

---

# On Revisiting Entropy for Identifying Mislabeled Medical Images

---

Anonymous Authors<sup>1</sup>

## Abstract

Mislabeled samples in training datasets severely degrade the performance of deep networks, as overparameterized models tend to memorize erroneous labels. We address this challenge by proposing a novel approach for mislabeled data detection that leverages training dynamics. Our method is grounded in the key observation that correctly labeled samples exhibit consistent entropy decrease during training, while mislabeled samples maintain relatively high entropy throughout the training process. Building on this insight, we introduce a signed entropy integral (SEI) statistic that captures both the magnitude and temporal trend of prediction entropy across training epochs. SEI is broadly applicable to classification networks and demonstrates particular effectiveness when integrated with contrastive language-image pre-training (CLIP) architectures. Through extensive experiments on three medical imaging datasets—a domain particularly susceptible to labeling errors due to diagnostic complexity—spanning diverse modalities and pathologies, we demonstrate that SEI achieves state-of-the-art performance in mislabeled data identification, outperforming existing methods while maintaining computational efficiency and implementation simplicity.

## 1. Introduction

Deep networks have achieved remarkable results in medical imaging, enabling applications from tumor segmentation to disease classification (Litjens et al., 2017). Yet, their performance hinges on the quality of training data (Shi et al., 2024). In practice, medical datasets often contain mislabeled samples due to the complexity of diagnosis, inter-observer variability, and limited annotation resources (Shi et al., 2024; Alderman et al., 2025). For instance, in dermatology, the vi-

sual appearance of skin lesions can overlap heavily between malignant and benign conditions; melanomas may resemble benign nevi in early stages, and fungal infections can mimic inflammatory skin disorders. Even experienced dermatologists may disagree without histopathological confirmation. Similarly, in ophthalmology, subtle retinal changes in early diabetic retinopathy can be challenging to grade consistently, especially when annotation guidelines differ across graders. Such challenges mean that noisy labels are a realistic concern in many medical imaging datasets.

Noisy labels pose a particular risk for overparameterized deep networks, which can fit even randomly assigned labels given sufficient capacity (Zhang et al., 2017). When this happens, a model effectively memorizes mislabeled samples by learning overly specific, non-generalizable features. For example, if an image of a benign skin nevus is mistakenly labeled as melanoma, the model may latch onto spurious visual patterns unique to that single image—such as lighting artifacts or sensor noise—rather than features indicative of melanoma. These spurious correlations will lead to overfitting and degraded performance.

Our aim is to automatically identify and remove mislabeled samples from training data. This not only uncovers systematic annotation errors but also improves overall label quality. Moreover, because large medical datasets are often too extensive for exhaustive manual inspection, automated methods are essential for isolating mislabeled samples without overburdening domain experts.

Prior works on mislabeled data detection typically involve multi-stage pipelines and tailored loss functions or modules (Chen et al., 2019; 2024; Huang et al., 2019; Li et al., 2020; Cheng et al., 2021; Wei et al., 2024), which can be tightly coupled to specific architectures or disrupt standard training workflows. We take a different approach: a simple yet effective plug-and-play method that works with existing training setups and requires minimal changes. Our method tracks training dynamics of a classification model with a contrastive loss (Radford et al., 2021) to align medical images and labels. We observe that the evolution of prediction entropy encodes subtle cues for distinguishing noisy from clean samples. Correctly labeled samples tend to exhibit a steady entropy decrease as learning progresses, whereas mislabeled samples often maintain relatively high entropy.

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

055 However, entropy alone is insufficient to separate mislabeled  
 056 data from hard but valuable clean samples, which may also  
 057 retain high entropy despite having correct labels. To address  
 058 this, we propose a signed entropy integral (SEI) statistic that  
 059 captures both the magnitude and trend of entropy evolution.  
 060 This signed formulation provides a richer characterization  
 061 of training dynamics: hard clean samples and mislabeled  
 062 samples exhibit distinct patterns in how their entropy changes  
 063 over time, allowing SEI to differentiate the two.

064 Our work makes three key contributions:

- 065 • Through large-scale analysis of training dynamics, we  
 066 identify two informative signals for distinguishing mis-  
 067 labeled from correctly labeled samples: entropy evolu-  
 068 tion and label-prediction consistency over time.
- 069 • Building on these insights, we introduce *signed entropy*,  
 070 a novel extension of Shannon entropy that incorporates  
 071 label consistency, and propose the SEI statistic that  
 072 captures the cumulative training behavior of different  
 073 sample types.
- 074 • we demonstrate that our simple yet effective method  
 075 achieves state-of-the-art performance on three medi-  
 076 cal imaging datasets spanning different modalities and  
 077 pathologies, without requiring architectural modifica-  
 078 tions or complex training procedures.

## 082 2. Methodology

### 083 2.1. Problem Formulation

084 We consider the task of  $K$ -class image classification, where  
 085 the objective is to train a model that predicts a label  $y$  for  
 086 an input image  $x$ . The training dataset  $\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}$   
 087 contains two types of samples. A *mislabeled* sample is one  
 088 whose assigned label does not match its underlying semantic  
 089 content (e.g.,  $x$  is an image of a melanoma but the label  $y$   
 090 is “nevus”). A *correctly labeled* sample is one where the  
 091 assigned label aligns with the true category. Our goal is to  
 092 identify mislabeled data in  $\mathcal{D}_{\text{train}}$  by exploiting differences  
 093 in their training dynamics compared to correctly labeled  
 094 data.

### 095 2.2. Preliminary: CLIP for Image Classification

096 CLIP (Radford et al., 2021) demonstrates strong zero-shot  
 097 performance by learning joint visual-textual representations  
 098 from large-scale image-text pairs using a contrastive objec-  
 099 tive. Classification is performed by measuring the similarity  
 100 between image features and text embeddings derived from  
 101 prompts such as “a photo of [CLS]”, where [CLS] denotes  
 102 a class name. The prediction probability is computed as

$$103 p(y = k|\mathbf{x}) = \frac{\exp(\text{sim}(\mathbf{v}, \mathbf{t}_k)/\tau)}{\sum_{j=1}^K \exp(\text{sim}(\mathbf{v}, \mathbf{t}_j)/\tau)}, \quad (1)$$

where  $\text{sim}(\mathbf{v}, \mathbf{t}_k)$  denotes the cosine similarity between im-  
 age feature  $\mathbf{v}$  and text embedding  $\mathbf{t}_k$ , and  $\tau$  is a temperature  
 parameter. In this framework, each image is compared  
 against all  $K$  class-specific text embeddings, allowing us  
 to incorporate dataset labels directly while preserving the  
 representational benefits of visual-textual alignment.

### 104 2.3. Findings

To understand training dynamics of mislabeled samples,  
 we first examine their prediction *entropy*. Figure 1 shows  
 entropy trajectories for nevus images in the ISIC dataset  
 and grade 4 diabetic retinopathy images in the DeepDRiD  
 dataset (Liu et al., 2022), comparing correctly labeled sam-  
 ples with mislabeled ones. Early in training, the two groups  
 are entangled, exhibiting similar levels of uncertainty. As  
 training progresses, however, a clear separation emerges:  
 entropy for correctly labeled data decreases steadily, while  
 mislabeled samples maintain elevated entropy. This suggests  
 that entropy provides a promising signal for distinguishing  
 mislabeled samples from clean ones. We observe similar  
 trends across other classes and datasets (see Appendix C).

Nevertheless, entropy alone proves insufficient for reliable  
 detection. Figure 2 illustrates this limitation by comparing a  
 mislabeled sample with a hard-to-learn but correctly labeled  
 sample. Despite their fundamentally different ground-truth  
 status, their entropy curves are nearly indistinguishable,  
 making it difficult to separate the two based solely on this  
 statistic.

To address this limitation, we investigate an additional as-  
 pect of training dynamics: the *alignment between given  
 labels and model predictions over time*. Figure 3 presents  
 alignment statistics for three distinct sample categories: easy  
 clean, hard clean, and mislabeled data, where circle size  
 and opacity reflect frequency. We observe that: (1) most  
 easy clean samples exhibit consistent alignment, with pre-  
 dictions typically matching their labels throughout training;  
 (2) hard clean samples generally demonstrate mixed dynam-  
 ics; (3) mislabeled samples are predominantly characterized  
 by persistent misalignment with their given labels. These  
 distinctive alignment patterns provide complementary in-  
 formation beyond entropy, enabling better discrimination  
 between mislabeled samples and challenging yet correctly  
 labeled ones.

### 105 2.4. Identifying Mislabeled Data

#### 106 2.4.1. SIGNED ENTROPY

The observations in Section 2.3 suggest that two com-  
 plementary cues—entropy dynamics and label-prediction  
 consistency—can be jointly exploited for identifying misla-  
 beled data. We therefore introduce *signed entropy*, a quantity  
 that extends Shannon entropy with a label consistency term.

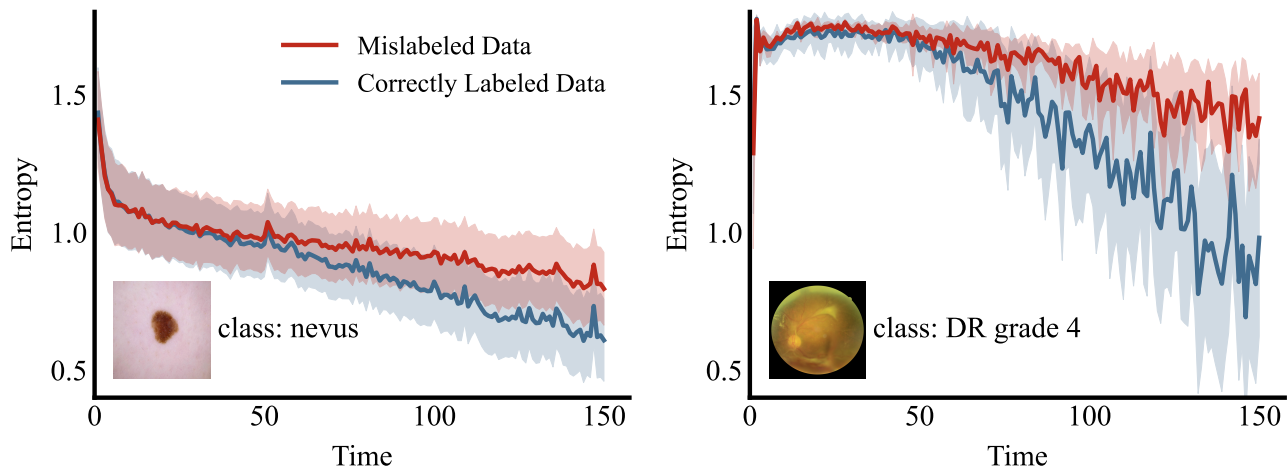


Figure 1. Training dynamics of prediction entropy for correctly labeled versus mislabeled samples. Left: Nevus images from the ISIC dataset, comparing correctly labeled samples (ground truth: nevus; given label: nevus) with mislabeled samples (ground truth: nevus; given label: other skin lesion categories). Right: Grade 4 diabetic retinopathy (DR) images from the DeepDRiD dataset. In both cases, correctly labeled samples exhibit steadily decreasing entropy throughout training, while mislabeled samples maintain persistently high entropy, demonstrating the potential of entropy as an indicator for noisy label detection.

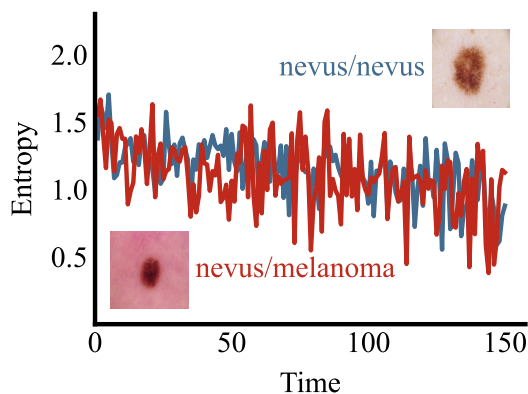


Figure 2. Entropy trajectories for a mislabeled sample (ground truth: nevus; given label: melanoma) and a hard clean sample (ground truth: nevus; given label: nevus). Despite differing label correctness, their entropy curves are nearly indistinguishable. This illustrates that entropy alone cannot reliably distinguish mislabeled data from challenging but clean examples.

Formally, for any  $(\mathbf{x}, y) \in \mathcal{D}_{\text{train}}$ , let  $\mathbf{p}(\mathbf{x}) = (p_1(\mathbf{x}), \dots, p_K(\mathbf{x}))$  denote the model’s posterior distribution over  $K$  classes (cf. Eq. 1). We define signed entropy as

$$\mathcal{H}(\mathbf{p}(\mathbf{x}), y) = (-1)^{\mathbb{1}[y = \arg \max_k p_k(\mathbf{x})]} \sum_{k=1}^K p_k(\mathbf{x}) \log p_k(\mathbf{x}), \quad (2)$$

where the exponent introduces a sign depending on whether the assigned label  $y$  matches the model’s prediction. In other words,  $\mathcal{H}$  reduces to Shannon entropy when  $y$  agrees with the prediction, but flips its sign under misalignment.

**Discussion** Shannon’s entropy is always nonnegative and only reflects distributional uncertainty, making it blind to label correctness: both hard clean samples and mislabeled ones may exhibit high entropy. In contrast, the signed entropy in Eq. 2 attaches a directionality:

- **Positive signed entropy** indicates both uncertainty and label consistency, as typically seen in clean samples (easy or hard).
- **Negative signed entropy** emerges when the model contradicts the given label, flagging potential annotation errors.

#### 2.4.2. SIGNED ENTROPY INTEGRAL

While signed entropy at a single epoch can provide useful information, training dynamics often fluctuate, making individual snapshots unreliable for detecting mislabeled samples. Moreover, it is unclear a priori which epochs contain the most discriminative signals. Aggregating *cumulative behavior* across training offers a more stable criterion, as it naturally smooths out such fluctuations (Huang et al., 2019; Pleiss et al., 2020; Jiang et al., 2022). To this end, we introduce SEI, which accumulates signed entropy values over the entire training trajectory.

Let  $\mathbf{p}^{(t)}(\mathbf{x})$  denote the posterior distribution at epoch  $t$ . SEI is defined as

$$\text{SEI}(\mathbf{x}, y) = \sum_{t=1}^T \mathcal{H}(\mathbf{p}^{(t)}(\mathbf{x}), y), \quad (3)$$

where  $T$  is the total number of training epochs.

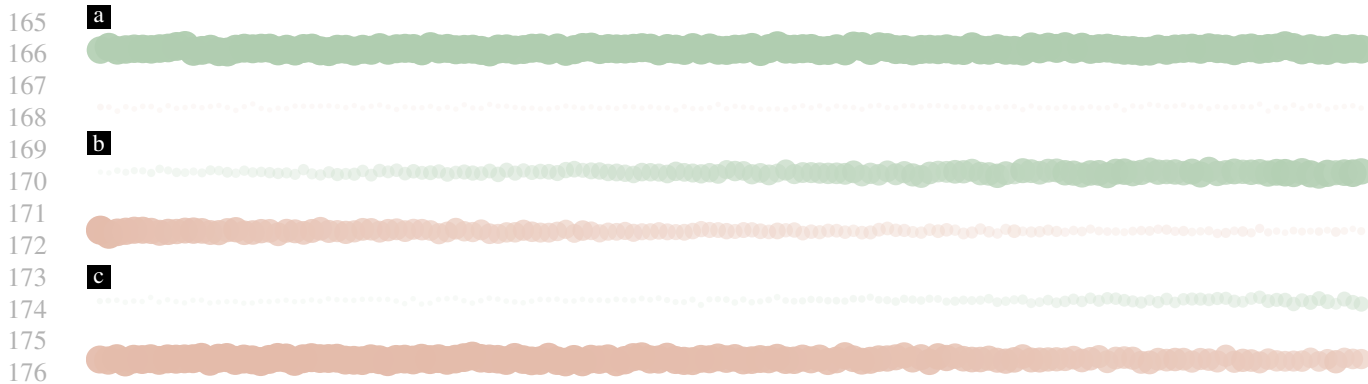


Figure 3. Training dynamics analysis through label-prediction alignment patterns over time. Bubble timeline charts illustrate the evolution of three sample categories during training: (a) easy clean samples, (b) hard clean samples, and (c) mislabeled samples. At each time step, green circles indicate alignment between predicted and given labels, while red circles denote misalignment. Circle size and opacity encode frequency.

**Analysis** Figure 4 illustrates how SEI separates different sample types. For an easy clean sample, the signed entropy remains mostly positive, as predictions stay consistent with its assigned label, resulting in a large positive integral. A hard clean sample may initially exhibit misalignment and accumulate negative contributions, but as training progresses, the signed entropy turns positive; the resulting positive and negative areas partly cancel, leading to a smaller integral. By contrast, a mislabeled sample typically shows persistent disagreement, so its signed entropy curve stays negative throughout training, yielding a strongly negative integral.

In this way, SEI produces a single scalar that naturally ranks samples: mislabeled ones cluster at the negative end, while correctly labeled ones occupy positive values. This simple measure proves effective across datasets and modalities for isolating annotation errors from both easy and hard clean samples (see Appendix D).

### 2.5. Thresholding

We require a threshold to distinguish between clean and mislabeled data. A fixed threshold is impractical since the distribution of SEI values varies across datasets and training configurations. Inspired by Yang et al. (2023); Pleiss et al. (2020); Kim et al. (2025), we use an adaptive thresholding strategy that leverages artificially mislabeled samples as reference points.

Concretely, given  $N$  training samples across  $K$  classes, we randomly select  $N/(K+1)$  instances and reassign their labels to an auxiliary class  $K+1$  that does not exist in the dataset. In the CLIP setting, this corresponds to using a synthetic prompt such as “a dermoscopic image showing other lesions”. Since the auxiliary class is designed to be semantically meaningless, these relabeled samples serve as natural surrogates for mislabeled data. Moreover, choosing  $N/(K+1)$  instances ensures that the auxiliary class

appears with a frequency comparable to the original classes, preventing imbalance in calibration.

We then compute SEIs for all auxiliary-class samples and use their average value as the decision threshold. Any sample whose SEI falls below this threshold is flagged as mislabeled. This simple scheme yields an adaptive, data-driven cutoff without requiring manual tuning.

## 3. Related Work

Addressing noisy labels in datasets has spawned two primary research directions: (1) Explicit mislabeled data detection aims to identify and remove incorrectly labeled instances to improve data quality, while (2) Noise robust learning develops algorithms that maintain performance despite label noise.

**Mislabeled Data Detection** Most methods exploit training signals as proxies for label correctness. Loss-based approaches leverage the intuition that higher losses often indicate incorrect labels. O2U-Net (Huang et al., 2019) alternates between overfitting and underfitting phases, identifying mislabeled samples through consistently higher normalized losses. CORES (Cheng et al., 2021) progressively filters incorrectly labeled instances using training loss in a proposed sample sieve framework. Beyond loss, various proxy measures have been developed, including gradient-based metrics (Zhang & Sabuncu, 2018) and prediction-based statistics (Northcutt et al., 2021; Pleiss et al., 2020; Song et al., 2019). For instance, Northcutt et al. (2021) filter low-confidence samples using class-conditional thresholds on predicted probabilities.

Recent work has explored training-free approaches. SIM-IFEAT (Zhu et al., 2022) detects noisy labels through  $k$ -nearest neighbors in the feature space of a pre-trained model, while DEFT (Wei et al., 2024) leverages CLIP’s image-text

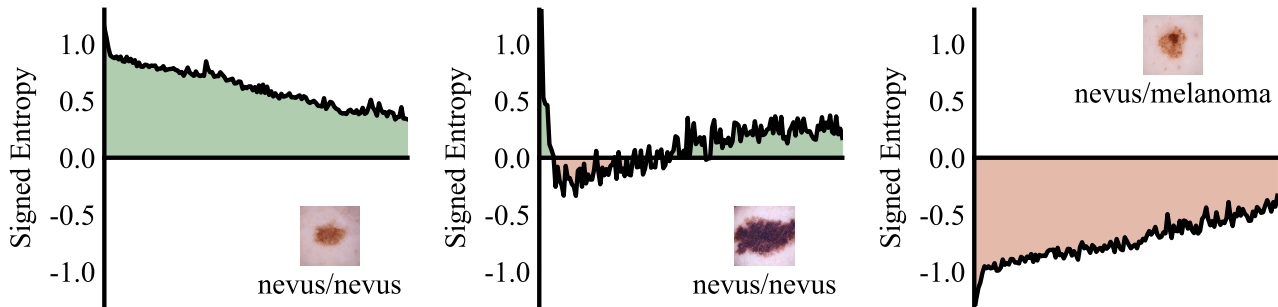


Figure 4. Illustration of SEI. The plots depict signed entropy curves over training epochs for easy clean (left), hard clean (middle), and mislabeled (right) samples. Each curve is averaged over 200 samples, and the signed area under each curve corresponds to the SEI. Correctly labeled samples yield larger SEIs than mislabeled ones.

alignment to learn class-specific prompts for mislabel detection. LEMoN (Zhang et al., 2025) exploits multimodal neighborhoods of image-caption pairs in the latent space of CLIP to automatically identify label errors. However, the effectiveness of these methods depends heavily on the generalization capacity of pre-trained models, which may be limited in specialized domains like medical imaging. ReCoV (Chen et al., 2024) identifies mislabeled medical data through cross-validation, though this increases computational overhead for large datasets. Some approaches (Wang et al., 2024) assume access to clean subsets, while we consider the more restrictive but realistic setting where no training data can be trusted.

**Noise Robust Learning** Rather than targeting specific noisy instances, robust learning methods design modules enabling effective training on noisy datasets. This includes novel architectures (Cheng et al., 2020; Xiao et al., 2015b), loss functions (Wang et al., 2019; Ye et al., 2023), regularization techniques (Cheng et al., 2022; Liu et al., 2020), and training strategies (Lukasik et al., 2020; Xia et al., 2021; Yuan et al., 2024). Some works integrate noisy label detection into training pipelines through loss re-weighting (Ren et al., 2018; Bae et al., 2024) or re-annotation (Han et al., 2019; Arazo et al., 2019; Engleson & Azizpour, 2024). Our work focuses on identifying reliably labeled data rather than recycling mislabeled samples. For simplicity, we discard instances flagged as mislabeled; nevertheless, our method is compatible with approaches that attempt to reuse them (cf. Section 4.3.4).

## 4. Experiments

We evaluate our approach through two complementary assessments: first, we measure mislabeled data detection performance on synthetically corrupted datasets to directly assess identification capability; second, we evaluate downstream classification performance.

### 4.1. Experimental Setup

#### 4.1.1. DATASETS

We conduct experiments on three medical imaging datasets spanning different modalities and diagnostic tasks.

**ISIC 2018.** We use the ISIC 2018 Challenge<sup>1</sup> Task 3 dataset for skin lesion diagnosis, containing dermoscopic images across seven categories: melanoma, nevus, basal cell carcinoma, actinic keratosis/intraepithelial carcinoma, benign keratosis, dermatofibroma, and vascular lesion. The dataset comprises 10,015 training images and 1,512 test images.

**DeepDRiD.** This fundus photography dataset (Liu et al., 2022) targets diabetic retinopathy severity grading using a 5-point scale (0-4) following the International Clinical Diabetic Retinopathy standard. Each patient contributes dual-view images (macula-centered and optic disc-centered) for both eyes. We use 1,200 training images and 400 test images from the official split.

**PANDA.** This dataset<sup>2</sup> contains 10,616 whole slide images (WSIs) of digitized H&E-stained prostate biopsies from Radboud University Medical Center and Karolinska Institute. We focus on the Radboud subset, which provides pixel-level annotations distinguishing background, stroma, benign epithelium, and cancerous epithelium (subdivided into Gleason patterns 3, 4, and 5). Since Gleason grading depends solely on epithelial architecture, we consider four classes: benign epithelium, Gleason 3, Gleason 4, and Gleason 5. Each WSI is tiled into  $224 \times 224$  patches, with background-dominated patches discarded. For remaining patches, we compute area proportions of epithelial categories from pixel-level masks and apply a dominant-label rule: patches are assigned the highest-grade category present (priority: Gleason 5 > Gleason 4 > Gleason 3 > benign epithelium) if it covers more

<sup>1</sup><https://challenge.isic-archive.com/data/#2018>

<sup>2</sup><https://kaggle.com/competitions/prostate-cancer-grade-assessment>

than 80% of non-background pixels; otherwise, patches are discarded. This yields 4,102 benign epithelium patches and 6,914, 6,413, and 6,956 patches for Gleason grades 3, 4, and 5, respectively. We use a 4:1 train/test split.

**Besides**, to further evaluate SEI on real-world noisy data, we conduct experiments on CIFAR-100N (Wei et al., 2022), a natural image dataset with human annotation noise. Although not a medical dataset, CIFAR-100N contains realistic instance-level label errors that serve our evaluation purpose. While other real-world noisy datasets exist, such as Clothing1M (Xiao et al., 2015a) and WebVision (Li et al., 2017), they do not provide clean ground-truth labels for all samples and are primarily designed for learning with noisy labels rather than noise detection. We therefore adopt CIFAR-100N with its official training and test splits.

#### 4.1.2. SYNTHETIC NOISE GENERATION

To simulate real-world annotation errors, we synthesize mislabeled samples at five noise rates  $\eta \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$  using two corruption strategies:

**Symmetric Noise.** In a  $K$ -class setting, each sample with ground-truth label  $a \in \{1, \dots, K\}$  retains its correct label with probability  $1 - \eta$ . With probability  $\eta$ , it is corrupted to a different class  $b \neq a$ , chosen uniformly from the  $K - 1$  alternatives:

$$p_{a \rightarrow b} = \frac{\eta}{K - 1}, \quad b \neq a. \quad (4)$$

**Confusion-Calibrated Noise.** To simulate realistic errors, we first train a reference ResNet-50 classifier and compute an empirical confusion matrix  $T \in \mathbb{R}^{K \times K}$  from its predictions. For a sample with ground-truth label  $a$ , the label remains unchanged with probability  $1 - \eta$ . With probability  $\eta$ , it is corrupted to class  $b$  ( $b \neq a$ ) according to:

$$p_{a \rightarrow b} = \eta \frac{\exp(T_{ab})}{\sum_{k \neq a} \exp(T_{ak})}, \quad b \neq a. \quad (5)$$

This preserves the target corruption rate  $\eta$  while aligning errors with observed class confusions.

#### 4.1.3. EVALUATION METRICS

For mislabeled data detection, we use F1 score as our primary metric, following prior works (Zhu et al., 2022; Kim et al., 2024). F1 score balances precision and recall, providing a more reliable assessment. For downstream image classification, we report accuracy, F1 score, and AUC for comprehensive evaluation.

#### 4.1.4. IMPLEMENTATION DETAILS

All experiments use PyTorch and run on a single NVIDIA RTX 4090 GPU. For mislabeled data identification, we

employ CLIP with a Transformer text encoder and ResNet-50 vision encoder. Training uses SGD with momentum 0.9, weight decay  $1 \times 10^{-4}$ , batch size 128, and initial learning rate  $1 \times 10^{-3}$  for 150 epochs. The learning rate decays by 0.1 at epochs 75 and 115. Images are resized to  $224 \times 224$  and normalized using ImageNet statistics (Russakovsky et al., 2015). For downstream classification tasks, we use the official dataset splits.

## 4.2. Comparison with State-of-the-Art Methods

We evaluate our approach against existing mislabeled data identification methods, including INCV (Chen et al., 2019), BMM (Arazo et al., 2019), GMM (Li et al., 2020), AUM (Pleiss et al., 2020), CORES (Cheng et al., 2021), CL (Northcutt et al., 2021), SIMIFEAT (Zhu et al., 2022), DeFT (Wei et al., 2024), ReCoV (Chen et al., 2024), and LEMoN (Zhang et al., 2025). As summarized in Tables 1 and 2, our approach consistently outperforms all competing methods across all datasets and noise levels.

Under the more challenging confusion-calibrated noise setting, our method delivers significant improvements in mislabeled data detection. On the ISIC dataset, it outperforms the second-best approach by 4.79%, 2.11%, 1.29%, 10.32%, and 8.49% across different noise rates. On DeepDRiD, the gains are 3.96%, 4.16%, 6.37%, 4.49%, and 6.22%, while on PANDA, our method achieves improvements of 11.87%, 9.92%, 9.55%, 5.88%, and 5.33%. Under the symmetric noise setting, we again observe consistent and notable performance boosts, further confirming the effectiveness of our approach.

## 4.3. Analysis and Discussion

### 4.3.1. EFFECTIVENESS OF THE SIGNED TERM

To assess the contribution of the signed component in Eq. 2, we compare SEI against an unsigned counterpart using standard Shannon entropy. As reported in Table 3, SEI consistently outperforms the unsigned variant, validating the effectiveness of signed entropy.

We further visualize score distributions for clean and mislabeled samples under both formulations in Appendix F.1.

### 4.3.2. EFFECTIVENESS OF TEMPORAL INTEGRATION

To evaluate the integral component, we compare SEI with two single-epoch baselines: SE@T (signed entropy at the final epoch) and SE@T/2 (signed entropy at mid-training). Table 3 demonstrates that SEI achieves superior F1 scores, indicating that single snapshots provide unreliable signals while temporal integration yields robust detection.

The integral accumulates directional evidence over time, with each epoch contributing a signed cue: positive when

Table 1. Comparison with state-of-the-art mislabeled data detection methods under symmetric noise. Results are reported in F1 score (%) across five noise rates ( $\eta \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ ) on three medical datasets (ISIC, DeepDRiD, and PANDA). The best results for each dataset and noise rate are highlighted in **bold**.

	ISIC					DeepDRiD					PANDA				
	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5
INCV	38.60	42.72	47.21	55.93	61.67	33.05	44.56	49.29	57.87	63.32	53.30	57.69	61.23	68.45	67.11
BMM	30.49	39.59	40.65	55.89	55.73	30.23	39.44	46.11	50.66	56.21	54.37	62.88	66.19	71.09	72.43
GMM	36.31	47.49	49.67	64.79	67.10	35.89	45.39	52.28	58.87	66.27	59.60	62.79	68.24	71.11	70.76
AUM	48.65	62.97	72.60	77.70	81.67	39.66	54.96	60.92	69.20	75.75	65.23	72.39	75.95	75.30	77.18
CORES	36.20	56.90	67.84	76.72	82.67	26.21	42.68	56.38	69.22	68.26	60.73	63.61	66.28	73.43	70.51
CL	34.15	39.98	43.53	44.55	43.56	33.99	47.12	53.84	51.28	62.17	54.81	59.87	62.40	69.09	72.30
SIMIFEAT	32.01	39.33	43.95	44.82	41.62	36.04	45.72	53.67	55.06	61.53	55.82	61.15	68.12	70.56	69.76
DEFT	25.67	38.06	44.01	52.78	55.40	29.84	42.36	45.14	52.08	62.95	53.69	59.87	62.25	67.99	69.59
ReCoV	42.15	46.78	52.90	61.12	64.31	37.09	50.94	54.75	57.69	63.30	69.24	68.84	73.06	72.09	72.49
LEMOn	38.41	55.86	61.97	73.70	75.35	31.82	41.64	48.20	64.35	69.55	59.60	62.54	73.95	74.30	75.58
Ours	<b>50.44</b>	<b>65.64</b>	<b>74.80</b>	<b>80.07</b>	<b>83.93</b>	<b>45.20</b>	<b>56.53</b>	<b>63.43</b>	<b>71.26</b>	<b>78.19</b>	<b>72.57</b>	<b>77.22</b>	<b>81.46</b>	<b>83.11</b>	<b>81.18</b>

Table 2. Comparison with state-of-the-art mislabeled data detection methods under confusion-calibrated noise.

	ISIC					DeepDRiD					PANDA				
	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5
INCV	32.99	40.44	44.03	51.62	56.64	33.15	43.50	46.54	48.60	50.68	41.19	50.44	59.96	63.54	66.90
BMM	37.01	26.92	36.10	57.36	52.22	31.50	36.63	45.24	50.18	52.09	46.59	63.23	55.36	60.51	65.10
GMM	31.17	36.38	42.14	48.94	59.09	37.82	47.33	54.34	56.46	57.72	54.15	58.11	64.58	67.93	69.24
AUM	42.96	57.24	65.96	64.54	71.42	41.35	48.21	56.47	60.55	61.84	61.30	68.29	72.31	76.08	76.52
CORES	35.75	55.64	65.56	62.49	42.87	26.59	39.93	44.01	36.17	47.67	45.66	54.72	60.25	63.09	65.64
CL	38.14	43.44	48.66	45.46	43.28	34.06	37.60	43.06	42.78	51.58	51.87	60.77	69.19	71.53	72.35
SIMIFEAT	38.54	48.21	49.66	50.65	46.17	32.13	40.98	41.57	42.38	48.57	59.97	61.08	68.25	69.10	71.94
DEFT	15.33	24.67	32.72	37.03	34.41	18.46	28.74	39.68	46.53	50.34	40.66	54.72	65.48	67.83	69.93
ReCoV	40.90	44.21	61.16	64.66	66.61	42.63	48.49	50.93	52.67	61.79	59.50	63.15	64.58	65.91	69.87
LEMOn	30.59	42.50	56.70	61.91	67.07	34.06	42.62	55.80	63.86	66.82	56.40	66.66	70.22	74.80	75.13
Ours	<b>47.75</b>	<b>59.35</b>	<b>67.25</b>	<b>74.98</b>	<b>79.91</b>	<b>46.59</b>	<b>52.65</b>	<b>62.84</b>	<b>68.35</b>	<b>73.04</b>	<b>73.17</b>	<b>78.21</b>	<b>81.86</b>	<b>81.96</b>	<b>81.85</b>

predictions align with assigned labels, negative otherwise. Mislabeled samples accumulate predominantly negative values, while hard clean samples eventually offset early negative contributions through later positive ones. Single-epoch measurements suffer from training fluctuations, which temporal integration effectively smooths. This enables the integral to better capture long-term consistency patterns and reduces false detection of hard clean samples.

Besides, we analyze the timing of temporal evidence in Appendix F.2, comparing windowed integrals: SEI@Early (epochs 1–75) and SEI@Late (epochs 76–150).

#### 4.3.3. ARCHITECTURE GENERALIZABILITY

To assess the broader applicability of our approach beyond CLIP, we evaluate SEI with standard classification networks including ResNet-50 (He et al., 2016) and ViT (Dosovitskiy et al., 2021). Table 4 shows that while performance decreases, results remain competitive. This demonstrates that: (1) strong mislabeled data detection performance stems primarily from our proposed SEI rather than the CLIP architecture itself; (2) CLIP nevertheless provides advantages, likely due to its contrastive learning objective.

We also evaluate CLIP with various vision encoder backbone to assess generalization across both CNN and Transformer architectures, confirming consistent performance.

#### 4.3.4. SYNERGY BETWEEN SEI AND LEARNING WITH NOISY LABELS

To further demonstrate the utility of our approach, we integrate SEI as a data cleaning module with four representative noisy label learning methods: SCE (Wang et al., 2019), M-correction (Arazo et al., 2019), DivideMix (Li et al., 2020), and ProMix (Xiao et al., 2023). The resulting variants, SCE+SEI, M-correction+SEI, DivideMix+SEI, and ProMix+SEI are evaluated against their respective baselines on the ISIC dataset under confusion-calibrated noise. As shown in Figure 5, incorporating SEI consistently improves performance in F1 score (results on accuracy and AUC are provided in Appendix G). These gains highlight the plug-and-play nature of SEI: it is architecture-agnostic and integrates seamlessly into diverse noisy label learning frameworks.

Table 3. Ablation study on the effectiveness of the signed term and temporal integration. EI denotes the unsigned entropy integral (standard Shannon entropy), SE@T represents signed entropy at the final epoch, SE@T/2 denotes signed entropy at mid-training, and SEI is our full signed entropy integral method.

	ISIC					DeepDRiD					PANDA				
	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5
EI	35.31	45.05	50.46	60.77	63.65	42.91	47.86	53.85	59.28	61.14	51.76	61.50	64.26	67.26	65.86
SE@T	35.73	44.64	52.01	57.89	63.20	39.16	44.33	51.96	57.93	61.29	41.69	56.26	62.27	62.91	56.85
SE@T/2	36.97	48.19	56.30	63.72	68.73	40.32	48.87	54.89	62.84	66.58	57.35	60.06	64.80	66.26	66.88
SEI	<b>47.75</b>	<b>59.35</b>	<b>67.25</b>	<b>74.98</b>	<b>79.91</b>	<b>46.59</b>	<b>52.65</b>	<b>62.84</b>	<b>68.35</b>	<b>73.04</b>	<b>73.17</b>	<b>78.21</b>	<b>81.86</b>	<b>81.96</b>	<b>81.85</b>

Table 4. Architecture generalizability of SEI. We compare performance using standard classification networks (ResNet-50, ViT-B/16) and CLIP with different visual backbones (ResNet-50, ViT-B/16). Results demonstrate that SEI remains effective across diverse architectures.

	ISIC					DeepDRiD					PANDA				
	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5
RNet-50	45.68	57.58	65.07	73.04	76.33	43.99	51.34	61.17	66.89	71.75	71.21	75.64	79.21	79.28	78.25
ViT-B/16	44.25	56.71	64.57	72.12	75.91	42.88	50.47	60.71	65.68	71.31	69.44	74.17	78.89	78.86	76.03
CLIP <sub>RNet-50</sub>	<b>47.75</b>	<b>59.35</b>	<b>67.25</b>	<b>74.98</b>	<b>79.91</b>	<b>46.59</b>	<b>52.65</b>	<b>62.84</b>	<b>68.35</b>	<b>73.04</b>	<b>73.17</b>	<b>78.21</b>	<b>81.86</b>	<b>81.96</b>	<b>81.85</b>
CLIP <sub>ViT-B/16</sub>	46.66	58.36	66.89	73.59	77.25	44.58	51.78	61.27	67.74	72.41	71.92	76.52	80.63	81.28	79.45

Table 5. Mislabeled sample detection on CIFAR-100N. We report F1 scores (%) for identifying mislabeled samples. The best result is highlighted in **bold**.

Method	INCV	BMM	GMM	AUM	CORES	CL	SIMIFEAT	DeFT	ReCoV	LEMOn	Ours
F1 (%)	59.77	63.55	63.83	74.54	38.52	67.64	79.21	75.03	67.59	78.40	<b>81.41</b>

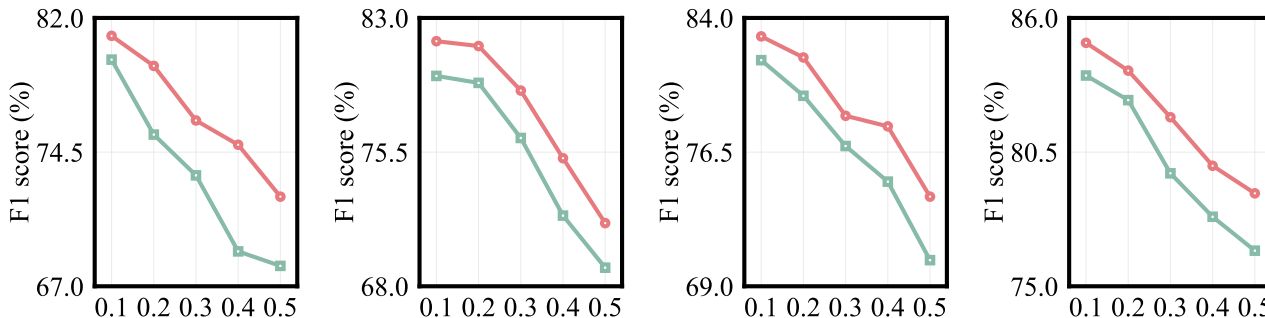


Figure 5. F1 score comparison between baseline noisy label learning methods (green) and their SEI-enhanced variants (red). From left to right: SCE, M-correction, DivideMix, and ProMix.

#### 4.3.5. GENERALIZATION OF SEI TO NATURAL IMAGES WITH REAL-WORLD NOISE

To assess the generalization ability of SEI beyond the medical domain and to evaluate its behavior under real-world human annotation noise, we conducted additional experiments on the CIFAR-100N. CIFAR-100N is a noisy-label variant of CIFAR-100 where each training image is re-annotated by human annotators, while the original CIFAR-100 labels are kept as clean ground truth.

On CIFAR-100N, we compare SEI against the same set of noisy label detection baselines used in the main paper. As shown in Table 5, SEI again achieves the best F1 score for mislabeled data detection, outperforming the second-best method by 2.2%. These results indicate that SEI is not

domain-specific. The signed entropy dynamics remain effective on a standard natural image dataset, and SEI robustly handles real human annotation noise—not only synthetic symmetric or confusion-calibrated noise.

## 5. Conclusion

We present SEI, a simple yet effective metric for detecting mislabeled data by leveraging signed entropy dynamics during training, which integrates seamlessly into standard training workflows. Extensive experiments on diverse medical imaging datasets demonstrate that SEI achieves state-of-the-art performance while remaining efficient and easy to apply.

## Impact Statement

This paper presents work whose goal is to advance the field of machine learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

## References

- Alderman, J. E., Palmer, J., Laws, E., McCradden, M. D., Ordish, J., Ghassemi, M., Pfohl, S. R., Rostamzadeh, N., Cole-Lewis, H., Glocker, B., et al. Tackling algorithmic bias and promoting transparency in health datasets: the standing together consensus recommendations. *The Lancet Digital Health*, 7:e64–e88, 2025.
- Arazo, E., Ortego, D., Albert, P., O’Connor, N. E., and McGuinness, K. Unsupervised label noise modeling and loss correction. In *Proceedings of the International Conference on Machine Learning*, pp. 312–321, 2019.
- Bae, H., Shin, S., Na, B., and Moon, I. Dirichlet-based per-sample weighting by transition matrix for noisy label learning. In *Proceedings of the 12th International Conference on Learning Representations*, 2024.
- Chen, J., Ramanathan, V., Xu, T., and Martel, A. L. Detecting noisy labels with repeated cross-validations. In *Proceedings of the Medical Image Computing and Computer Assisted Intervention*, pp. 197–207, 2024.
- Chen, P., Liao, B., Chen, G., and Zhang, S. Understanding and utilizing deep neural networks trained with noisy labels. In *Proceedings of the International Conference on Machine Learning*, pp. 1062–1070, 2019.
- Cheng, D., Ning, Y., Wang, N., Gao, X., Yang, H., Du, Y., Han, B., and Liu, T. Class-dependent label-noise learning with cycle-consistency regularization. In *Proceedings of the Advances in Neural Information Processing Systems*, pp. 11104–11116, 2022.
- Cheng, H., Zhu, Z., Li, X., Gong, Y., Sun, X., and Liu, Y. Learning with instance-dependent label noise: A sample sieve approach. In *Proceedings of the 9th International Conference on Learning Representations*, 2021.
- Cheng, L., Zhou, X., Zhao, L., Li, D., Shang, H., Zheng, Y., Pan, P., and Xu, Y. Weakly supervised learning with side information for noisy labeled images. In *Proceedings of the European Conference on Computer Vision*, pp. 306–321, 2020.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proceedings of the 9th International Conference on Learning Representations*, 2021.
- Engleson, E. and Azizpour, H. Robust classification via regression for learning with noisy labels. In *Proceedings of the 12th International Conference on Learning Representations*, 2024.
- Han, J., Luo, P., and Wang, X. Deep Self-Learning from noisy labels. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5138–5147, 2019.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Huang, J., Qu, L., Jia, R., and Zhao, B. O2U-Net: A simple noisy label detection approach for deep neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3325–3333, 2019.
- Jiang, S., Li, J., Wang, Y., Huang, B., Zhang, Z., and Xu, T. Delving into sample loss curve to embrace noisy and imbalanced data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 7024–7032, 2022.
- Kim, D., Ryoo, K., Cho, H., and Kim, S. Splitnet: Learnable clean-noisy label splitting for learning with noisy labels. *International Journal of Computer Vision*, 133(2):549–566, 2025.
- Kim, S., Lee, D., Kang, S., Chae, S., Jang, S., and Yu, H. Learning discriminative dynamics with label corruption for noisy label detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22477–22487, 2024.
- Kohl, S. A. A., Romera-Paredes, B., Meyer, C., Fauw, J. D., Ledsam, J. R., Maier-Hein, K. H., Eslami, S. M. A., Rezende, D. J., and Ronneberger, O. A probabilistic U-Net for segmentation of ambiguous images. In *NeurIPS*, 2018.
- Li, J., Socher, R., and Hoi, S. C. H. DivideMix: Learning with noisy labels as semi-supervised learning. In *Proceedings of the 8th International Conference on Learning Representations*, 2020.
- Li, W., Wang, L., Li, W., Agustsson, E., and Van Gool, L. Webvision database: Visual learning and understanding from web data. *arXiv preprint arXiv:1708.02862*, 2017.
- Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian, M., van der Laak, J. A. W. M., van Ginneken, B., and Sánchez, C. I. A survey on

- 495 deep learning in medical image analysis. *Medical Image*  
496 *Analysis*, 42:60–88, 2017.
- 497
- 498 Liu, R., Wang, X., Wu, Q., Dai, L., Fang, X., Yan, T., Son,  
499 J., Tang, S., Li, J., Gao, Z., Galdran, A., Poorneshwaran,  
500 J. M., Liu, H., Wang, J., Chen, Y., Porwal, P., Tan, G.  
501 S. W., Yang, X., Dai, C., Song, H., Chen, M., Li, H.,  
502 Jia, W., Shen, D., Sheng, B., and Zhang, P. Deepdrid:  
503 Diabetic retinopathy-grading and image quality estimation  
504 challenge. *Patterns*, 3:100512, 2022.
- 505
- 506 Liu, S., Niles-Weed, J., Razavian, N., and Fernandez-Granda,  
507 C. Early-learning regularization prevents memorization  
508 of noisy labels. In *Proceedings of the Advances in Neural*  
509 *Information Processing Systems*, pp. 20331–20342, 2020.
- 510
- 511 Lukasik, M., Bhojanapalli, S., Menon, A. K., and Kumar,  
512 S. Does label smoothing mitigate label noise? In  
513 *Proceedings of the International Conference on Machine*  
514 *Learning*, pp. 6448–6458, 2020.
- 515
- 516 Northcutt, C. G., Jiang, L., and Chuang, I. L. Confident  
517 learning: Estimating uncertainty in dataset labels. *Journal*  
518 *of Artificial Intelligence Research*, 70:1373–1411, 2021.
- 519
- 520 Pleiss, G., Zhang, T., Elenberg, E. R., and Weinberger,  
521 K. Q. Identifying mislabeled data using the area under the  
522 margin ranking. In *Proceedings of the Advances in Neural*  
523 *Information Processing Systems*, pp. 17044–17056, 2020.
- 524
- 525 Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G.,  
526 Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark,  
527 J., Krueger, G., and Sutskever, I. Learning transferable  
528 visual models from natural language supervision. In  
529 *Proceedings of the International Conference on Machine*  
530 *Learning*, pp. 8748–8763, 2021.
- 531
- 532 Ren, M., Zeng, W., Yang, B., and Urtasun, R. Learning  
533 to reweight examples for robust deep learning. In *Pro-*  
534 *ceedings of the International Conference on Machine*  
535 *Learning*, pp. 4334–4343, 2018.
- 536
- 537 Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S.,  
538 Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein,  
539 M., et al. Imagenet large scale visual recognition challenge.  
540 *International Journal of Computer Vision*, 115:211–252,  
541 2015.
- 542
- 543 Shi, J., Zhang, K., Guo, C., Yang, Y., Xu, Y., and Wu, J.  
544 A survey of label-noise deep learning for medical image  
545 analysis. *Medical Image Analysis*, 95:103166, 2024.
- 546
- 547 Song, H., Kim, M., and Lee, J. SELFIE: Refurbishing  
548 unclean samples for robust deep learning. In *Proceedings*  
549 *of the International Conference on Machine Learning*, pp.  
5907–5915, 2019.
- 500
- 501 Wang, H., Huang, Z., Lin, Z., and Liu, T. Noisegpt: La-  
502 bel noise detection and rectification through probability  
503 curvature. In *Proceedings of the Advances in Neural*  
504 *Information Processing Systems*, pp. 120159–120183,  
505 2024.
- 506
- 507 Wang, Y., Ma, X., Chen, Z., Luo, Y., Yi, J., and Bailey, J.  
508 Symmetric cross entropy for robust learning with noisy  
509 labels. In *Proceedings of the IEEE/CVF International*  
510 *Conference on Computer Vision*, pp. 322–330, 2019.
- 511
- 512 Wei, J., Zhu, Z., Cheng, H., Liu, T., Niu, G., and Liu,  
513 Y. Learning with noisy labels revisited: A study using  
514 real-world human annotations. In *Proceedings of the 10th*  
515 *International Conference on Learning Representations*,  
516 2022.
- 517
- 518 Wei, T., Li, H., Li, C., Shi, J., Li, Y., and Zhang, M.  
519 Vision-Language models are strong noisy label detectors.  
520 In *Proceedings of the Advances in Neural Information*  
521 *Processing Systems*, 2024.
- 522
- 523 Xia, X., Liu, T., Han, B., Gong, C., Wang, N., Ge, Z.,  
524 and Chang, Y. Robust early-learning: Hindering the  
525 memorization of noisy labels. In *Proceedings of the 9th*  
526 *International Conference on Learning Representations*,  
527 2021.
- 528
- 529 Xiao, R., Dong, Y., Wang, H., Feng, L., Wu, R., Chen,  
530 G., and Zhao, J. ProMix: Combating label noise via  
531 maximizing clean sample utility. In *Proceedings of the*  
532 *International Joint Conference on Artificial Intelligence*,  
533 pp. 4442–4450, 2023.
- 534
- 535 Xiao, T., Xia, T., Yang, Y., Huang, C., and Wang, X.  
536 Learning from massive noisy labeled data for image  
537 classification. In *Proceedings of the IEEE conference on*  
538 *Computer Vision and Pattern Recognition*, pp. 2691–2699,  
539 2015a.
- 540
- 541 Xiao, T., Xia, T., Yang, Y., Huang, C., and Wang, X.  
542 Learning from massive noisy labeled data for image  
543 classification. In *Proceedings of the IEEE conference on*  
544 *Computer Vision and Pattern Recognition*, pp. 2691–2699,  
545 2015b.
- 546
- 547 Yang, H., Jin, Y.-Z., Li, Z.-Y., Wang, D.-B., Geng, X., and  
548 Zhang, M.-L. Learning from noisy labels via dynamic  
549 loss thresholding. *IEEE Transactions on Knowledge and*  
550 *Data Engineering*, 2023.
- 551
- 552 Ye, X., Li, X., Dai, S., Liu, T., Sun, Y., and Tong, W. Active  
553 negative loss functions for learning with noisy labels.  
554 In *Proceedings of the Advances in Neural Information*  
555 *Processing Systems*, pp. 6917–6940, 2023.

550 Yuan, S., Feng, L., and Liu, T. Early stopping against label  
551 noise without validation data. In *Proceedings of the 12th*  
552 *International Conference on Learning Representations*,  
553 2024.

554 Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O.  
555 Understanding deep learning requires rethinking general-  
556 ization. In *Proceedings of the 5th International Confer-*  
557 *ence on Learning Representations*, 2017.

559 Zhang, H., Balagopalan, A., Oufattole, N., Jeong, H., Wu,  
560 Y., Zhu, J., and Ghassemi, M. LEMoN: Label error  
561 detection using multimodal neighbors. In *Proceedings of*  
562 *the International Conference on Machine Learning*, 2025.

564 Zhang, Z. and Sabuncu, M. R. Generalized cross entropy  
565 loss for training deep neural networks with noisy labels.  
566 In *Proceedings of the Advances in Neural Information*  
567 *Processing Systems*, pp. 8792–8802, 2018.

568 Zhu, Z., Dong, Z., and Liu, Y. Detecting corrupted labels  
569 without training a model to predict. In *Proceedings of*  
570 *the International Conference on Machine Learning*, pp.  
571 27412–27427, 2022.

572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604

## A. Reproducibility Statement

Our code is available at <https://anonymous.4open.science/r/SEI-14123>.

## B. Datasets

Figure 6 presents representative samples from each class across the three datasets employed in our study: ISIC, DeepDRiD, and PANDA. We display one exemplar image per class, organized with rows corresponding to individual datasets and columns representing distinct classes. This visualization facilitates direct comparison of class-specific visual characteristics. The corresponding text prompts utilized for training CLIP models are detailed in Table 6, including auxiliary class prompts for each dataset.

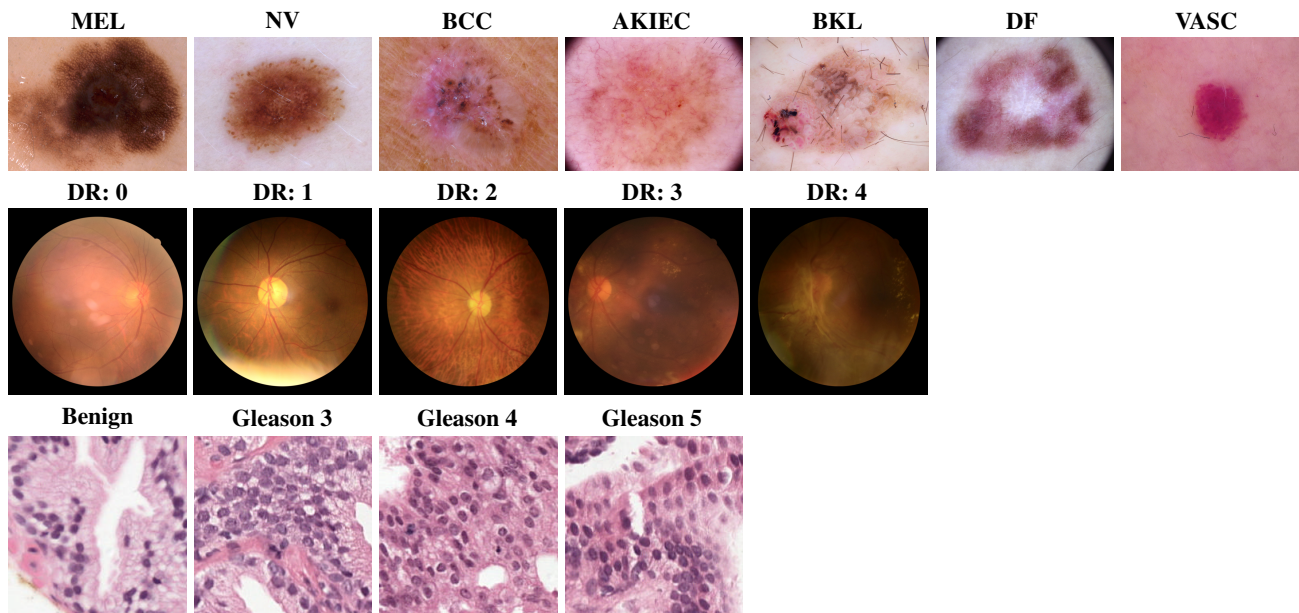


Figure 6. Representative images from the datasets used in this work. Rows correspond to ISIC, DeepDRiD, and PANDA datasets (top to bottom) and columns to class labels.

## C. Additional Analysis of Entropy Trajectories

In this section, we visualize entropy trajectories for more representative categories across the three datasets. For ISIC and DeepDRiD, we additionally show trajectories for other classes, comparing correctly labeled samples with mislabeled ones. As shown in Figure 7, we plot entropy trajectories for melanoma cases from the ISIC dataset, grade 0 diabetic retinopathy images from the DeepDRiD dataset, benign epithelium samples from the PANDA dataset, and Gleason 5 cancerous epithelium images from the PANDA dataset. Across all examined cases, we consistently observe the regularity described in Section 2.3.

## D. Extended Evaluation of SEI for Mislabeled Sample Detection

We present additional empirical evidence demonstrating the effectiveness of SEI. Figures 8, 9, and 10 illustrate the discriminative power of SEI in separating different sample types within the ISIC, DeepDRiD, and PANDA datasets, respectively. The results consistently validate the theoretical framework outlined in Section 2.4.2: samples with correct labels that are easily classified exhibit large positive SEI values, challenging but correctly labeled samples demonstrate moderate SEI values, while mislabeled samples consistently display strongly negative SEI values.

Table 6. Text prompts for each class in the datasets. Prompts highlighted in gray are auxiliary class prompts.

Dataset	Text Prompts
ISIC	A dermoscopic image showing melanoma.
	A dermoscopic image showing nevus.
	A dermoscopic image showing basal cell carcinoma.
	A dermoscopic image showing actinic keratosis/intraepithelial carcinoma.
	A dermoscopic image showing benign keratosis.
	A dermoscopic image showing dermatofibroma.
	A dermoscopic image showing vascular lesion.
	A dermoscopic image showing other lesions.
DeepDRiD	A fundus image showing no evidence of diabetic retinopathy.
	A fundus image exhibiting mild diabetic retinopathy.
	A fundus image exhibiting moderate diabetic retinopathy.
	A fundus image exhibiting severe diabetic retinopathy.
	A fundus image exhibiting proliferative diabetic retinopathy.
	A fundus image showing other retinal conditions.
PANDA	A histology image showing benign glandular epithelium.
	A histology image showing Gleason pattern 3 adenocarcinoma.
	A histology image showing Gleason pattern 4 adenocarcinoma.
	A histology image showing Gleason pattern 5 adenocarcinoma.
	A histology image showing other conditions.

## E. Theoretical Properties of Signed Entropy

In this section, we provide a short theoretical analysis of the proposed signed entropy (Eq. 2). Recall that for  $(\mathbf{x}, y) \in \mathcal{D}_{\text{train}}$  with posterior  $\mathbf{p}(\mathbf{x})$ , we define

$$\mathcal{H}(\mathbf{p}(\mathbf{x}), y) = (-1)^{\mathbb{1}[y = \arg \max_k p_k(\mathbf{x})]} \sum_{k=1}^K p_k(\mathbf{x}) \log p_k(\mathbf{x}).$$

### E.1. Relation to Shannon Entropy

**Proposition E.1** (Reduction to Shannon Entropy). *If  $y = \arg \max_k p_k(\mathbf{x})$ , then*

$$\mathcal{H}(\mathbf{p}(\mathbf{x}), y) = H(\mathbf{p}(\mathbf{x})),$$

where  $H(\mathbf{p}) = -\sum_k p_k \log p_k$  is Shannon’s entropy.

*Proof.* By definition, the sign exponent equals  $(-1)^1 = -1$  when the prediction agrees with  $y$ , yielding the standard Shannon entropy.  $\square$

### E.2. Concavity and Sign Symmetry

**Proposition E.2** (Concavity up to Sign). *Let  $\mathcal{P}$  denote the probability simplex in  $\mathbb{R}^K$ . For a fixed label  $y$ , the signed entropy  $\mathcal{H}(\cdot, y)$  is concave on  $\mathcal{P}$  if  $y$  matches the prediction, and is convex on  $\mathcal{P}$  if  $y$  disagrees with the prediction.*

*Proof.* The Shannon entropy  $H(\mathbf{p})$  is strictly concave on  $\mathcal{P}$  (classical result). Multiplying by  $-1$  flips concavity to convexity. Since the sign of  $\mathcal{H}$  depends only on alignment, the stated property follows.  $\square$

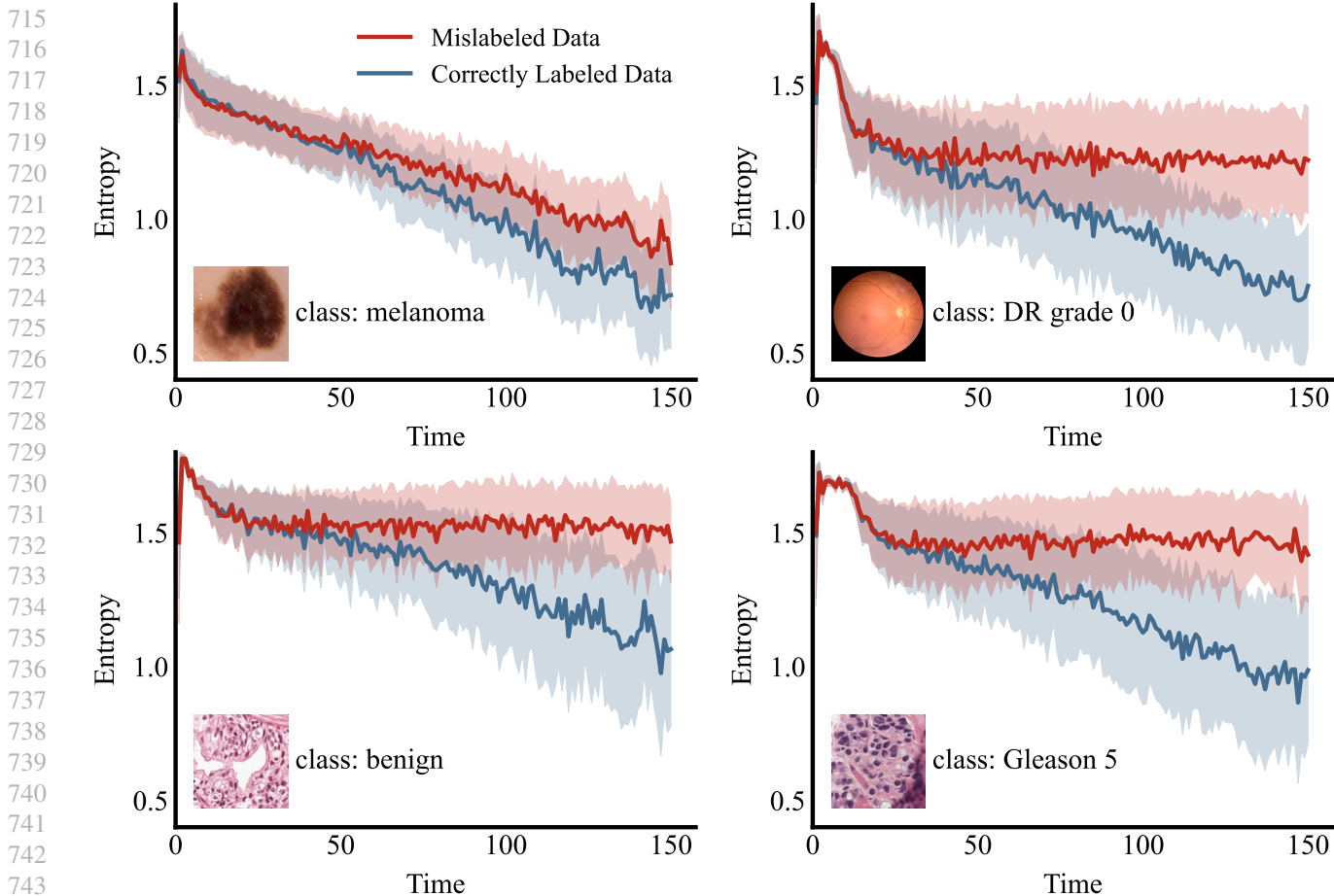


Figure 7. Training dynamics of prediction entropy for correctly labeled versus mislabeled samples. Top-left: melanoma images from the ISIC dataset. Top-right: grade 0 diabetic retinopathy (DR) images from the DeepDRiD dataset. Bottom-left: benign glandular epithelium images from the PANDA dataset. Bottom-right: Gleason 5 cancerous epithelium images from the PANDA dataset.

### E.3. Implication for SEI

These properties imply that SEI (Eq. 3) can be interpreted as a signed, temporally averaged measure of prediction uncertainty. Its sign encodes long-term label alignment, while its magnitude captures how confidently the model reaches this alignment (or misalignment). This dual role is what enables SEI to separate mislabeled data from both easy and hard clean samples.

## F. Additional Ablation Results

### F.1. Visual Evidence for the Signed Term

We further visualize score distributions for clean and mislabeled samples under both formulations (see Figure 11). The unsigned variant exhibits heavily overlapping positive-only distributions, making separation difficult. In contrast, SEI introduces a clear negative tail for mislabeled samples, creating a bimodal structure with reduced overlap and more distinguishable groups.

### F.2. Full-Trajectory Integration vs. Windowed Integrals

Table 7 shows that both windowed variants (SEI@Early and SEI@Late) perform worse than the full-trajectory SEI. Restricting to either an early or late window discards complementary cues present in the other phase.

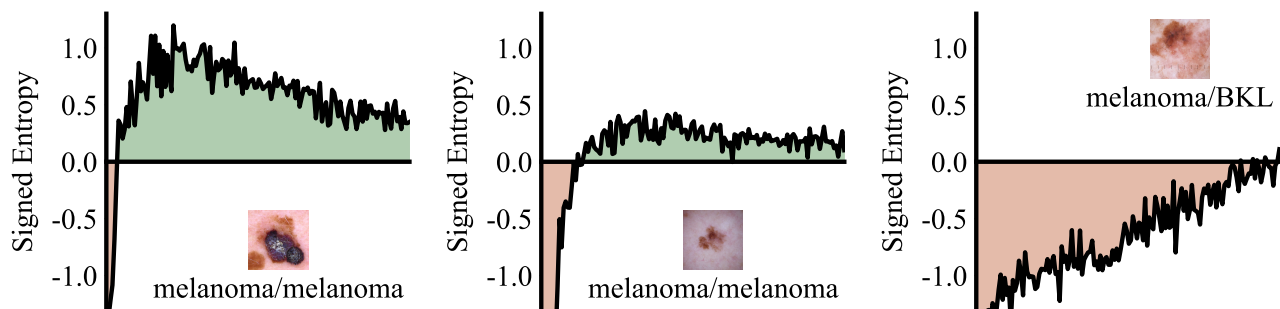


Figure 8. Illustration of SEI using melanoma images from the ISIC dataset. The plots show signed entropy curves across training epochs for easy clean (left), hard clean (middle), and mislabeled (right) samples. Each curve is averaged over 200 samples, and the signed area under the curve represents the SEI. Correctly labeled samples consistently exhibit larger SEIs than mislabeled ones.

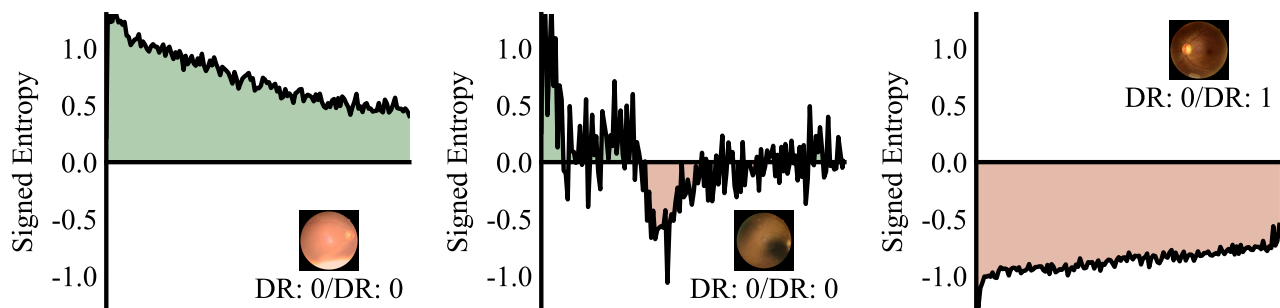


Figure 9. Illustration of SEI using Grade 0 diabetic retinopathy (DR) images from the DeepDRiD dataset.

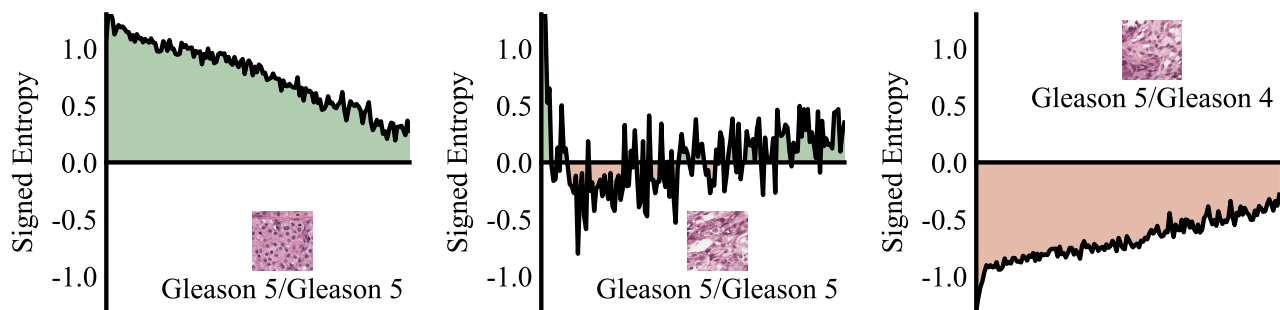


Figure 10. Illustration of SEI using Gleason 5 images from the PANDA dataset.

Table 7. Evaluation of windowed integrals—SEI@Early (epochs 1–75) and SEI@Late (epochs 76–150)—compared to the full SEI.

	ISIC					DeepDRiD					PANDA				
	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5	0.1	0.2	0.3	0.4	0.5
SE@Early	44.88	56.15	63.50	69.88	73.80	36.05	42.69	54.30	60.63	66.48	68.01	70.89	75.05	76.36	76.11
SE@Late	37.98	51.23	62.57	65.77	69.56	41.11	47.08	58.81	63.96	70.81	70.46	73.50	77.10	77.91	78.66
SEI	<b>47.75</b>	<b>59.35</b>	<b>67.25</b>	<b>74.98</b>	<b>79.91</b>	<b>46.59</b>	<b>52.65</b>	<b>62.84</b>	<b>68.35</b>	<b>73.04</b>	<b>73.17</b>	<b>78.21</b>	<b>81.86</b>	<b>81.96</b>	<b>81.85</b>

## G. Additional Results on Downstream Image Classification

Beyond F1 score, we also report accuracy and AUC for SCE, M-correction, DivideMix, and ProMix, both with and without SEI, under the same protocol as Section 4.3.4. Results are presented in Figure 12 and Figure 13.

In addition, we provide comprehensive results (F1 score, accuracy, and AUC) under symmetric noise, which are reported in Figures 14, 15, and 16.

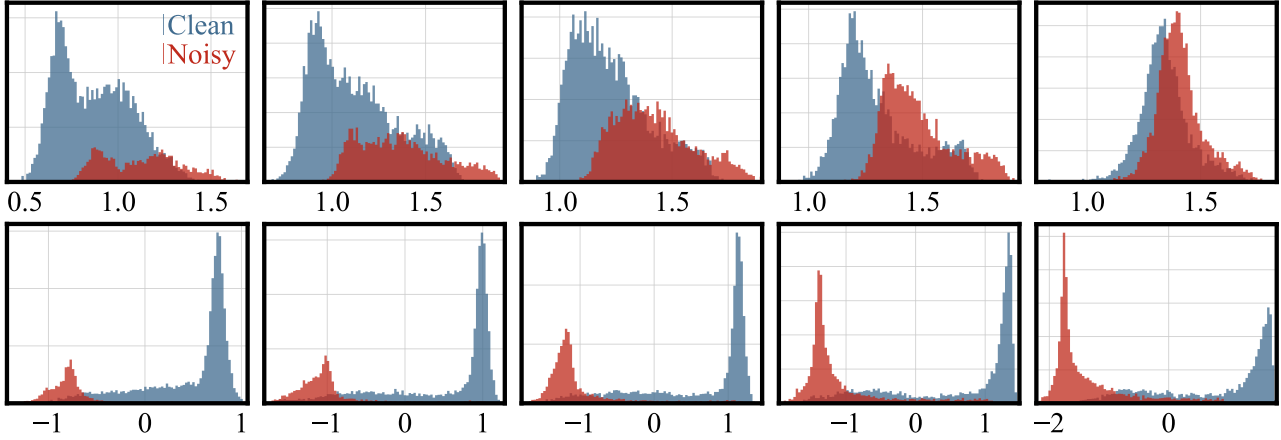


Figure 11. Score distributions for clean and noisy samples on the PANDA dataset. We compare our proposed SEI statistic (bottom row) against the unsigned Shannon entropy integral baseline (top row) for noise rates  $\eta \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ . Each column corresponds to increasing noise levels from left to right. The SEI statistic demonstrates better separation between clean and noisy sample distributions at all noise levels.

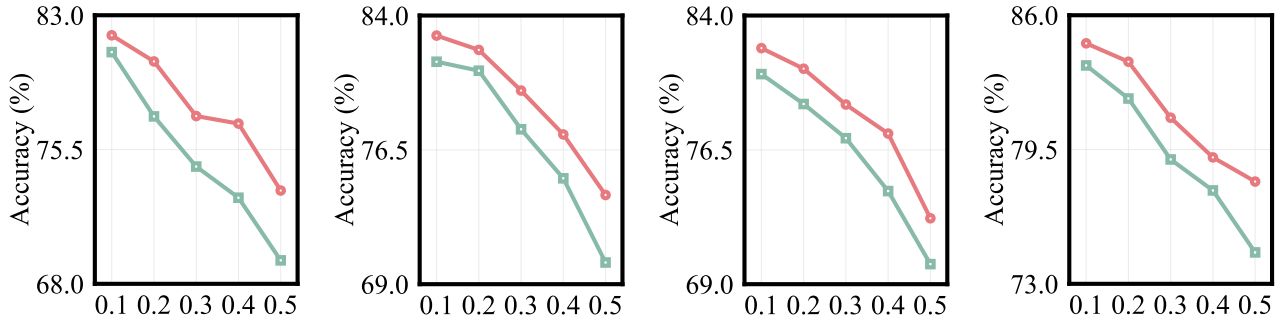


Figure 12. Accuracy comparison under confusion-calibrated noise between baseline noisy label learning methods (green) and their SEI-enhanced variants (red). From left to right: SCE, M-correction, DivideMix, and ProMix.

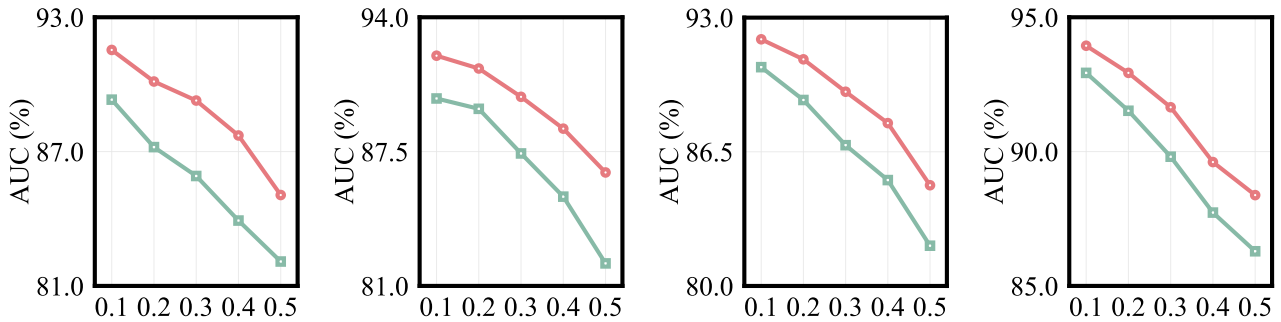


Figure 13. AUC comparison under confusion-calibrated noise between baseline noisy label learning methods (green) and their SEI-enhanced variants (red). From left to right: SCE, M-correction, DivideMix, and ProMix.

## H. Sensitivity to Training Configuration

We conducted a series of hyperparameter sensitivity experiments on the PANDA dataset under the confusion-calibrated 50% noise setting to assess how SEI behaves under different training configurations.

**Training duration and early stopping.** We first varied the number of training epochs while keeping all other hyperparameters fixed. As shown in Table 8, the small standard deviation (0.42%) indicates that SEI is insensitive to moderate changes in training duration as long as the model reaches convergence. In contrast, early stopping at epoch 75 leads to a noticeable drop (Table 7), confirming that SEI assumes a reasonably converged model rather than an undertrained one.

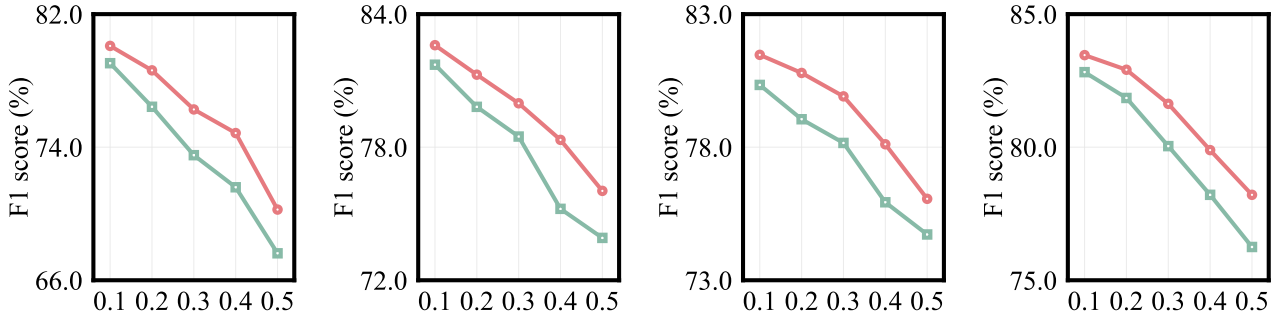


Figure 14. F1 score comparison under symmetric noise between baseline noisy label learning methods (green) and their SEI-enhanced variants (red). From left to right: SCE, M-correction, DivideMix, and ProMix.

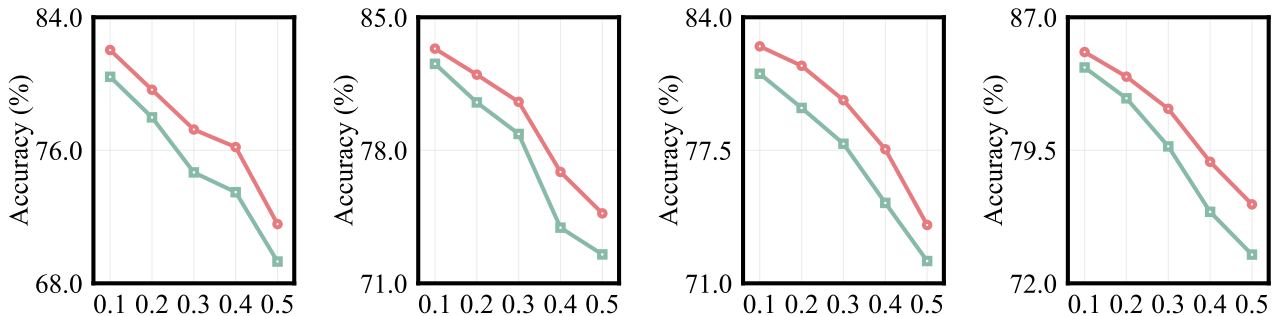


Figure 15. Accuracy comparison under symmetric noise between baseline noisy label learning methods (green) and their SEI-enhanced variants (red). From left to right: SCE, M-correction, DivideMix, and ProMix.

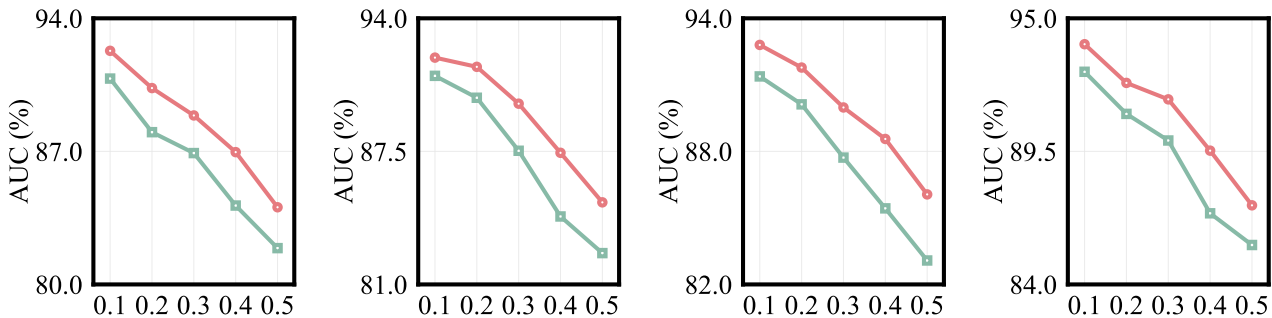


Figure 16. AUC comparison under symmetric noise between baseline noisy label learning methods (green) and their SEI-enhanced variants (red). From left to right: SCE, M-correction, DivideMix, and ProMix.

Table 8. Sensitivity of SEI to training duration on PANDA (confusion-calibrated noise,  $\eta = 0.5$ ). The default setting is 150 epochs.

Epochs	100	150	200	mean	std
F1 (%)	81.34	81.85	81.02	81.40	0.42

**Learning rate.** We investigate the sensitivity of SEI to the choice of initial learning rate by sweeping it from  $1e-4$  to  $1e-2$ . As shown in Table 9, we observe that the performance of SEI remains very stable within the range typically used for training CLIP and standard classifiers ( $5e-4$ – $1e-3$ ), and only starts to degrade when the learning rate is excessively large ( $5e-3$  or  $1e-2$ ). This trend mirrors standard classification training, where an appropriate learning rate is needed for effective learning of the base model itself.

**Data augmentation, label smoothing, and weight decay.** We also examined several common regularization choices. For data augmentation, as shown in Table 10, “Strong aug.” adds random affine transforms and random erasing. SEI remains

Table 9. Sensitivity of SEI to learning rate on PANDA (confusion-calibrated noise,  $\eta = 0.5$ ). The default learning rate is  $1e-3$ .

Learning rate	1e-4	5e-4	1e-3	5e-3	1e-2
F1 (%)	80.06	81.50	81.85	78.86	77.00

robust under both weak and strong augmentations, with a slight gain under stronger augmentation. MixUp—whose heavy mixing disrupts label–prediction alignment—slightly reduces performance. For label smoothing, as shown in Table 11, moderate label smoothing slightly degrades performance, and stronger smoothing further weakens entropy signals, as expected. As shown in Table 12, weight decay within [5e-5, 5e-4] has only mild impact, with slightly higher F1 at the upper end.

Table 10. Effect of data augmentation on SEI (PANDA, confusion-calibrated noise,  $\eta = 0.5$ ). The default setting is weak augmentation.

Augmentation setting	Strong aug.	Weak aug.	MixUp
F1 (%)	82.38	81.85	80.88

Table 11. Effect of label smoothing on SEI (PANDA, confusion-calibrated noise,  $\eta = 0.5$ ). The default setting is 0.0.

Label smoothing	0.0	0.1	0.2
F1 (%)	81.85	81.45	80.36

Table 12. Effect of weight decay on SEI (PANDA, confusion-calibrated noise,  $\eta = 0.5$ ). The default setting is  $1e-4$ .

Weight decay	5e-5	1e-4	5e-4
F1 (%)	81.36	81.85	82.21

**Temporal integration window.** We compare the full-trajectory SEI with variants that integrate signed entropy over fixed 50-epoch windows with a stride of 20. As reported in Table 13, all window-based variants perform significantly worse than full-trajectory SEI, consistent with the SEI@Early/SEI@Late results in Table 7. This confirms that SEI benefits from integrating the entire training trajectory rather than relying on a narrow early or late slice.

In summary, SEI is robust under standard training-to-convergence settings and typical hyperparameter choices. The ranking is stable with respect to reasonable training lengths, LR schedules, and regularization. Our recommendation is to use the full training trajectory; extreme early stopping or highly atypical hyperparameters are not ideal conditions for SEI.

## I. Sensitivity to the Auxiliary-Class Sampling Ratio

To evaluate the robustness of the auxiliary-class–based threshold, we performed a sensitivity analysis on the PANDA dataset with 50% confusion-calibrated noise by varying the sampling ratio of auxiliary-class samples. In our default setting, the number of auxiliary samples is  $N/(K+1)$ . We scaled this number using factors of  $0.5\times$ ,  $0.75\times$ ,  $1.0\times$ ,  $1.5\times$ , and  $2.0\times$ . As shown in Table 14, the F1 scores remain within a tight range, with a mean of 81.87% and a standard deviation of only 0.72%, even when the number of auxiliary samples varies by a factor of four. This indicates that SEI is not sensitive to the exact sampling ratio: as long as a reasonable number of auxiliary samples is used, the estimated mean SEI for the auxiliary class remains stable. It is also worth noting that although the sampling ratio is manually specified, the threshold itself is entirely data-driven and learned adaptively.

## J. Inter-Obse Variability vs. Label Noise

Inter-observer variability and label noise are related but not identical. Variability across annotators often reflects inherent uncertainty, whereas noisy label detection assumes a single hard label per sample and seeks to identify cases where that label is incorrect. Learning with uncertain or probabilistic labels—such as modeling annotator distributions or ambiguity—is an important but distinct problem setting (Kohl et al., 2018) and is not the focus of this work.

Table 13. SEI with different 50-epoch integration windows (PANDA, confusion-calibrated noise,  $\eta = 0.5$ ).

Window (epochs)	[0, 49]	[20, 69]	[40, 89]	[60, 109]	[80, 129]	[100, 149]
F1 (%)	75.34	50.78	49.92	51.03	54.09	56.81

Table 14. Sensitivity of SEI to the auxiliary-class sampling ratio on PANDA (confusion-calibrated noise,  $\eta = 0.5$ ). The ratio is expressed as a multiplier of the default  $N/(K + 1)$  setting. The default setting is 1.0 $\times$ .

Ratio factor	0.5 $\times$	0.75 $\times$	1.0 $\times$	1.5 $\times$	2.0 $\times$
F1 (%)	81.94	83.01	81.85	81.41	81.12

Our work operates strictly under the standard hard-label noise detection setting, where each training sample is associated with one label, and the goal is to detect mislabeled instances under this assumption. SEI is therefore designed and evaluated within this framework.

## K. SEI under Extreme Noise Rates

In clinical practice, datasets with more than 50% label disagreement are typically considered unreliable for training. Accordingly, our main experiments focus on noise rates up to 0.5. To further assess robustness, we additionally evaluate SEI on the PANDA dataset under more extreme noise levels  $\eta \in \{0.6, 0.7, 0.8\}$ , for both confusion-calibrated noise and symmetric noise. The F1 scores for mislabeled-data detection are summarized in Table 15.

Table 15. SEI performance on PANDA under extreme noise rates. We report F1 scores (%) for mislabeled data detection.

Noise rate	0.6	0.7	0.8
Confusion-calibrated	81.07	79.93	81.83
Symmetric	81.66	80.95	80.51

Even at very high noise rates (60–80%), SEI remains stable around 80% F1, without a significant performance collapse. Performance does not degrade monotonically; instead, it fluctuates slightly within a narrow band, suggesting that SEI can still capture useful training-dynamics signals even when a large fraction of labels is corrupted.

## L. Per-Class Analysis of Noisy-Label Detection

We perform a per-class analysis of noisy label detection performance on PANDA with 50% confusion-calibrated noise. For each class, we report the class-wise false positive rate (FPR) and per-class F1 score as shown in Table 16.

These results show that SEI does not systematically over-filter any particular class. Notably, benign epithelium—the smallest class—has the lowest false positive rate (13.66%), suggesting that minority patterns are not disproportionately removed. Overall, SEI achieves consistently strong detection quality across all classes.

This outcome aligns with the design of SEI: by integrating signed entropy over the entire training trajectory, SEI naturally separates hard-but-correct samples from mislabeled ones. Hard, correctly labeled samples may have high entropy early in training but eventually align with their labels and accumulate positive signed contributions. In contrast, mislabeled samples remain misaligned for most of training and accumulate negative contributions. This reduces the risk of misclassifying intrinsically hard or minority-subtype samples as noisy.

1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099

Table 16. Per-class false positive rate and F1 for noisy-label detection on PANDA (confusion-calibrated noise,  $\eta = 0.5$ ).

<b>Class</b>	<b>Per-class FPR (%)</b>	<b>Per-class F1 (%)</b>
Benign epithelium	13.66	85.47
Gleason 3	22.87	80.38
Gleason 4	25.87	80.14
Gleason 5	21.63	82.91