1: Test-set performance of LEMON<sub>OPT</sub> compared to the CLIP Similarity for varying levels of the synthetic noise.

## I. Additional Experimental Results

### I.1. Label Error Detection in Classification Settings

Full results on classification datasets using the noise types bolded in Table 1 (including AUPRC) can be found in Table I.1.

The performance of all baselines and our method on the two types of synthetic errors are shown in Table I.2, all at a noise level of 40% (comparable to the amount of error in the noisy CIFAR datasets).

### I.2. Label Error Detection in Captioning Settings

Full results on classification datasets using the noise types bolded in Table 1 (including AUPRC) can be found in Table I.3.

Results on the remaining synthetic noise types (at 40%) can be found in: flickr30k: Table I.4, mscoco: Table I.5, mmimdb: Table I.6, and mimc-cxr: Table I.7. Across all datasets and noising types, we find that our model outperforms other non-oracle/supervised baselines.

### I.3. Varying Noise Level

We show the AUROC for varying noise levels in Figure I.1.

### I.4. Robustness to Hyperparameters

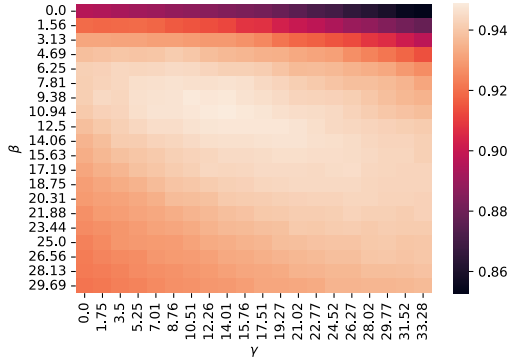
We show the test-set F1 of LEMON for varying  $\beta$  and  $\gamma$ , keeping all other hyperparameters at their fixed optimal values, in Figure I.2. In Table I.8, we show the performance of LEMON when hyperparameters are fixed (at  $k = 30$ , cosine distance,  $\beta = \gamma = 5$ ,  $\tau_{1,n} = \tau_{1,m} = 0.1$ , and  $\tau_{2,n} = \tau_{2,m} = 5$ ) versus when they are optimized using a labeled validation set. Note that F1 is not computed as it requires external information to select a threshold.

### I.5. Ablations of our Method

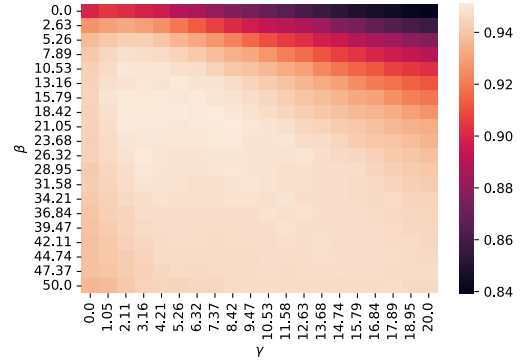
Ablations of our method can be found in Table I.10 and Table I.11.

### I.6. Runtime Comparison

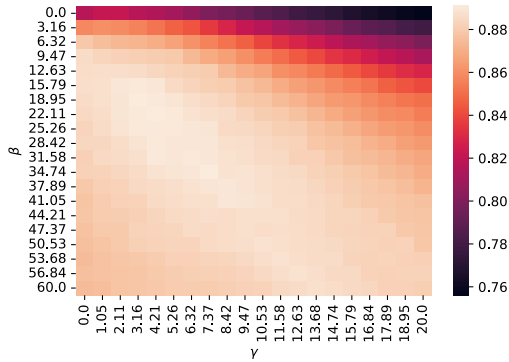
We compare the runtime of LEMON with baselines in Table I.12.



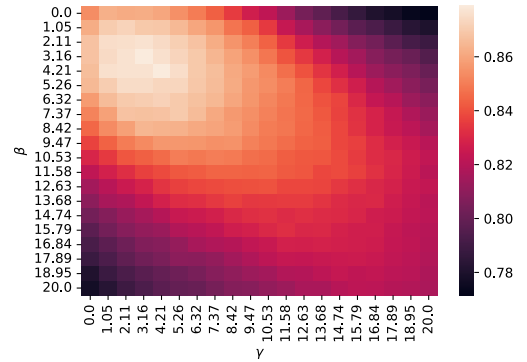
(a) cifar10, asymmetric noise



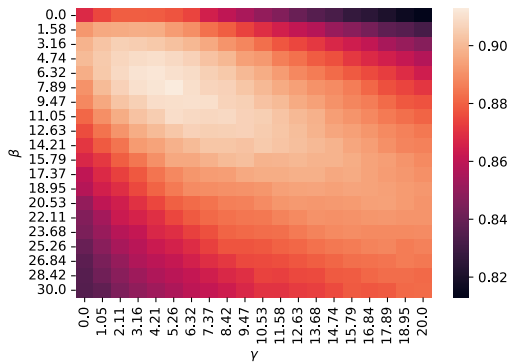
(b) cifar10, symmetric noise



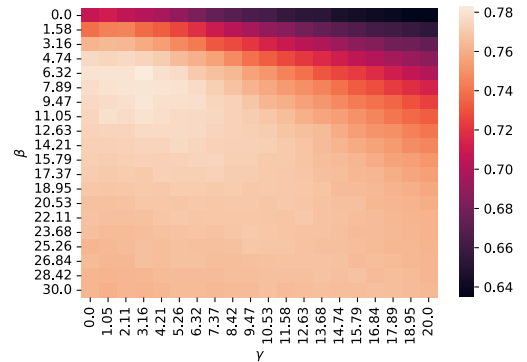
(c) cifar10, real noise



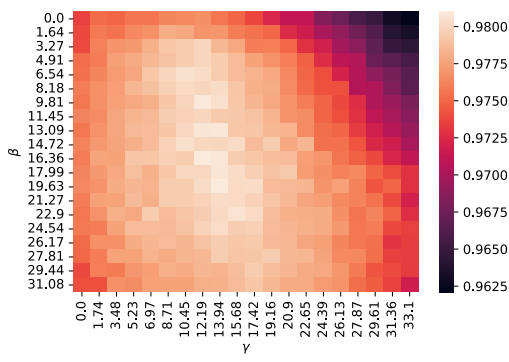
(d) cifar100, asymmetric noise



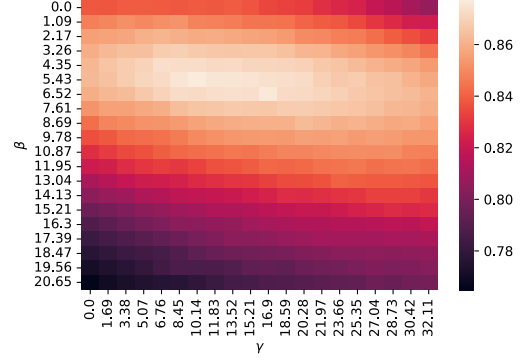
(e) cifar100, symmetric noise



(f) cifar100, real noise



(g) mscoco, random noise



(h) mscoco, cat noise

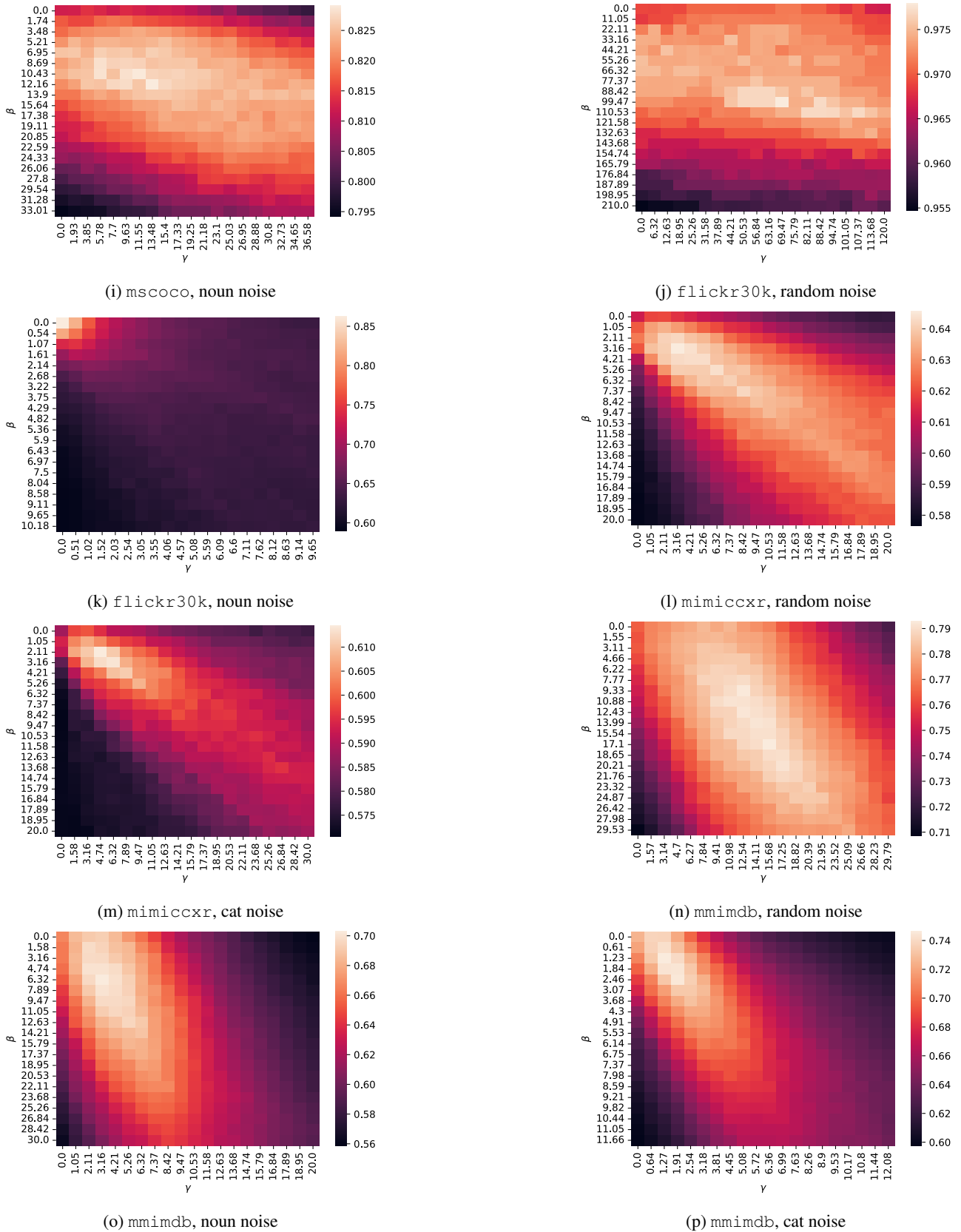


Figure I.2: F1 of our method for varying  $\beta$  and  $\gamma$ , keeping all other hyperparameters their fixed optimal values.

Table I.1: Label error detection performance across classification datasets, for the bolded noise types in Table 1. We separate AUM, Datamap, and Confident learning, as they require training a classifier from scratch. Bold denotes best score within each training approach.

Dataset	Method	Training-free	AUROC (%)	AUPRC (%)	F1 (%)
cifar10	AUM		<b>98.3</b> (0.1)	<b>97.9</b> (0.1)	<b>92.9</b> (0.1)
	Datamap	$\times$	98.2 (0.1)	97.6 (0.1)	92.2 (0.5)
	Confident		89.6 (1.4)	86.1 (1.8)	88.2 (1.7)
	CLIP Logits		95.5 (0.2)	93.9 (0.3)	86.2 (0.6)
	CLIP Sim.		92.2 (0.2)	90.2 (0.4)	82.3 (0.3)
	Simifeat-V		90.8 (0.3)	88.0 (0.7)	88.2 (0.5)
	Simifeat-R	$\checkmark$	88.8 (0.1)	85.1 (0.5)	85.9 (0.1)
	Discrepancy		77.1 (1.9)	70.4 (2.7)	66.4 (2.2)
	Deep k-NN		97.8 (0.1)	96.5 (0.2)	91.4 (0.6)
	LEMOn <sub>FIX</sub> (Ours)		97.7 (0.2)	96.8 (0.3)	90.9 (0.1)
LEMOn <sub>OPT</sub> (Ours)		<b>98.1</b> (0.0)	<b>97.4</b> (0.1)	<b>92.0</b> (0.2)	
cifar100	AUM		<b>92.3</b> (0.3)	<b>90.0</b> (0.5)	<b>81.1</b> (0.3)
	Datamap	$\times$	91.8 (0.3)	89.5 (0.5)	80.8 (0.5)
	Confident		78.6 (0.4)	68.8 (0.9)	73.7 (0.5)
	CLIP Logits		84.9 (0.7)	80.3 (1.2)	72.0 (0.9)
	CLIP Sim.		80.8 (0.9)	75.2 (1.3)	68.7 (1.1)
	Simifeat-V		79.5 (0.1)	71.1 (0.2)	73.3 (0.2)
	Simifeat-R	$\checkmark$	81.4 (0.1)	70.5 (0.5)	76.1 (0.1)
	Discrepancy		66.0 (1.5)	57.4 (2.3)	58.9 (0.8)
	Deep k-NN		87.4 (0.3)	77.9 (1.0)	75.7 (0.3)
	LEMOn <sub>FIX</sub> (Ours)		88.9 (0.7)	84.6 (1.1)	75.4 (0.6)
LEMOn <sub>OPT</sub> (Ours)		<b>90.8</b> (0.0)	<b>87.4</b> (0.3)	<b>78.4</b> (0.0)	
miniImageNet	AUM		83.1 (0.2)	73.2 (0.5)	68.3 (0.4)
	Datamap	$\times$	<b>85.1</b> (0.3)	<b>70.1</b> (0.8)	<b>70.6</b> (0.2)
	Confident		59.5 (0.7)	42.0 (0.8)	37.7 (1.5)
	CLIP Logits		<b>90.0</b> (0.2)	<b>80.9</b> (0.5)	<b>77.1</b> (0.2)
	CLIP Sim.		89.3 (0.2)	80.7 (0.3)	76.1 (0.4)
	Simifeat-V		68.1 (0.2)	52.9 (0.4)	54.8 (0.3)
	Simifeat-R	$\checkmark$	65.1 (0.3)	50.0 (0.5)	48.6 (0.5)
	Discrepancy		79.5 (0.2)	65.9 (0.4)	64.0 (0.1)
	Deep k-NN		83.2 (0.2)	70.9 (0.6)	68.5 (0.2)
	LEMOn <sub>FIX</sub> (Ours)		89.5 (0.2)	81.5 (0.3)	74.7 (0.2)
LEMOn <sub>OPT</sub> (Ours)		<b>90.0</b> (0.4)	79.7 (3.1)	76.9 (0.2)	
stanfordCars	AUM		70.4 (2.3)	<b>42.0</b> (1.0)	47.2 (3.1)
	Datamap	$\times$	<b>72.2</b> (1.7)	39.5 (0.2)	<b>50.4</b> (2.1)
	Confident		60.7 (0.3)	29.9 (0.2)	39.9 (0.6)
	CLIP Logits		68.8 (0.1)	39.7 (0.9)	47.3 (0.5)
	CLIP Sim.		69.8 (0.4)	40.7 (1.0)	46.6 (0.5)
	Simifeat-V		63.6 (1.3)	33.5 (1.4)	43.5 (1.6)
	Simifeat-R	$\checkmark$	66.8 (0.3)	34.9 (0.2)	47.9 (0.5)
	Discrepancy		65.7 (0.3)	33.3 (0.6)	44.3 (0.7)
	Deep k-NN		71.5 (0.6)	41.7 (0.7)	49.1 (0.6)
	LEMOn <sub>FIX</sub> (Ours)		72.6 (0.7)	<b>44.9</b> (1.4)	47.7 (2.0)
LEMOn <sub>OPT</sub> (Ours)		<b>73.1</b> (0.5)	40.5 (0.4)	<b>51.3</b> (0.5)	

Table I.2: Label error detection performance on synthetic errors for classification datasets

Dataset	Flip Type	Method	AUROC		AUPRC		F1	
			mean	std	mean	std	mean	std
cifar10	asymmetric	AUM	<b>93.7</b>	0.6	<b>86.6</b>	0.6	<b>86.6</b>	1.1
		Datamap	93.6	0.4	86.3	0.8	86.0	0.9
		Confident	87.3	4.7	78.0	7.3	84.7	5.7
		CLIP Logits	<b>98.8</b>	0.2	<b>97.9</b>	0.3	93.2	0.4
		CLIP Sim.	97.0	0.2	95.3	0.1	89.4	0.5
		Simifeat-V	70.0	0.6	58.7	0.8	60.7	1.1
		Simifeat-R	71.4	1.0	57.6	1.3	64.3	1.4
		Discrepancy	63.2	3.3	51.3	1.4	58.5	1.2
		Deep k-NN	85.2	0.7	66.2	0.9	79.5	0.8
	LEMON <sub>FIX</sub>	97.5	0.2	94.8	0.6	90.2	0.8	
	LEMON <sub>OPT</sub>	<b>98.8</b>	0.2	97.8	0.5	<b>93.9</b>	0.3	
	symmetric	AUM	<b>99.8</b>	0.0	<b>99.7</b>	0.0	<b>98.1</b>	0.3
		Datamap	<b>99.8</b>	0.0	<b>99.7</b>	0.0	98.0	0.2
		Confident	97.5	0.2	94.8	0.7	96.9	0.4
		CLIP Logits	98.5	0.0	97.9	0.1	92.2	0.1
		CLIP Sim.	97.1	0.1	95.9	0.2	89.5	0.1
		Simifeat-V	96.5	0.2	93.9	0.5	94.3	0.2
		Simifeat-R	95.4	0.1	91.8	0.2	93.1	0.2
Discrepancy		84.5	2.3	78.2	3.0	73.2	1.1	
Deep k-NN		99.2	0.1	98.1	0.2	96.1	0.3	
LEMON <sub>FIX</sub>	99.5	0.1	99.2	0.1	95.9	0.3		
LEMON <sub>OPT</sub>	<b>99.6</b>	0.1	<b>99.4</b>	0.1	<b>96.8</b>	0.2		
cifar100	asymmetric	AUM	<b>82.4</b>	2.0	<b>67.5</b>	2.6	<b>75.2</b>	1.5
		Datamap	74.4	1.0	59.3	1.2	66.1	1.5
		Confident	71.7	2.5	56.6	3.0	67.3	2.7
		CLIP Logits	96.6	0.3	94.8	0.5	88.3	0.7
		CLIP Sim.	95.1	0.4	93.1	0.5	85.6	0.6
		Simifeat-V	65.4	1.0	52.8	1.2	57.0	1.3
		Simifeat-R	60.1	0.5	45.9	0.7	56.2	1.3
		Discrepancy	62.0	0.6	50.8	1.2	58.1	0.9
		Deep k-NN	62.6	0.3	47.7	0.8	63.2	0.5
	LEMON <sub>FIX</sub>	94.9	0.3	92.1	0.4	84.3	0.4	
	LEMON <sub>OPT</sub>	<b>96.7</b>	0.3	<b>95.3</b>	0.4	<b>88.4</b>	0.5	
	symmetric	AUM	<b>99.2</b>	0.3	<b>98.9</b>	0.4	<b>95.2</b>	1.1
		Datamap	<b>99.2</b>	0.4	98.8	0.7	95.0	1.5
		Confident	87.8	0.6	77.8	1.0	85.2	0.5
		CLIP Logits	96.8	0.1	95.2	0.3	89.1	0.4
		CLIP Sim.	95.5	0.2	93.5	0.3	86.6	0.9
		Simifeat-V	91.0	0.7	84.8	1.4	84.5	1.0
		Simifeat-R	90.5	0.3	79.8	1.1	84.1	0.9
Discrepancy		82.8	0.9	75.8	1.0	70.7	0.5	
Deep k-NN		96.7	0.1	91.7	0.3	91.0	0.4	
LEMON <sub>FIX</sub>	98.4	0.1	97.7	0.2	92.0	0.1		
LEMON <sub>OPT</sub>	<b>99.0</b>	0.0	<b>98.7</b>	0.1	<b>94.1</b>	0.1		

Table I.3: Label error detection performance on captioning datasets, for the bolded noise types in Table 1.

Dataset	Method	AUROC (%)	AUPRC (%)	F1 (%)
flickr30k	LLaVA	79.3 (0.8)	58.5 (0.2)	65.0 (1.1)
	Datamap	80.6 (1.3)	63.3 (1.5)	66.6 (0.6)
	Discrepancy	73.0 (0.6)	59.2 (1.8)	59.0 (0.3)
	VDC	92.9 (1.1)	87.2 (0.3)	81.1 (1.6)
	Deep k-NN	71.1 (0.4)	52.0 (1.0)	59.2 (0.8)
	Confident	63.1 (0.9)	42.1 (1.2)	54.0 (0.9)
	CLIP Sim.	<b>94.8</b> (0.5)	<b>92.8</b> (0.5)	<b>84.2</b> (0.9)
	LEMOn <sub>FIX</sub> (Ours)	93.6 (0.2)	92.0 (0.2)	-
	LEMOn <sub>OPT</sub> (Ours)	94.5 (0.2)	<b>92.8</b> (0.3)	83.6 (1.4)
	CapFilt (Oracle)	98.6 (0.1)	98.1 (0.1)	93.1 (0.7)
mscoco	LLaVA	80.3 (0.1)	63.4 (0.3)	74.9 (0.3)
	Datamap	85.8 (0.4)	80.0 (0.6)	74.3 (0.3)
	Discrepancy	72.7 (0.3)	67.2 (0.4)	62.5 (0.3)
	VDC	94.1 (0.2)	91.8 (0.2)	86.3 (0.4)
	Deep k-NN	76.6 (0.4)	70.3 (0.6)	67.5 (0.8)
	Confident	71.5 (0.5)	56.4 (0.5)	66.5 (0.5)
	CLIP Sim.	93.8 (0.2)	93.0 (0.4)	84.5 (0.4)
	LEMOn <sub>FIX</sub> (Ours)	92.0 (0.1)	91.8 (0.3)	-
	LEMOn <sub>OPT</sub> (Ours)	<b>95.6</b> (0.2)	<b>94.6</b> (0.3)	<b>87.0</b> (0.2)
	CapFilt (Oracle)	99.3 (0.0)	99.1 (0.0)	95.4 (0.4)
mmimdb	LLaVA	58.4 (0.2)	46.4 (0.2)	58.5 (0.1)
	Discrepancy	57.8 (0.4)	46.1 (0.9)	57.4 (0.2)
	VDC	80.5 (0.3)	67.1 (0.3)	69.3 (0.6)
	Datamap	55.4 (0.5)	44.1 (0.6)	57.2 (0.0)
	deep k-NN	61.2 (0.4)	47.2 (0.5)	58.3 (0.4)
	Confident	52.8 (1.1)	41.4 (0.6)	51.8 (1.8)
	CLIP Sim.	85.1 (0.3)	77.8 (0.7)	72.7 (0.6)
	LEMOn <sub>FIX</sub> (Ours)	84.3 (0.3)	77.7 (0.8)	-
	LEMOn <sub>OPT</sub> (Ours)	<b>86.0</b> (0.1)	<b>79.4</b> (0.6)	<b>73.5</b> (0.3)
	CapFilt	82.7 (0.7)	73.3 (1.2)	71.3 (0.3)
mimiccxr	LLaVA	53.9 (0.5)	42.7 (0.7)	57.0 (0.1)
	Datamap	51.0 (0.5)	40.5 (0.7)	57.0 (0.1)
	Discrepancy	60.0 (0.7)	50.2 (0.5)	57.2 (0.1)
	VDC	50.8 (0.4)	40.3 (0.2)	57.0 (0.1)
	deep k-NN	62.9 (0.4)	48.0 (0.3)	59.2 (0.1)
	Confident	61.8 (0.3)	47.0 (0.2)	58.1 (0.6)
	CLIP Sim.	64.1 (0.4)	51.7 (0.5)	59.2 (0.0)
	LEMOn <sub>FIX</sub> (Ours)	66.3 (0.4)	54.6 (0.5)	-
	LEMOn <sub>OPT</sub> (Ours)	<b>70.4</b> (1.6)	<b>60.4</b> (1.6)	<b>61.1</b> (0.8)
	CapFilt	49.7 (0.3)	40.0 (0.2)	57.0 (0.0)

Table I.4: flickr30k: Label Error Detection

Dataset	Noise Type	Method	AUROC		AUPRC		F1	
			mean	std	mean	std	mean	std
flickr30k	noun	LLaVA	79.3	0.8	58.5	0.2	65.0	1.1
		Datamap	80.6	1.3	63.3	1.5	66.6	0.6
		Discrepancy	73.0	0.6	59.2	1.8	59.0	0.3
		VDC	92.9	1.1	87.2	0.3	81.1	1.6
		Deep kNN	71.1	0.4	52.0	1.0	59.2	0.8
		Confident	63.1	0.9	42.1	1.2	54.0	0.9
		CLIP Sim.	94.8	0.5	92.8	0.5	84.2	0.9
		LEMON <sub>FIX</sub>	93.6	0.2	92.0	0.2	83.4	0.6
		LEMON <sub>OPT</sub>	94.5	0.2	92.8	0.3	83.6	1.4
		CapFilt	98.6	0.1	98.1	0.1	93.1	0.7
	random	LLaVA	81.3	1.0	65.6	1.4	72.2	1.1
		Datamap	90.5	0.8	83.9	0.6	78.7	0.6
		Discrepancy	83.8	1.2	75.3	2.3	72.2	1.4
		VDC	98.2	0.2	96.4	0.3	92.1	0.6
		Deep kNN	81.1	1.6	65.3	1.7	72.7	1.2
		Confident	71.2	0.6	55.3	0.5	67.1	0.8
		CLIP Sim.	99.5	0.1	99.3	0.1	95.7	0.5
		LEMON <sub>FIX</sub>	99.4	0.2	99.3	0.2	96.3	0.7
		LEMON <sub>OPT</sub>	99.5	0.2	99.4	0.3	96.3	1.0
		CapFilt	99.9	0.0	99.8	0.0	97.9	0.2

Table I.5: mscoco: Label Error Detection

Dataset	Noise Type	Method	AUROC		AUPRC		F1	
			mean	std	mean	std	mean	std
mscoco	cat	LLaVA	80.3	0.1	63.4	0.3	74.9	0.3
		Datamap	85.8	0.4	80.0	0.6	74.3	0.3
		Discrepancy	72.7	0.3	67.2	0.4	62.5	0.3
		VDC	94.1	0.2	91.8	0.2	86.3	0.4
		Deep kNN	76.6	0.4	70.3	0.6	67.5	0.8
		Confident	71.5	0.5	56.4	0.5	66.5	0.5
		CLIP Sim.	93.8	0.2	93.0	0.4	84.5	0.4
		LEMON <sub>FIX</sub>	92.0	0.1	91.8	0.3	82.8	0.4
		LEMON <sub>OPT</sub>	95.6	0.2	94.6	0.3	87.0	0.2
		CapFilt	99.3	0.0	99.1	0.0	95.4	0.4
	noun	LLaVA	79.4	0.2	61.3	0.3	72.6	0.2
		Datamap	82.3	0.4	73.8	0.9	69.4	0.7
		Discrepancy	72.4	0.6	64.0	0.4	60.0	0.9
		VDC	91.9	0.5	88.3	0.7	82.7	0.5
		Deep kNN	75.7	1.3	66.8	1.4	65.8	1.3
		Confident	69.8	1.1	52.8	1.5	63.3	1.3
		CLIP Sim.	92.1	0.2	90.5	0.2	80.8	0.6
		LEMON <sub>FIX</sub>	90.4	0.5	89.5	0.4	80.2	0.5
		LEMON <sub>OPT</sub>	92.9	0.5	91.5	0.5	82.3	0.7
		CapFilt	98.7	0.2	98.4	0.2	93.6	0.5
	random	LLaVA	82.6	0.3	65.1	0.6	76.7	0.2
		Datamap	94.5	0.1	91.0	0.2	85.5	0.2
		Discrepancy	90.8	0.4	84.2	0.6	80.1	0.9
		VDC	99.1	0.2	98.3	0.2	95.8	0.3
		Deep kNN	94.5	0.4	88.4	0.7	87.9	0.7
		Confident	85.0	0.8	71.4	1.1	81.8	0.9
		CLIP Sim.	99.5	0.1	99.4	0.1	97.1	0.1
LEMON <sub>FIX</sub>	99.5	0.2	99.4	0.1	97.3	0.2		
LEMON <sub>OPT</sub>	99.6	0.1	99.5	0.1	97.5	0.1		
CapFilt	99.9	0.0	99.9	0.0	98.9	0.1		

Table I.6: mmimdb: Label Error Detection

Dataset	Noise Type	Method	AUROC		AUPRC		F1	
			mean	std	mean	std	mean	std
mmimdb	cat	LLaVA	58.4	0.2	46.4	0.2	58.5	0.1
		Datamap	55.4	0.5	44.1	0.6	57.2	0.0
		Discrepancy	57.8	0.4	46.1	0.9	57.4	0.2
		VDC	80.5	0.3	67.1	0.3	69.3	0.6
		deep k-nn	61.2	0.4	47.2	0.5	58.3	0.4
		Confident	52.8	1.1	41.4	0.6	51.8	1.8
		CLIP Sim.	85.1	0.3	77.8	0.7	72.7	0.6
		LEMON <sub>FIX</sub>	84.3	0.3	77.7	0.8	69.9	0.8
		LEMON <sub>OPT</sub>	86.0	0.1	79.4	0.6	73.5	0.3
	CapFilt	82.7	0.7	73.3	1.2	71.3	0.3	
	noun	LLaVA	59.1	0.3	44.2	0.6	55.2	0.2
		Datamap	53.9	0.7	38.9	0.8	53.9	0.2
		Discrepancy	58.8	1.0	43.3	1.3	54.6	0.7
		VDC	79.0	0.2	62.3	0.3	65.9	0.4
		deep k-nn	61.4	0.1	44.2	0.3	55.8	0.3
		Confident	53.7	1.3	38.8	0.8	48.6	1.9
		CLIP Sim.	82.8	0.4	72.8	0.5	68.4	1.1
		LEMON <sub>FIX</sub>	82.1	0.4	72.7	0.6	65.2	1.1
		LEMON <sub>OPT</sub>	84.2	0.5	75.9	0.5	69.5	0.4
	CapFilt	79.9	0.1	66.2	0.4	67.1	0.3	
	random	LLaVA	58.5	0.8	46.7	0.5	58.5	0.1
		Datamap	55.1	0.5	43.7	0.6	57.0	0.3
		Discrepancy	60.0	0.7	47.4	0.8	58.0	0.5
		VDC	82.7	0.2	69.6	0.4	72.1	0.4
		deep k-nn	64.0	0.2	49.0	0.1	60.4	0.1
		Confident	52.0	1.0	41.0	0.5	52.1	3.1
		CLIP Sim.	88.1	0.1	82.0	0.2	75.7	0.4
LEMON <sub>FIX</sub>		87.6	0.1	81.9	0.3	73.8	0.2	
LEMON <sub>OPT</sub>		89.3	0.9	84.2	1.3	77.0	1.2	
CapFilt	84.9	0.4	76.4	0.7	73.1	0.3		

Table I.7: mimiccxr: Label Error Detection

Dataset	Noise Type	Method	AUROC		AUPRC		F1	
			mean	std	mean	std	mean	std
mimiccxr	cat	LLaVA	53.9	0.5	42.7	0.7	57.0	0.1
		Datamap	51.0	0.5	40.5	0.7	57.0	0.1
		Discrepancy	60.0	0.7	50.2	0.5	57.2	0.1
		VDC	50.8	0.4	40.3	0.2	57.0	0.1
		Deep k-NN	62.9	0.4	48.0	0.3	59.2	0.1
		Confident	61.8	0.3	47.0	0.2	58.1	0.6
		CLIP Sim.	64.1	0.4	51.7	0.5	59.2	0.0
		LEMON <sub>FIX</sub>	66.3	0.4	54.6	0.5	55.5	0.3
		LEMON <sub>OPT</sub>	70.4	1.6	60.4	1.6	61.1	0.8
	CapFilt	49.7	0.3	40.0	0.2	57.0	0.0	
	random	LLaVA	50.8	0.4	40.6	0.2	57.1	0.0
		Datamap	50.5	0.4	40.3	0.4	57.1	0.0
		Discrepancy	62.5	0.5	52.2	1.0	57.2	0.5
		VDC	52.4	0.9	41.6	0.7	57.2	0.1
		Deep k-NN	66.6	0.6	53.8	1.1	59.4	0.2
		Confident	65.0	1.0	49.5	0.7	61.6	1.1
		CLIP Sim.	66.8	0.8	54.4	0.9	60.1	0.4
		LEMON <sub>FIX</sub>	69.5	0.7	57.8	1.0	57.7	1.2
		LEMON <sub>OPT</sub>	73.7	1.7	64.1	2.2	63.5	0.8
		CapFilt	50.6	0.4	40.1	0.5	57.1	0.0

Table I.8: We show the AUROC and AUPRC of LEMOn when we search for the optimal hyperparameters using a labeled validation set (LEMOn<sub>OPT</sub>) and when we use fixed hyperparameters (LEMOn<sub>FIX</sub>:  $k = 30$ , cosine distance,  $\beta = \gamma = 5$ ,  $\tau_{1,n} = \tau_{1,m} = 0.1$ , and  $\tau_{2,n} = \tau_{2,m} = 5$ ). The mean gap in AUROC is -1.6 (1.3), and the mean gap in AUPRC is -1.6 (2.4). Note that F1 is not computed as it requires external information to select a threshold.

Dataset	Noise Type	AUROC			AUPRC		
		LEMOn <sub>OPT</sub>	LEMOn <sub>FIX</sub>	Gap	LEMOn <sub>OPT</sub>	LEMOn <sub>FIX</sub>	Gap
cifar10	asymmetric	98.8 (0.2)	97.5 (0.2)	-1.4 (0.1)	97.8 (0.5)	94.8 (0.6)	-3.0 (0.1)
	real	98.1 (0.0)	97.7 (0.2)	-0.5 (0.2)	97.4 (0.1)	96.8 (0.3)	-0.5 (0.2)
	symmetric	99.6 (0.1)	99.5 (0.1)	-0.2 (0.1)	99.4 (0.1)	99.2 (0.1)	-0.2 (0.1)
cifar100	asymmetric	96.7 (0.3)	94.9 (0.3)	-1.9 (0.1)	95.3 (0.4)	92.1 (0.4)	-3.2 (0.1)
	real	90.8 (0.0)	88.9 (0.7)	-1.8 (0.7)	87.4 (0.3)	84.6 (1.1)	-2.8 (0.9)
	symmetric	99.0 (0.0)	98.4 (0.1)	-0.7 (0.1)	98.7 (0.1)	97.7 (0.2)	-1.0 (0.1)
miniImageNet	human	90.0 (0.4)	89.5 (0.2)	-0.5 (0.2)	79.7 (3.0)	81.5 (0.3)	+1.8 (2.7)
StanfordCars	human	73.1 (0.5)	72.6 (0.7)	-0.5 (0.5)	40.5 (0.4)	44.9 (1.4)	+4.4 (1.0)
flickr30k	noun	94.5 (0.2)	93.6 (0.2)	-0.9 (0.3)	92.8 (0.3)	92.0 (0.2)	-0.8 (0.1)
	random	99.5 (0.2)	99.4 (0.2)	-0.0 (0.1)	99.4 (0.3)	99.3 (0.2)	-0.1 (0.2)
mimiccxr	cat	70.4 (1.6)	66.3 (0.4)	-4.1 (1.5)	60.4 (1.6)	54.6 (0.5)	-5.8 (1.5)
	random	73.7 (1.7)	69.5 (0.7)	-4.1 (1.5)	64.1 (2.2)	57.8 (1.0)	-6.3 (1.7)
mmimdb	cat	86.0 (0.1)	84.3 (0.3)	-1.6 (0.3)	79.4 (0.6)	77.7 (0.8)	-1.7 (0.2)
	noun	84.2 (0.5)	82.1 (0.4)	-2.1 (0.6)	75.9 (0.5)	72.7 (0.6)	-3.2 (0.5)
	random	89.3 (0.9)	87.6 (0.1)	-1.6 (0.8)	84.2 (1.3)	81.9 (0.3)	-2.3 (1.2)
mscoco	cat	95.6 (0.2)	92.0 (0.1)	-3.6 (0.1)	94.6 (0.3)	91.8 (0.3)	-2.8 (0.1)
	noun	92.9 (0.5)	90.4 (0.5)	-2.5 (0.2)	91.5 (0.5)	89.5 (0.4)	-2.0 (0.3)
	random	99.6 (0.1)	99.5 (0.2)	-0.1 (0.0)	99.5 (0.1)	99.4 (0.1)	-0.1 (0.0)

Table I.9: Downstream captioning performance when removing 40% samples with highest mislabel scores, and models are trained without early stopping for 10 epochs. We find that filtering noisy data with LEMOn<sub>OPT</sub> improves captioning.

Dataset	Method	B@4	CIDER	ROUGE
flickr30k	No Filtering	31.7 (0.2)	75.6 (0.6)	51.7 (0.1)
	CLIP Sim.	32.8 (0.2)	80.6 (0.8)	53.1 (0.3)
	LEMOn <sub>OPT</sub>	32.5 (0.4)	79.2 (0.6)	52.7 (0.2)
	Clean	33.2 (0.4)	79.8 (1.0)	53.1 (0.2)
mscoco	No Filtering	36.1 (0.3)	119.2 (0.6)	56.9 (0.2)
	CLIP Sim.	39.6 (0.1)	130.8 (0.2)	59.3 (0.0)
	LEMOn <sub>OPT</sub>	39.6 (0.1)	131.0 (0.1)	59.4 (0.1)
	Clean	39.6 (0.1)	131.0 (0.3)	59.4 (0.1)

Table I.10: Performance of our method after ablating various components. We find that mislabel detection performance almost decreases monotonically as we remove additional components, with the exception of two metrics on mmimdb where one ablation is statistically comparable to the original method.

	mmimdb			mscoco		
	AUROC	AUPRC	F1	AUROC	AUPRC	F1
LEMOn <sub>OPT</sub> (Ours)	86.0 (0.1)	79.4 (0.6)	73.5 (0.3)	<b>95.6</b> (0.2)	<b>94.6</b> (0.3)	<b>87.0</b> (0.2)
- $\tau_1$	85.3 (0.3)	78.2 (1.1)	72.9 (0.5)	94.6 (0.3)	93.8 (0.4)	85.2 (0.5)
- $\tau_2$	85.6 (0.1)	78.6 (0.5)	73.3 (0.3)	94.7 (0.2)	93.8 (0.5)	85.4 (0.8)
- $\tau_1, \tau_2$	85.4 (0.2)	78.1 (0.7)	73.0 (0.7)	94.7 (0.3)	93.8 (0.5)	85.3 (0.9)
- $s_n$	<b>86.1</b> (0.2)	<b>79.6</b> (0.6)	<b>73.7</b> (0.2)	94.6 (0.3)	93.6 (0.5)	84.7 (0.7)
- $s_m$	85.3 (0.2)	77.9 (0.7)	73.1 (0.4)	94.9 (0.2)	94.0 (0.4)	86.5 (0.6)
- $s_n, s_m$ (CLIP Sim.)	85.1 (0.3)	77.8 (0.7)	72.7 (0.6)	93.8 (0.2)	93.0 (0.4)	84.5 (0.4)

Table I.11: AUROC of label error detection for each component of our score. We find that  $d_{mm}$  is the most critical component of the score. Of the two nearest neighbor terms, we find that  $s_n$  (nearest image neighbors) is the more important term for most datasets.

	cifar10	cifar100	miniImageNet	stanfordCars	flickr30k	mscoco	mmimdb	mimiccxr
$d_{mm}$ (CLIP Sim.)	92.2 (0.2)	80.8 (0.1)	89.3 (0.2)	69.8 (0.4)	94.8 (0.5)	93.8 (0.2)	85.1 (0.3)	64.1 (0.4)
$s_m$	80.2 (1.1)	65.4 (2.0)	80.8 (0.3)	66.1 (0.6)	75.9 (3.0)	75.9 (0.2)	60.4 (0.7)	59.7 (0.5)
$s_n$	98.1 (0.0)	88.4 (0.1)	84.3 (0.2)	73.0 (0.6)	69.9 (2.3)	76.3 (0.8)	53.7 (0.2)	57.7 (0.7)
$d_{mm} + s_m$	92.5 (0.5)	81.3 (1.1)	89.6 (0.2)	70.0 (0.6)	<b>95.0</b> (0.5)	94.6 (0.3)	<b>86.0</b> (0.2)	64.5 (0.6)
$s_n + s_m$	98.0 (0.2)	88.8 (0.2)	84.5 (0.4)	73.0 (0.7)	83.2 (0.5)	86.1 (0.6)	67.5 (1.0)	64.7 (1.0)
$d_{mm} + s_n$	<b>98.2</b> (0.1)	90.7 (0.2)	<b>90.0</b> (0.4)	<b>73.1</b> (0.5)	94.9 (0.3)	94.9 (0.2)	85.3 (0.2)	64.4 (2.0)
$d_{mm} + s_n + s_m$ (LEMON)	98.1 (0.0)	<b>90.8</b> (0.0)	<b>90.0</b> (0.4)	<b>73.1</b> (0.5)	94.5 (0.2)	<b>95.6</b> (0.2)	<b>86.0</b> (0.0)	<b>70.4</b> (1.6)

Table I.12: Average ( $\sim$ ) per-sample runtime (milliseconds) of each method for label error detection. Standard deviation across 3 random data seeds are shown in parentheses. Experiments were run using job scheduling with GTX A6000 or A100 GPUs.

	cifar10	cifar100	miniImageNet	stanfordCars	mscoco	flickr30k	mimiccxr	mmimdb
LEMON	10.1 (0.5)	9.6 (0.5)	7.8 (1.6)	11.0 (2.0)	18.8 (1.8)	35.9 (1.2)	52.2 (2.7)	21.1 (1.4)
CLIP Sim.	2.9 (0.0)	2.9 (0.0)	6.5 (0.1)	6.8 (0.1)	11.5 (0.1)	9.8 (0.0)	17.1 (0.0)	24.7 (1.0)
Deep kNN	7.0 (1.4)	6.3 (1.8)	7.5 (1.8)	5.4 (0.4)	50.5 (1.4)	63.3 (3.4)	337.6 (13.9)	29.6 (2.0)
Datamap	19.0 (0.7)	19.3 (0.5)	39.2 (1.4)	41.2 (2.7)	33.1 (0.2)	34.6 (1.3)	47.1 (1.9)	52.5 (0.5)
CaptFilt	-	-	-	-	13.5 (0.1)	14.8 (0.0)	29.6 (0.0)	35.3 (1.6)
VDC	-	-	-	-	4460 (880)	5160 (503)	7932 (3.4)	4672 (357)

### I.7. Varying Validation Set Size

In Figure I.3, we examine the effect of varying validation set size (by random subsampling) on  $LEMON_{OPT}$ .

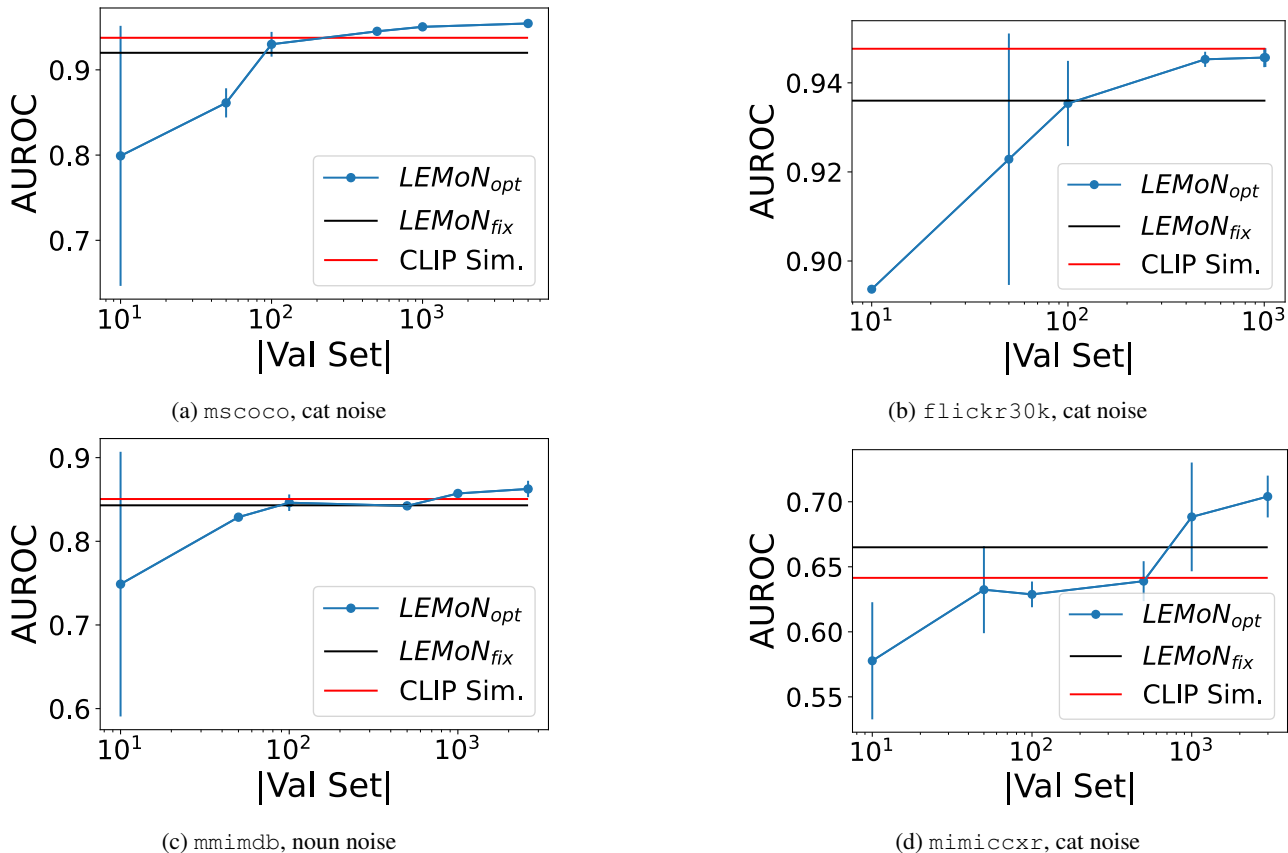


Figure I.3: Test-set AUROC of mislabel detection with varying size of the labeled validation set for  $LEMON_{OPT}$ . Note that  $LEMON_{FIX}$  and CLIP Sim. do not have any hyperparameters and as such do not rely on a labeled validation set.

### I.8. Empirical Comparison with Thomas & Kovashka (2022)

In Table I.13, we compare the performance of  $LEMON_{OPT}$  against the four individual scores and two combined scores proposed in Thomas & Kovashka (2022), using the datasets and noise types shown in Table 1. For the Comb-Val strategy, as there are four terms, we sweep over weights in  $\{1, 2, 3, 4, 5\}^4$ , selecting the best combination using a labeled validation set, identically to LEMON. For the Comb-Stat strategy, we use the mean and standard deviations, as in Equation (8) in Thomas & Kovashka (2022). We find that none of the combined scores significantly outperform  $\Upsilon_X^{DIS}$ . This is because in both combination strategies, a non-zero weight is placed on the other terms, which essentially adds noise to the final score without contributing any signal.

### I.9. Real-World Web Scale Corpus (CC3M)

We conduct an experiment of  $LEMON_{FIX}$  on CC3M (Changpinyo et al., 2021), a large web-scraped dataset of images and annotations, where we demonstrate the utility of LEMON filtered data on CLIP pretraining. We download CC3M, which contains 2.9 million valid URLs to image-caption pairs. We then pretrain a CLIP model (ViT-B/16) from scratch on this dataset for 20 epochs, with a batch size of 128, and using a cyclic learning rate scheduler with a learning rate of  $10^{-4}$ .

We then use this CLIP model as the basis to compute distances for  $LEMON_{FIX}$ , using the reasonable hyperparameters from the main paper:  $k = 30$ , cosine distance,  $\tau_{1,n} = \tau_{1,m} = 0.1$ , and  $\tau_{2,n} = \tau_{2,m} = 5$ . We then select the 1 million samples with the lowest mislabel scores, filtering out the 1.9 million samples most suspected to be mislabels. We pre-train another CLIP model from scratch on this subset using the same architecture and setup as above. We evaluate the resulting model on

Table I.13: Comparison of label error detection performance of LEMOn versus baselines from Thomas & Kovashka (2022).

Dataset	Metric	$\Upsilon_X^{DIS}$	$\Upsilon_Y^{DIS}$	$\Upsilon_X^{DIV}$	$\Upsilon_Y^{DIV}$	Comb-Val	Comb-Stat	LEMOn <sub>OPT</sub>
cifar10	AUROC	77.1 (1.9)	48.2 (1.2)	50.3 (3.5)	45.0 (1.9)	59.6 (19.5)	59.6 (19.5)	<b>98.1 (0.0)</b>
	AUPRC	70.4 (2.7)	41.2 (1.1)	41.6 (1.6)	38.9 (2.1)	51.6 (20.5)	51.6 (20.5)	<b>97.4 (0.1)</b>
	F1	66.4 (2.2)	58.4 (0.8)	58.4 (0.8)	58.4 (0.8)	62.9 (8.0)	62.9 (8.0)	<b>92.0 (0.2)</b>
cifar100	AUROC	66.0 (1.5)	49.7 (0.9)	51.1 (1.1)	49.9 (1.9)	51.5 (6.2)	51.8 (7.0)	<b>90.8 (0.0)</b>
	AUPRC	57.4 (2.3)	40.0 (1.3)	42.8 (1.6)	40.9 (1.1)	41.9 (5.6)	42.3 (6.3)	<b>87.4 (0.3)</b>
	F1	58.9 (0.8)	57.2 (0.3)	57.3 (0.2)	57.1 (0.2)	57.8 (0.7)	57.8 (0.7)	<b>78.4 (0.0)</b>
miniImageNet	AUROC	79.5 (0.3)	47.4 (0.5)	64.6 (0.2)	48.3 (1.1)	75.4 (0.2)	76.6 (0.3)	<b>90.0 (0.4)</b>
	AUPRC	65.9 (0.4)	32.5 (0.0)	46.3 (0.2)	33.5 (0.2)	59.8 (0.2)	61.7 (0.3)	<b>79.7 (3.1)</b>
	F1	64.0 (0.1)	50.9 (0.1)	53.7 (0.3)	50.9 (0.1)	60.5 (0.4)	61.2 (0.6)	<b>76.9 (0.2)</b>
stanfordCars	AUROC	65.7 (0.3)	50.8 (1.1)	51.9 (0.9)	50.1 (0.5)	62.0 (0.1)	64.1 (0.1)	<b>73.1 (0.5)</b>
	AUPRC	33.3 (0.6)	23.3 (0.7)	23.3 (0.4)	23.4 (0.2)	30.1 (0.2)	31.9 (0.4)	<b>40.5 (0.4)</b>
	F1	44.3 (0.7)	38.0 (0.1)	38.0 (0.5)	38.2 (0.2)	41.4 (0.1)	43.3 (0.2)	<b>51.3 (0.5)</b>
flickr30k	AUROC	73.0 (0.6)	53.3 (1.4)	49.9 (2.9)	52.9 (0.2)	63.9 (0.6)	70.5 (0.2)	<b>94.5 (0.2)</b>
	AUPRC	59.2 (1.8)	37.1 (1.8)	32.8 (1.9)	37.0 (0.8)	47.2 (2.1)	53.7 (2.7)	<b>92.8 (0.3)</b>
	F1	59.0 (0.3)	50.9 (0.3)	50.8 (0.4)	50.8 (0.3)	59.4 (3.7)	62.6 (3.3)	<b>83.6 (1.4)</b>
mimiccxr	AUROC	60.0 (0.7)	49.6 (0.4)	49.7 (1.1)	49.1 (1.3)	51.6 (3.1)	52.5 (4.7)	<b>70.4 (1.6)</b>
	AUPRC	50.2 (0.5)	39.3 (0.5)	39.8 (0.1)	39.4 (1.0)	40.5 (2.7)	42.2 (4.8)	<b>60.4 (1.6)</b>
	F1	57.2 (0.1)	57.0 (0.0)	57.0 (0.1)	57.0 (0.1)	57.2 (0.1)	57.3 (0.1)	<b>61.1 (0.8)</b>
mmimdb	AUROC	57.8 (0.4)	50.1 (0.4)	48.6 (0.4)	50.4 (0.3)	53.4 (2.0)	54.7 (2.6)	<b>86.0 (0.1)</b>
	AUPRC	46.1 (0.9)	40.2 (0.6)	38.9 (0.5)	40.2 (0.5)	41.0 (3.1)	42.0 (3.6)	<b>79.4 (0.6)</b>
	F1	57.4 (0.2)	57.1 (0.0)	57.1 (0.0)	57.1 (0.0)	56.5 (1.8)	56.6 (2.0)	<b>73.5 (0.3)</b>
mscoco	AUROC	72.7 (0.3)	48.5 (0.8)	52.9 (0.8)	48.8 (0.2)	49.7 (0.5)	50.1 (0.7)	<b>95.6 (0.2)</b>
	AUPRC	67.2 (0.4)	39.1 (0.5)	42.3 (1.0)	39.4 (0.0)	38.9 (0.3)	39.5 (0.2)	<b>94.6 (0.3)</b>
	F1	62.5 (0.3)	57.0 (0.2)	57.1 (0.0)	57.0 (0.2)	57.3 (0.1)	57.3 (0.1)	<b>87.0 (0.2)</b>

Table I.14: Performance of each method on the Datacomp (Gadre et al., 2024) small benchmark from the filtering track. As of 2024/11/14, only 9.96M images (“Data Available”) out of 12.8M are accessible. We compare the performance of LEMON<sub>OPT</sub> versus the CLIP score baseline after filtering to 3.5M images.

	Method	ImageNet	ImageNet Dist. Shifts	VTAB	Retrieval	Avg (38 datasets)
Data Available (9.96M Samples)	LEMÓN <sub>FIX</sub>	<b>0.045</b>	<b>0.053</b>	<b>0.188</b>	0.116	<b>0.168</b>
	CLIP score	0.043	0.049	0.177	<b>0.119</b>	0.160
From Gadre et al. (2024) (12.8M Samples)	No filtering	0.025	0.033	0.145	0.114	0.132
	Basic filtering	0.038	0.043	0.150	0.118	0.142
	Text-based	0.046	0.052	0.169	<b>0.125</b>	0.157
	Image-based	0.043	0.047	0.178	0.121	0.159
	LAION-2B filtering	0.031	0.040	0.136	0.092	0.133
	CLIP score	<b>0.051</b>	<b>0.055</b>	<b>0.190</b>	0.119	<b>0.173</b>
	Image-based + CLIP score	0.039	0.045	0.162	0.094	0.144

zero-shot classification using the VTAB benchmark (Zhai et al., 2019), and compare it with CLIP models trained using data filtered to 1 million examples using the CLIP similarity baseline, and the original unfiltered model.

In Table I.15, we find that LEMON<sub>FIX</sub> marginally outperforms the CLIP similarity baseline on average zero-shot accuracy, though both underperform pretraining on the full corpus. Similar results can be found for few-shot linear probing (Table I.16) and full finetuning (Table I.17). One likely explanation of this is that although a large proportion of images in the CC3M dataset are technically “mislabelled” in that the caption is not a precisely correct description of the image, some substrings of these noisy captions may, on aggregate, contain useful word associations which the model learns, and thus may be useful to downstream tasks.

We examine images of images selected to be mislabels by our method in Figure I.4. We find that our method identifies images that are completely mislabeled – one cause of which is images changing after they have been indexed. In addition, our method also identifies samples which are ambiguous or imprecise.

### I.10. Real-World Web Scale Corpus (Datacomp)

We conduct an experiment of LEMON<sub>FIX</sub> on Datacomp (Gadre et al., 2024). We use the small dataset from the filtering track, which originally consisted of 12.8M images. As these images are accessed directly from the web, only 9.96M images were able to be downloaded as of 2024/11/14. We apply LEMON<sub>FIX</sub> to this dataset using OpenAI CLIP ViT-L/14 embeddings provided by Datacomp. We select the 3.5M images with lowest mislabel scores, and use the default hyperparameters from Datacomp to train a CLIP model, and evaluate it on the same 38 zero-shot classification datasets. We compare with filtering using only the CLIP score (equivalent to CLIP Sim.) to the same number of images. In Table I.14, we find that given the available images, LEMON<sub>FIX</sub> outperforms the baseline on average, and on three of four individual evaluations. However, neither method outperforms the scores reported in the original paper due to their dataset being larger.

### I.11. Hyperparameters Used for Real-World

We show the hyperparameters used for the real-world experiment in Table I.18. We use  $k = 30$ , cosine distance, and these hyperparameters, which originate from a hyperparameter search on synthetically noised data. We note that flickr30k has some negative hyperparameters, which we attribute to overfitting to a relatively small validation set during hyperparameter selection.

### I.12. Examples of Detected Real Label Errors

We show additional examples of label errors in Figure I.5.

Table I.15: Zero-shot accuracy (%) of various CLIP models on the VTAB benchmark (Zhai et al., 2019). CLIP models (ViT-B/16) are pretrained from scratch on a subset of CC3M (Changpinyo et al., 2021) which has been filtered to 1 million samples using LEMON<sub>FIX</sub> and the CLIP similarity baseline, using a version of CLIP pretrained on the entire dataset.

	CLIP Sim.	LEMON <sub>FIX</sub>	Unfiltered
caltech101	28.25	<b>28.99</b>	51.43
cifar100	<b>11.02</b>	6.79	18.65
clevr_closest_object_distance	18.11	<b>22.58</b>	25.76
clevr_count_all	<b>12.98</b>	12.65	12.05
dmlab	14.78	<b>16.22</b>	16.62
dsprites_label_orientation	<b>2.44</b>	1.34	1.98
dsprites_label_x_position	3.06	<b>3.20</b>	3.13
dsprites_label_y_position	<b>3.11</b>	2.89	3.20
dtd	<b>6.60</b>	3.94	12.34
eurosat	14.37	<b>22.07</b>	9.93
flowers	<b>6.11</b>	5.19	6.83
food101	4.94	<b>5.31</b>	9.02
pets	<b>7.63</b>	4.69	8.23
sun397	13.89	<b>14.22</b>	24.02
svhn	7.80	<b>12.35</b>	8.00
<b>Average</b>	10.34	<b>10.83</b>	14.08

Table I.16: Few-shot linear-probing accuracy (%) using 5 samples per class of various CLIP models on the VTAB benchmark.

	CLIP Sim.	LEMON <sub>FIX</sub>	Unfiltered
caltech101	<b>53.50</b>	53.45	60.82
cifar100	<b>18.66</b>	17.11	23.59
clevr_closest_object_distance	23.39	<b>25.48</b>	25.97
clevr_count_all	19.19	<b>20.43</b>	23.22
dmlab	18.43	<b>21.01</b>	19.72
dsprites_label_orientation	<b>9.11</b>	8.40	9.91
dsprites_label_x_position	3.78	<b>4.52</b>	4.20
dsprites_label_y_position	6.93	<b>7.90</b>	8.20
dtd	28.88	<b>29.84</b>	38.78
eurosat	<b>66.17</b>	65.20	64.76
flowers	<b>57.31</b>	54.53	62.35
food101	11.94	<b>13.50</b>	17.89
pets	<b>18.92</b>	17.12	24.42
sun397	21.43	<b>21.60</b>	32.53
svhn	12.05	<b>12.60</b>	12.99
<b>Average</b>	24.65	<b>24.85</b>	28.62

Table I.17: Full finetuning linear-probing accuracy (%) using 5 samples per class of various CLIP models on the VTAB benchmark.

	CLIP Sim.	LEMON <sub>FIX</sub>	Unfiltered
caltech101	65.91	<b>65.94</b>	74.80
cifar100	<b>49.74</b>	48.53	56.54
clevr_closest_object_distance	52.35	<b>55.83</b>	53.74
clevr_count_all	<b>55.03</b>	54.05	57.99
dmlab	39.10	<b>39.31</b>	43.00
dsprites_label_orientation	38.58	<b>38.65</b>	42.75
dsprites_label_x_position	27.08	<b>36.61</b>	32.78
dsprites_label_y_position	46.31	<b>47.63</b>	49.56
dtd	<b>44.20</b>	42.07	51.91
eurosat	<b>92.70</b>	92.09	93.00
flowers	<b>65.72</b>	63.65	72.45
food101	47.75	<b>48.62</b>	55.23
pets	<b>40.86</b>	40.69	51.57
sun397	<b>49.58</b>	49.45	57.84
svhn	38.08	<b>39.25</b>	42.74
<b>Average</b>	50.20	<b>50.83</b>	55.73



fresh milk in the glass on colour background, illustration



a very young baby girl playing with toys in a white studio



portrait of a stock photo



homes for sale and luxury real estate including horse farms and property in the areas



tangled tree roots on a forest trail

[Visit homeyhomey.club](http://homeyhomey.club)

a park covered in yellow leaves and lined with tall trees turning bright yellow during an autumn day



face of people -- stock photo #



begin your exercise with a jump rope easy and funny



evil looking person sitting atop a hay bale royalty - free

Figure I.4: Sample images and captions from CC3M which have been identified as mislabeled by LEMON<sub>FIX</sub>.

Table I.18: Hyperparameters used for the real-world experiment. We use  $k = 30$ , cosine distance, and the hyperparameters below, which originate from a hyperparameter search on synthetically noised data.

	$\beta$	$\gamma$	$\tau_{1,n}$	$\tau_{2,n}$	$\tau_{1,m}$	$\tau_{2,m}$
cifar10	20	10	0	5	0	5
cifar100	15	0	0	5	0	0
mscoco	5.324	11.057	5.143	10.498	7.233	15.637
mmimdb	15	5	5	10	5	10
flickr30k	0.092	-0.177	-0.274	-0.074	-0.072	0.000
mimiccxr	5	10	5	10	5	10

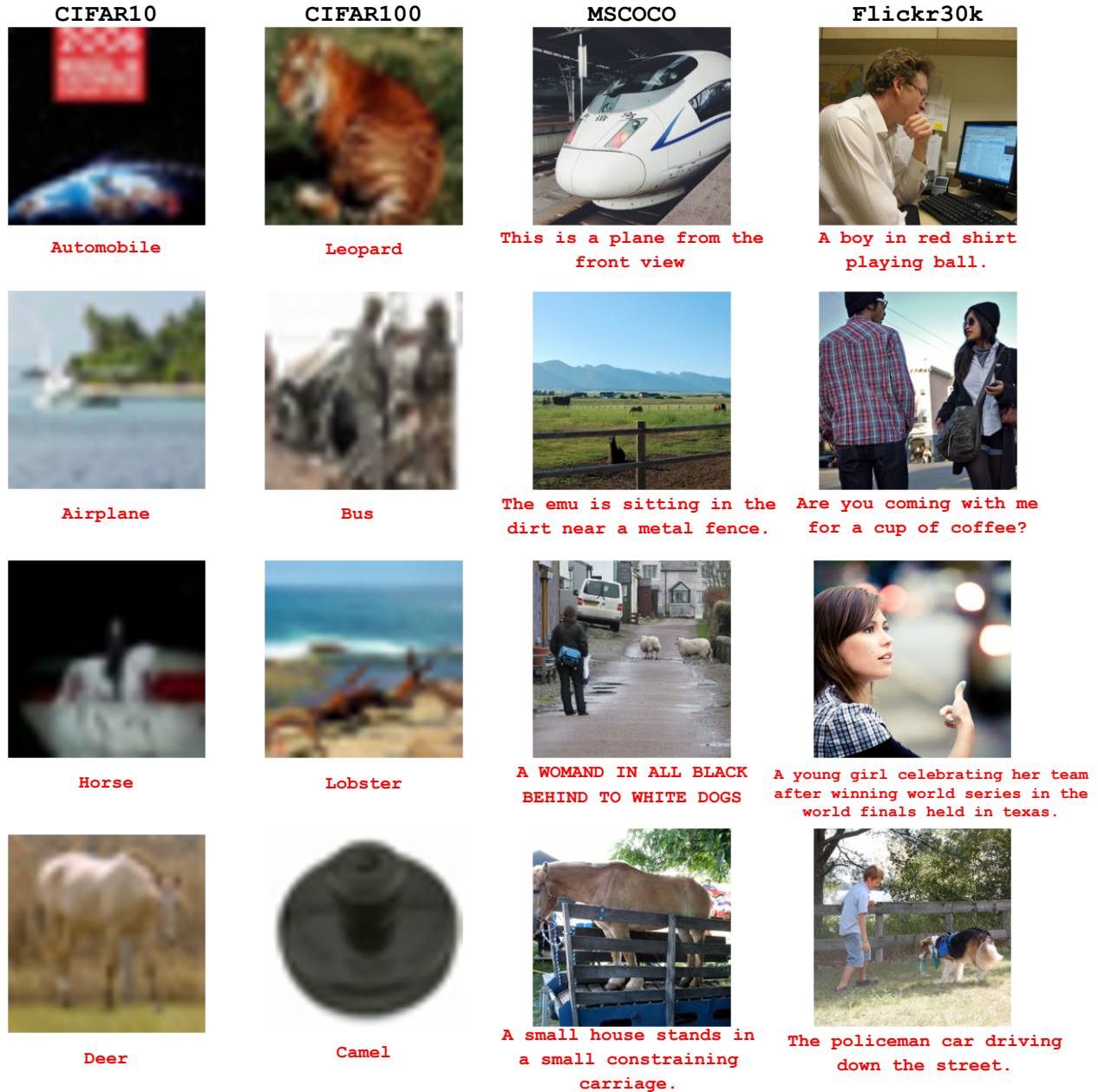


Figure I.5: Example images in each dataset identified by our method to be mislabels, and labeled as errors by a human annotator.

### I.13. Comparison with Northcutt et al., 2021 (Northcutt et al., 2021b)

In Northcutt et al., 2021 (Northcutt et al., 2021b), the authors utilized confident learning (Northcutt et al., 2021a) to identify suspected errors in the test sets of `cifar10` and `cifar100`. They then obtained 5 human labels for each suspected error using Amazon Mechanical Turk, and confirmed the image to be a mislabel if at least 3 of 5 workers stated so. This amounts to 54 confirmed mislabels in `cifar10` (out of 221 suspected), and 585 confirmed mislabels in `cifar100` (out of 1650 suspected). In this section, we compare the performance of  $\text{LEMÓN}_{\text{FIX}}$  versus the CLIP similarity baseline on this set. As this set is a subset of the images identified to be mislabels by confident learning, we are not able to compare our model performance with confident learning itself. In addition, this presents a pessimistic view (lower bound) of the performance of our method, as there are many images identified by  $\text{LEMÓN}$  which *are* mislabeled, but were not selected by confident learning in (Northcutt et al., 2021b). We demonstrate examples of these images in Figure I.6.

In Table I.19, we compare the performance of  $\text{LEMÓN}_{\text{FIX}}$  with the CLIP similarity baseline on the error set from Northcutt et al., 2021 (Northcutt et al., 2021b). First, we compute the mean ranking of all error set samples as ranked by each method, out of 10,000 test-set samples. We find that our method ranks error set samples higher on average than the baseline, though the variance is large. Next, we subset to the top —CL Set— ranked samples for each method, and compute the percentage of which are actually in the error set. We note that this precision metric is upper bounded by the precision of the reference method (confident learning). Again, we find that  $\text{LEMÓN}_{\text{FIX}}$  outperforms the baseline, and is able to identify more actual label errors than CLIP similarity at this threshold.



Figure I.6: Demonstrative examples of mislabeled samples in `cifar10` and `cifar100` which have been identified by our method in the top —CL Set—, but was not identified by confident learning in Northcutt et al., 2021 (Northcutt et al., 2021b) and thus was not a part of their error set.

### I.14. Downstream Classification with Label Error Detection-based Filtering

Here, we show the impact of filtering out different proportions of the training data based on label error predictions, and obtaining test performance.

#### I.14.1. AVERAGE ACCURACY

We observe that the gap in performance is low between  $\text{LEMÓN}_{\text{OPT}}$  and the best method (less than 1% on `stanfordCars`) in terms of downstream accuracy (see Figure I.7).  $\text{LEMÓN}_{\text{OPT}}$  is the best method on `miniImageNet`.

Table I.19: Comparison of LEMOn<sub>FIX</sub> (Ours) with the CLIP similarity baseline on the human labeled error set from Northcutt et al., 2021 (Northcutt et al., 2021b). In this prior work, the authors used confident learning to identify |CL Set| candidate label errors in `cifar10` and `cifar100`, |Error Set| of which are confirmed to be mislabels by Mechanical Turkers. Mean Ranking denotes the average ranking of all error set samples as ranked by each method. Precision @ Top |CL Set| involves taking the top |CL Set| samples as ranked by each method, and computing the percentage of which are in the error set. Note that each dataset’s test set consists of 10,000 samples. Numbers in parentheses represent one standard deviation.

Dataset	CL Set	Error Set	Mean Ranking		Precision @ Top  CL Set		
			LEMOn <sub>FIX</sub>	CLIP Sim.	Oracle	LEMOn <sub>FIX</sub>	CLIP Sim.
<code>cifar10</code>	275	54	1269.7 (1905.1)	2681.0 (2507.1)	19.64%	6.55%	1.45%
<code>cifar100</code>	2235	585	2357.5 (1981.5)	3642.1 (2719.5)	26.17%	14.41%	10.16%

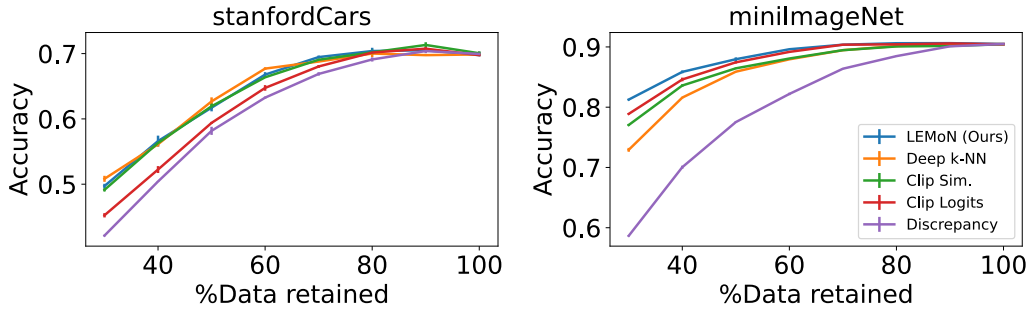


Figure I.7: Downstream accuracy on `stanfordCars` and `miniImageNet`.

**I.15. Area Under Test Error vs % Data Retained Curve**

We compute the area under the test error (i.e., 1-accuracy) vs % data retained curve in Table I.20. Note that the minimum data retained is 30% (i.e., the minimum amount of data required for training the downstream model).

On `cifar10` and `miniImageNet`, we observe that LEMOn performs the best in terms of AUC (i.e., lowest test error). On `stanfordCars`, Deep k-NN performs better. However, the gap in performance is low between LEMOn<sub>OPT</sub> and the best method on all datasets.

Table I.20: Area under the curve: test error vs % data retained for all four classification datasets. Lower is better, and bold denotes best method.

Method	<code>cifar10</code>	<code>cifar100</code>	<code>stanfordCars</code>	<code>miniImageNet</code>
CLIP Sim.	4.58	<b>16.89</b>	24.52	8.84
CLIP Logits	4.54	16.90	25.72	8.27
Discrepancy	5.98	20.22	26.57	13.07
Deep k-NN	4.37	17.54	<b>24.46</b>	9.29
Ours	<b>4.23</b>	17.02	24.48	<b>7.91</b>

**I.16. Out-of-Domain Robustness**

We report the test performance on an Out-of-Domain (OOD) dataset CIFAR-10C (Hendrycks & Dietterich, 2018), when models are trained and validated on the `cifar10` noisy train set. The CIFAR-10C dataset contains 19 corruptions applied to the `cifar10` test set, with varying severity of corruption. Then, robustness is measured as the average test top-1 class accuracy performance on the CIFAR-10C dataset (across all corruption types and severities), following prior work (Diffenderfer et al., 2021). We observe that our method outperforms the CLIP similarity baseline on robustness, when percentage of data retained is less than 60%.

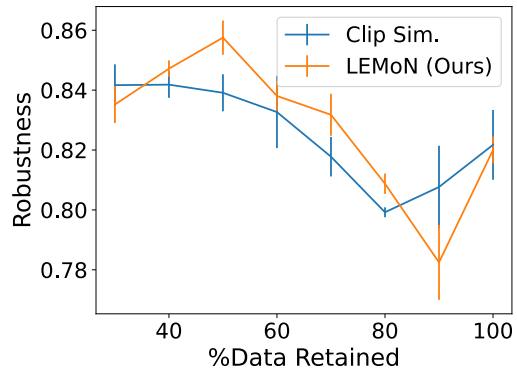


Figure I.8: Downstream accuracy on CIFAR-10C, averaged across all corruption types.