

TiDAL: Learning Training Dynamics for Active Learning

Seong Min Kye^{1,†} Kwanghee Choi^{2,†} Hyeonmin Byun¹ Buru Chang^{3,*}
¹Hyperconnect ²Carnegie Mellon University ³Sogang University

{harris,hyeonmin.byun}@hpcnt.com kwanghec@andrew.cmu.edu buru@sogang.ac.kr

Abstract

Active learning (AL) aims to select the most useful data samples from an unlabeled data pool and annotate them to expand the labeled dataset under a limited budget. Especially, uncertainty-based methods choose the most uncertain samples, which are known to be effective in improving model performance. However, previous methods often overlook training dynamics (TD), defined as the ever-changing model behavior during optimization via stochastic gradient descent, even though other research areas have empirically shown that TD provides important clues for measuring the data uncertainty. In this paper, we first provide theoretical and empirical evidence to argue the usefulness of utilizing the ever-changing model behavior rather than the fully trained model snapshot. We then propose a novel AL method, Training Dynamics for Active Learning (TiDAL), which efficiently predicts the training dynamics of unlabeled data to estimate their uncertainty. Experimental results show that our TiDAL achieves better or comparable performance on both balanced and imbalanced benchmark datasets compared to state-of-the-art AL methods, which estimate data uncertainty using only static information after model training.

1. Introduction

“There is a tide in the affairs of men. Which taken at the flood, leads on to fortune.” — William Shakespeare

Active learning (AL) [5, 31] aims to solve the real-world problem of selecting the most useful data samples from large-scale unlabeled data pools and annotating them to expand labeled data under a limited budget. Since the current deep neural networks are data-hungry, AL has increasingly gained attention in recent years. Existing AL methods can be divided into two mainstream categories: diversity- and uncertainty-based methods. Diversity-based methods [42, 14] focus on constructing a subset that fol-

[†]Equal contribution.

^{*}Corresponding author.

This work was done while all authors were affiliated with Hyperconnect.

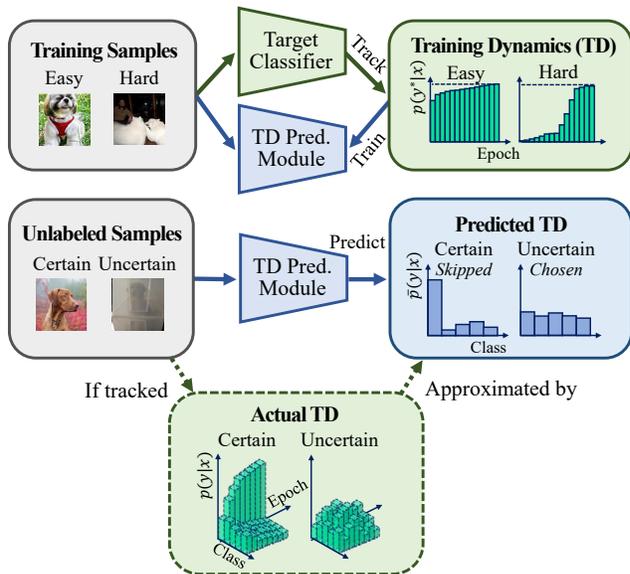


Figure 1: Our proposed TiDAL. TD of training samples x may differ even if they converge to the same final predicted probability $p(y^*|x)$ (Upper row). Hence, we are motivated to utilize the readily available rich information generated during training, *i.e.*, leveraging TD. We estimate TD of large-scale unlabeled data using a prediction module instead of tracking the actual TD of all the unlabeled samples to avoid the computational overhead (Lower row).

lows the target data distribution. Uncertainty-based methods [13, 6, 52] choose the most uncertain samples, which are known to be effective in improving model performance. Hence, the most critical question for the latter becomes, “How can we quantify the data uncertainty?”

In this paper, we leverage **training dynamics** (TD) to quantify data uncertainty. TD is defined as the ever-changing model behavior on each data sample during optimization via stochastic gradient descent. Recent studies [9, 29, 48, 47] have provided empirical evidence that TD provides important clues for measuring the contribution of each data sample to model performance improvement. Inspired by these studies, we argue that the data uncertainty

of unlabeled data can be estimated with TD. However, most uncertainty-based methods quantify data uncertainty based on static information (*e.g.*, loss [52] or predicted probability [45]) from a fully-trained model “*snapshot*,” neglecting the valuable information generated during training. We further argue that TD is more effective in separating uncertain and certain data than static information from a model snapshot captured after model training. In §3, we provide both theoretical and empirical evidence to support our argument that TD is a valuable tool for quantifying data uncertainty.

Despite its huge potential, TD is not yet actively explored in the domain of AL. This is because AL assumes a massive unlabeled data pool. Previous studies track TD only for the training data every epoch as it can be recorded easily during model optimization. On the other hand, AL targets a large number of unlabeled data, where ***tracking the TD for each unlabeled sample requires an impractical amount of computation*** (*e.g.*, inference all the unlabeled samples every training epoch).

Therefore, we propose TiDAL (*Training Dynamics for Active Learning*), a novel AL method that efficiently quantifies the uncertainty of unlabeled data by estimating their TD. We avoid tracking the TD of large-scale unlabeled data every epoch by predicting the TD of unlabeled samples with a TD prediction module. The module is trained with the TD of labeled data, which is readily available during model optimization. During the data selection phase, we predict the TD of unlabeled data with the trained module to quantify their uncertainties. We efficiently obtain TD using the module, which avoids inferring all the unlabeled samples every epoch. Experimental results demonstrate that our TiDAL achieves better or comparable performance to existing AL methods on both balanced and imbalanced datasets. Additional analyses show that our prediction module successfully predicts TD, and the predicted TD is useful in estimating uncertainties of unlabeled data. Our proposed method are illustrated in Figure 1.

Contributions of our study: (1) We bridge the concept of training dynamics and active learning with the theoretical and experimental evidence that training dynamics is effective in estimating data uncertainty. (2) We propose a new method that efficiently predicts the training dynamics of unlabeled data to estimate their uncertainty. (3) Our proposed method achieves better or comparable performance on both balanced and imbalanced benchmark datasets compared to existing active learning methods. For reproducibility, we release the source code¹.

2. Preliminaries

To better understand our proposed method, we first summarize key concepts, including uncertainty-based active

learning, quantification of uncertainty, and training dynamics.

Uncertainty-based active learning. In this work, we focus on uncertainty-based AL for multi-class classification problems. We define the predicted probabilities of the given sample x for C classes as:

$$\mathbf{p} = [p(1|x), p(2|x), \dots, p(C|x)]^T \in [0, 1]^C, \quad (1)$$

where we denote the true label of x as y and the classifier as f . \mathcal{D} and \mathcal{D}_u denote a labeled dataset and an unlabeled data pool, respectively. The general cycle of uncertainty-based AL is in two steps: (1) train the target classifier f on the labeled dataset \mathcal{D} and (2) select top- k uncertain data samples from the unlabeled data pool \mathcal{D}_u . Selected samples are then given to the human annotators to expand the labeled dataset \mathcal{D} , cycling back to the first step.

Quantifying uncertainty. The objective of this study is to establish a connection between the concept of TD and the field of AL. In order to clearly demonstrate the effectiveness of utilizing TD to quantify data uncertainty, we have employed two of the most prevalent and straightforward estimators, *entropy* [43] and *margin* [41], to measure data uncertainty in this paper. Entropy H is defined as follows:

$$H(\mathbf{p}) = - \sum_{c=1}^C p(c|x) \log p(c|x), \quad (2)$$

where the sample x is from the unlabeled data pool \mathcal{D}_u . Entropy concentrates on the level of the model’s confidence on the given sample x and gets bigger when the prediction across the classes becomes uniform (*i.e.*, uncertain). Margin M measures the difference between the probability of the true label and the maximum of the others:

$$M(\mathbf{p}) = p(y|x) - \max_{c \neq y} p(c|x), \quad (3)$$

where y denotes the true label. The smaller the margin, the lower the model’s confidence in the sample, so it can be considered uncertain. Both entropy and margin are computed with the predicted probabilities \mathbf{p} of the fully trained classifier f , only taking the snapshot of f into account.

Defining training dynamics. Our TiDAL targets to leverage TD of unlabeled data to estimate their uncertainties. TD can be defined as any model behavior during optimization, such as the area under the margin between logit values of the target class and the other largest class [39] or the variance of the predicted probabilities generated at each epoch [47]. In this work, we define the TD $\bar{\mathbf{p}}^{(t)}$ as the area under the predicted probabilities of each data sample x obtained

¹<https://github.com/hyperconnect/TiDAL>

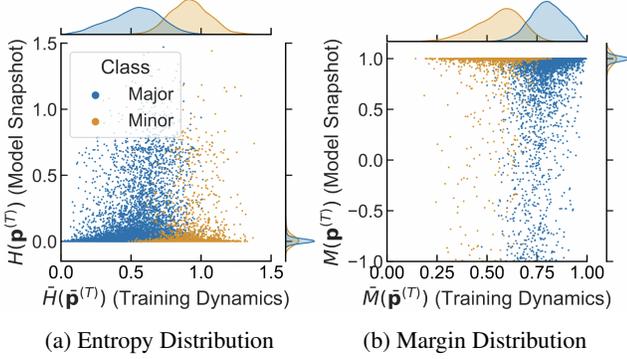


Figure 2: Score distribution after long-tailed training. We plot the marginal distributions using kernel density estimation (KDE). It is difficult to separate major (certain) and minor (uncertain) samples by the model snapshot-based scores (horizontal), unlike the TD-driven scores (vertical) that enable clearly separating the certain and uncertain samples.

during the t time steps of optimizing the target classifier f :

$$\begin{aligned} \mathbf{p}^{(i)} &= [p^{(i)}(1|x), p^{(i)}(2|x), \dots, p^{(i)}(C|x)]^T, \quad (4) \\ \bar{\mathbf{p}}^{(t)} &= [\bar{p}^{(t)}(1|x), \bar{p}^{(t)}(2|x), \dots, \bar{p}^{(t)}(C|x)]^T \\ &= \sum_{\tau} \mathbf{p}^{(\tau)} \Delta\tau \simeq \sum_{i=1}^t \mathbf{p}^{(i)} / t, \quad (5) \end{aligned}$$

where $\mathbf{p}^{(i)}$ is the predicted probabilities of a target classifier f at the i -th time step. $\Delta\tau$ is the unit time step to normalize the predicted probabilities. For simplicity, we record $\mathbf{p}^{(i)}$ every epoch and choose $\Delta\tau = 1/t$, namely, averaging the predicted probabilities during t epochs [47, 46]. The TD $\bar{\mathbf{p}}^{(t)}$ takes all the predicted probabilities during model optimization into account. Hence, it encapsulates the overall tendency of the model during t epochs of optimization, avoiding being solely biased towards the snapshot of $\mathbf{p}^{(t)}$ in the final epoch t .

3. Is TD Useful for Quantifying Uncertainty?

In this section, we provide empirical and theoretical evidence to support our argument: *TD is more effective in separating uncertain data from certain data than the model snapshot*, where the latter is often utilized to quantify data uncertainty in previous works [52, 45].

3.1. Motivating Observation

Settings. We aim to observe and compare the behavior of TD and the model snapshot for different sample difficulties. However, it is nontrivial to directly measure sample-wise difficulty, inhibiting the quantitative analysis of data uncertainty. To avoid this, we borrow the theoretical and empirical results of long-tailed visual recognition [33, 8, 19]: it is hard for the deep neural network-based model to train with fewer samples. Hence, we regard major and minor class

samples to contain many certain and uncertain samples for the model, respectively. We train the target classifier f on the long-tailed dataset during T epochs to obtain the TD and the model snapshot. We apply both approaches to the common estimators, entropy and margin. We denote entropy and margin scores from the model snapshot as H and M . In opposition, we denote the TD-driven scores as \bar{H} and \bar{M} . More details and discussions are described in Appendix B.

Results. Figure 2 shows the distribution the scores calculated with TD (x -axis) and model snapshot (y -axis). We can observe that scores from TD (\bar{H}, \bar{M}) successfully separate the major and the minor class samples, whereas scores from the model snapshot (H, M) fail to do so. We conclude that compared to model snapshots, TD is more helpful in separating uncertain samples from certain samples.

3.2. Theoretical Evidence

Theorem 1. (Informal) *Under the LE-SDE framework [54], with the assumption of local elasticity [17], certain samples and uncertain samples reveal different TD; especially, certain samples converge quickly than uncertain samples.*

The above theorem discusses different model behaviors depending on the difficulty of the sample. Compared to the uncertain sample, the certain sample has the same class samples nearby, which is the fundamental idea of level set estimation [22] and nearest neighbor [36] literature. We suspect that, due to the local elasticity of deep nets, samples close by have a bigger impact on the certain sample, hence changing its predicted probability more rapidly. As the certain sample is quicker to converge, its TD is larger than that of the uncertain sample. Intuitively, slower to train, struggling the classifier is to learn, hence TD capturing the uncertainty in the classifier’s perspective.

Theorem 2. (Informal) *Estimators such as Entropy (Equation 10) and Margin (Equation 11) successfully capture the difference of TD between easy and hard samples even for the case where it cannot be distinguished via the predicted probabilities of the model snapshot.*

The above theorem discusses the validity of entropy and margin on whether they can successfully differentiate between two samples of different TD but with the same final prediction. With Theorem 1, one can conclude that the common estimators’ scores calculated with TD are effective in capturing the data uncertainty. Due to the space constraints, we provide the details of the above results in Appendix A.

4. Utilizing TD for Active Learning

As tracking the TD of all the unlabeled data is computationally infeasible, we devise an efficient method to estimate the TD of unlabeled samples. We train the module

that directly predicts the TD of each sample by feeding the training samples, where its TD are freely available during training. Then, based on the predicted TD of each unlabeled sample, we use the common estimators, entropy or margin, to determine which sample is the most uncertain so that human annotators can label it. Hence, in this section, we describe the details of the module that estimates TD (§4.1) and how to train the module (§4.2). Finally, calculating the uncertainties using the module predictions for active learning is illustrated (§4.3).

4.1. Training Dynamics Prediction Module

As mentioned, it is not computationally feasible to track TD for the large-scale unlabeled data as it requires model inference on all the unlabeled data every training epoch. Thus, we propose the TD prediction module m to efficiently predict the TD of unlabeled data at the t -th epoch. Being influenced by the previous studies [11, 52, 45, 25] that use additional modules to predict useful values such as loss or confidence by the target model outputs, multi-scale feature maps are aggregated and passed into our TD prediction module. The module produces the C -dimensional predictions:

$$\tilde{\mathbf{p}}_m^{(t)} = [\tilde{p}_m^{(t)}(1|x), \dots, \tilde{p}_m^{(t)}(C|x)]^T \in [0, 1]^C \quad (6)$$

estimating the actual TD $\bar{\mathbf{p}}^{(t)}$ of the given sample x in Equation 5. TD prediction module is jointly trained with the target classifier using a handful of parameters, having a negligible computational cost during training. The detailed architecture of the module is described in Appendix C.

Even though the architecture is similar to previous works [52, 45, 25], we observed that ours were much more stable during optimization and easier to train. We suspect that it is due to the target task difference; previous works trained the module that outputs only a single value via regression, whereas our module outputs C -dimensional probability distribution, which is similar to the main task of classifying images.

4.2. Training Objectives

To train the target classifier f at the t -th epoch, we use the cross-entropy loss function $\mathcal{L}_{\text{target}}$ on the predicted probability $\mathbf{p}^{(t)}$ and a one-hot encoded vector $\mathbf{y} \in \{0, 1\}^C$ of the true label y :

$$\mathcal{L}_{\text{target}} = \mathcal{L}_{\text{CE}}(\mathbf{p}^{(t)}, \mathbf{y}) = -\log p^{(t)}(y|x). \quad (7)$$

Meanwhile, the prediction module m learns the TD of a sample x by minimizing the Kullback–Leibler (KL) divergence between the predicted TD $\tilde{\mathbf{p}}_m^{(t)}$ and the actual TD $\bar{\mathbf{p}}^{(t)}$:

$$\begin{aligned} \mathcal{L}_{\text{module}} &= \mathcal{L}_{\text{KL}}(\bar{\mathbf{p}}^{(t)} || \tilde{\mathbf{p}}_m^{(t)}) \\ &= \sum_{c=1}^C \bar{p}^{(t)}(c|x) \log \left(\frac{\bar{p}^{(t)}(c|x)}{\tilde{p}_m^{(t)}(c|x)} \right). \end{aligned} \quad (8)$$

The final objective function of our proposed method is defined as follows:

$$\mathcal{L} = \mathcal{L}_{\text{target}} + \lambda \mathcal{L}_{\text{module}} \quad (9)$$

where λ is a balancing factor to control the effect of $\mathcal{L}_{\text{module}}$ during model training.

4.3. Quantifying Uncertainty with TD

We argue that uncertain samples can be effectively distinguished from unlabeled data using the predicted TD. To verify the effectiveness of leveraging TD, we feed the predicted TD to entropy and margin (§2) by replacing snapshot probability \mathbf{p} with the predicted TD $\bar{\mathbf{p}}$. We choose these estimators as they are widely used for quantifying uncertainty. We feed $\bar{\mathbf{p}}$, replacing \mathbf{p} , to the entropy \bar{H} :

$$\bar{H}(\bar{\mathbf{p}}) = -\sum_{c=1}^C \bar{p}(c|x) \log \bar{p}(c|x). \quad (10)$$

Entropy \bar{H} is maximized when $\bar{\mathbf{p}}$ is uniform, *i.e.*, the sample is uncertain for the target classifier. Margin \bar{M} is also similarly employed:

$$\bar{M}(\bar{\mathbf{p}}) = \bar{p}(\hat{y}|x) - \max_{c \neq \hat{y}} \bar{p}(c|x). \quad (11)$$

Since we do not have true labels of unlabeled samples, we use the predicted labels \hat{y} of the target classifier instead of the true labels. There are several possible variants of \bar{M} depending on the definition of \hat{y} . We conduct experiments to compare \bar{M} with its variants. The experimental details and results are in Appendix D.4.

At the data selection phase, we use the predicted TD $\tilde{\mathbf{p}}_m^{(T)}$ instead of the actual TD $\bar{\mathbf{p}}^{(T)}$ as in Equation 10 & 11 to estimate the TD-driven uncertainties of the unlabeled sample x at the final epoch T . By using the estimated uncertainty with the predicted TD, we select the most informative samples for model training.

5. Experiments

In this section, we experimentally verify the effectiveness of our method, TiDAL, which utilizes the estimated training dynamics from the prediction module to discern uncertain samples from unlabeled data. We describe the detailed settings and the baseline methods for our experiments (§5.1) and show the results on both balanced (§5.2) and imbalanced datasets (§5.3). We further analyze whether the TD prediction module is effective for AL performance and can successfully estimate the TD (§5.4). We end the section by discussing the potential limitations of our method (§5.5).

5.1. Experimental Setup

Datasets. To assess the performance of our proposed method and baseline methods, we conduct experiments on

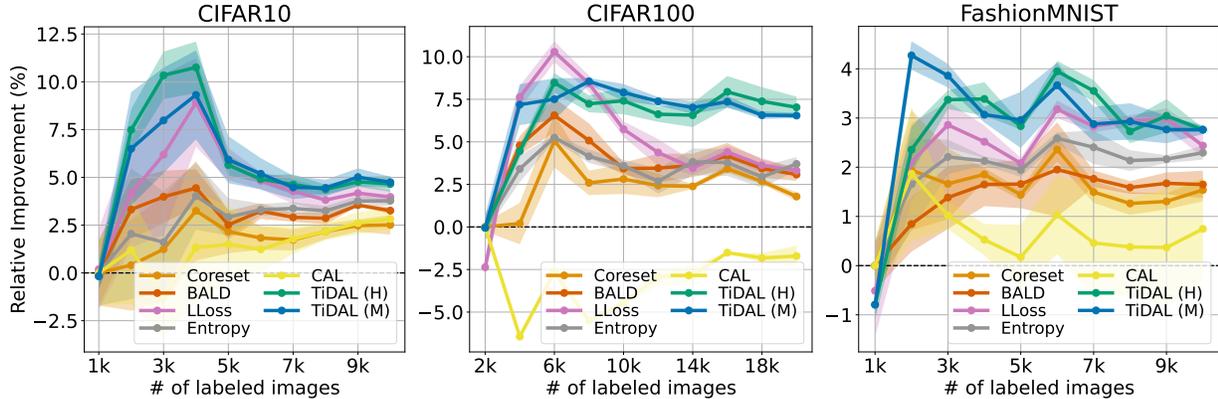


Figure 3: Averaged relative accuracy improvement curves and their 95% confidence interval (shaded) of AL methods over the number of labeled samples on balanced datasets. TiDAL (\bar{H}) and TiDAL (\bar{M}) denote the performance of TiDAL when with entropy \bar{H} and margin \bar{M} as the data uncertainty estimation strategy, respectively.

the following five datasets: CIFAR10/100 [27], FashionMNIST [51], SVHN [34], and iNaturalist2018 [50]. Since CIFAR and FashionMNIST are both balanced, we further modify them to simulate the data imbalance in the real world, following the previous long-tail visual recognition studies [8, 33, 56, 19]. The imbalance ratio is defined as N_{\max}/N_{\min} where N is the number of samples in each class. We make two variants with data imbalance ratios 10 and 100 for each dataset. Unlike the above, SVHN and iNaturalist18 are already imbalanced. Especially, iNaturalist2018 is commonly chosen to demonstrate how methods work in imbalanced real-world settings. The dataset statistics are summarized in Appendix D.

Baselines. For a fair comparison, we compare our TiDAL with the following baselines which train a target classifier with only labeled data. **Random sampling:** a simple baseline that randomly selects data samples from the unlabeled dataset. **Entropy sampling [43]:** an uncertainty-based method that selects data samples based on the maximum entropy. **BALD [13]:** an uncertainty-based method that selects data samples based on the mutual information between the model prediction and the posterior. **Core-Set [42]:** a diversity-based method that selects representative data samples covering all data through a minimum radius. **LLoss [52]:** an uncertainty-based method that learns to estimate the errors of the predictions (loss) made by the learner and select data samples based on the predicted loss. **CAL [55]:** recent work on using TD, gathering sample-wise TD information on whether the classifier was consistently correct or not during training. CAL splits the samples into two classes by applying a heuristic threshold to the TD information to train a binary classifier that outputs uncertainty score. To verify the effectiveness of TiDAL, we further compare it with the two semi-supervised AL methods, **VAAL [45]** and **TA-VAAL [25]** in Appendix D.3. Note that

these methods further utilize unlabeled data for training the selection module, thus it is unfair for our TiDAL.

Active learning setting. We follow the same setting from [6, 52] for the detailed AL settings. For the initial step, we randomly select initial samples to be annotated from the unlabeled dataset, where we use them to train the initial target classifier. Then, we obtain a random subset from the unlabeled data pool \mathcal{D}_u to choose the top- k samples based on the criterion of each method, where those samples will be annotated. We repeat the above cycle, training a classifier from scratch from the continuously expanding labeled set.

Implementation details. For a fair comparison, we use the same backbone network ResNet-18 [18] except for iNaturalist2018, where we use ResNet-50 [18] pretrained on ImageNet [12]. All models are trained with SGD optimizer with momentum 0.9, weight decay $5 \cdot 10^{-4}$, and learning rate (LR) decay of 0.1. For CIFAR10/100 and SVHN, we train the model for 200 epochs with an initial LR of 0.1 and decay at epoch 160. For FashionMNIST, 100 epochs with an initial LR of 0.1 and decay at epoch 80. For iNaturalist2018, 50 epochs with an initial LR of 0.01 and decay at epoch 40. For CIFAR10/100, SVHN and FashionMNIST, we set the batch size and the unlabeled subset size to be 128 and 10^4 , respectively. For iNaturalist2018, which is much larger than other datasets, we set the batch size and the unlabeled subset size to 256 and 10^6 , respectively. We set the balancing factor to 1.0.

Evaluation details. To compare with other state-of-the-art baselines, we show the average accuracy and 95% confidence interval with three trials. We mainly compare the model performances with relative accuracy improvement to random sampling, demonstrating how much it improves

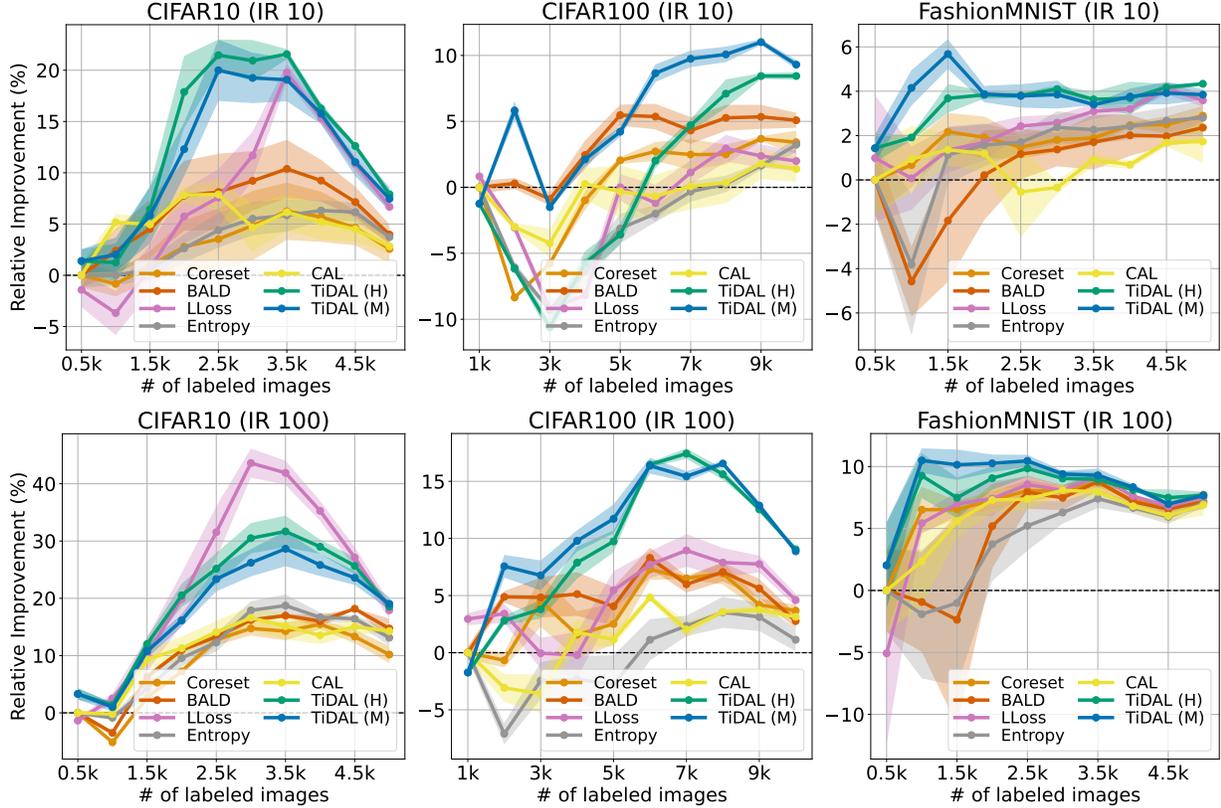


Figure 4: Averaged relative accuracy improvement curves and their 95% confidence interval (shaded) of AL methods over the number of labeled samples on synthetically imbalanced datasets. We use the imbalance ratio (IR) of 10 and 100 on CIFAR10, CIFAR100, and FashionMNIST.

upon the naive approach on each cycle. Additionally, absolute accuracy is also plotted in Appendix D.

5.2. Results on Balanced Datasets

Figure 3 and 10 compare our TiDAL against the state-of-the-art methods on various balanced datasets: CIFAR10, CIFAR100, and FashionMNIST. For all the datasets, the two variants of TiDAL outperform all the baselines at all AL cycles except for LLoss, which shows better improvement than TiDAL (\bar{M}) on CIFAR10 with an imbalance ratio of 100. Nonetheless, our TiDAL achieves the best final performance compared to all the baselines. CAL, which uses training dynamics, generally underperforms compared to others. We suspect that CAL is sensitive to its threshold hyperparameter.

5.3. Results on Imbalanced Datasets

Synthetically imbalanced datasets. Similar to the above, Figure 4, 9, and 11 shows the performance improvements on the synthetically imbalanced datasets with the two imbalance ratios, 10 and 100. Except for the CIFAR10 with an imbalance ratio of 100, our methods show superb performance across all the imbalanced settings. TiDAL per-

forms especially well with a small variance in imbalanced CIFAR100, where the number of classes is the largest. In imbalanced FashionMNIST, the performance quickly rises to 2.5k labeled images and then saturates. This implies that FashionMNIST is easier than other datasets, and needs to focus more on the early training steps to compare with other models. TiDAL also shows overall better performance on FashionMNIST, especially in the early steps.

Real-world imbalanced datasets. Figure 5 and 10 shows evaluation results on real-world imbalanced datasets. For iNaturalist2018, which is the large-scale long-tailed classification dataset, TiDAL shows outstanding performance compared to other methods. For SVHN, TiDAL shows the best improvements with low variance as the number of labeled images increases except for the initial stage. LLoss shows outstanding performance only in the initial stage, where we presume that the loss prediction module of LLoss acts as a regularizer during model optimization.

5.4. Analysis on the TD Prediction Module

Effectiveness of the TD prediction module. In order to verify the efficacy of using the predicted TD \tilde{p}_m , we con-

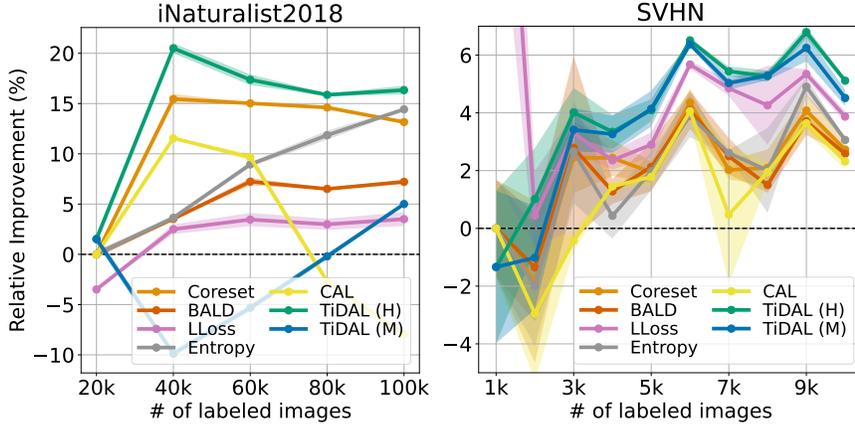
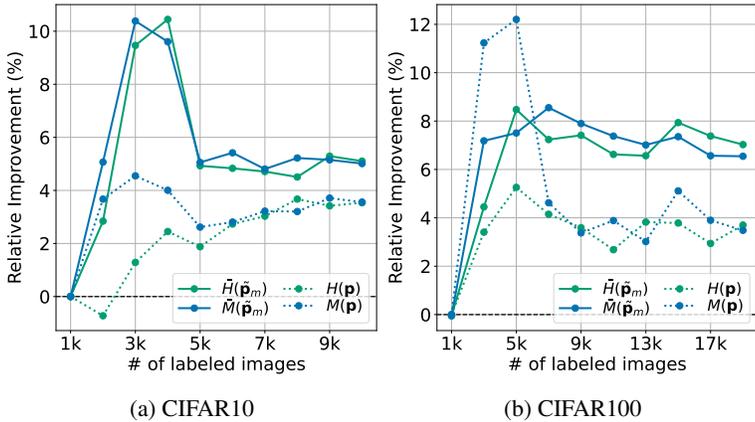


Figure 5: Averaged relative accuracy improvement curves and its 95% confidence interval (shaded) of AL methods over the number of labeled samples on real-world imbalanced datasets: iNaturalist2018 and SVHN. For SVHN, LLoss shows a substantial improvement of $20.02\% \pm 6.77\%$ at the initial phase (1k), but we clip the plot to show the performance afterward more clearly.



(a) CIFAR10

(b) CIFAR100

Figure 6: Ablation test results. $\bar{H}(\tilde{\mathbf{p}}_m)$ and $\bar{M}(\tilde{\mathbf{p}}_m)$ use the predicted TD $\tilde{\mathbf{p}}_m$ of the prediction module m . In contrast, $H(\mathbf{p})$ and $M(\mathbf{p})$ use the predicted probability of the model snapshot \mathbf{p} . TD shows better performance than the model snapshot, implying that TD is better at quantifying data uncertainty.

duct an ablation test that compares the performance between when using and not using the TD prediction module m . Figure 6 shows the results on balanced CIFAR10/100. We observe that $\bar{H}(\tilde{\mathbf{p}}_m)$ and $\bar{M}(\tilde{\mathbf{p}}_m)$ using the predicted TD $\tilde{\mathbf{p}}_m$ to estimate the data uncertainty significantly outperform the methods $H(\mathbf{p})$ and $M(\mathbf{p})$ that use only the final predicted probabilities \mathbf{p} of the target classifier f , showing better performance in the whole training cycle. Even $M(\mathbf{p})$ shows temporary improvement in earlier steps on CIFAR100, $\bar{H}(\tilde{\mathbf{p}}_m)$ and $\bar{M}(\tilde{\mathbf{p}}_m)$ maintain stable improvement, eventually winning over $M(\mathbf{p})$. This indicates that the predicted TD $\tilde{\mathbf{p}}_m$ of the TD prediction module m produces better data uncertainty estimation than the predicted

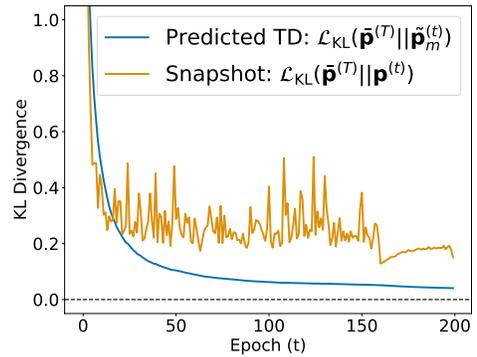


Figure 7: KL divergence scores of the actual TD $\tilde{\mathbf{p}}^{(T)}$ with the predicted TD $\tilde{\mathbf{p}}_m^{(t)}$ and the predicted probability of the model snapshot $\mathbf{p}^{(t)}$, respectively, during model optimization. Our predicted TD can accurately approximate the actual TD.

probability \mathbf{p} of the target classifier f .

Predictive performance of the TD prediction module.

We verify whether the TD prediction module m accurately predicts the actual TD $\tilde{\mathbf{p}}$. Its prediction performance is crucial as we use the predicted TD $\tilde{\mathbf{p}}_m$ of the module m to quantify uncertainties of unlabeled data. Using the KL divergence \mathcal{L}_{KL} , we analyze that the predicted TD $\tilde{\mathbf{p}}_m$ converges to the actual TD $\tilde{\mathbf{p}}$ at the data selection phase. We calculate $\mathcal{L}_{\text{KL}}(\tilde{\mathbf{p}}^{(T)} || \tilde{\mathbf{p}}_m^{(t)})$ and compare it with $\mathcal{L}_{\text{KL}}(\tilde{\mathbf{p}}^{(T)} || \mathbf{p}^{(t)})$ which is set as a baseline computed with the actual TD $\tilde{\mathbf{p}}$ and the predicted probabilities \mathbf{p} (snapshot) of the target classifier f . In this analysis, we use the bal-

anced CIFAR10 where the sample-wise averaged KL divergence scores are computed on the test set. Figure 7 shows that the final predicted TD successfully approximates the actual TD, while the predicted probability is highly different from the actual TD. We conclude that the TD prediction module m can produce the TD efficiently, leading to performance improvement, and the predicted TD acts as a better approximation of the actual TD than the predicted probability of a model snapshot captured at each epoch.

5.5. Limitations

We found two potential limitations of our TiDAL derived from the fact that it relies on the outputs of the target classifier to compute the TD. First, TiDAL is designed only for classification tasks, and thus it cannot be applied to AL targeting other tasks, such as regression [10, 15]. Second, TiDAL is highly influenced by the performance of the target classifier, especially when the target classifier wrongly classifies the hard negative samples with a high confidence during model optimization. These samples can be treated as certain samples (i.e. will not be selected for annotation) because they have low estimated uncertainties from the predicted TD, even though the target classifier fails to predict the true label of the samples correctly. As a future work, we will study extending our TiDAL in the task-agnostic ways with a safeguard combating the wrongly classified samples.

6. Related Work

6.1. Active Learning

AL methods target to construct a dataset with the most useful samples based on the assumption that each sample has different importance in model training [40]. Two mainstream AL approaches exist for efficiently querying the unlabeled data: pool-based methods [31, 52, 45] use various ways to extract samples from an unlabeled data pool effectively, and synthesis-based methods [1, 58, 49] generate informative samples for the model. Pool-based methods can be roughly divided based on query strategies: uncertainty-based [13, 52, 45, 20] and diversity-based [42, 14, 38] methods, where some methods use the hybrid of both [4, 44, 25]. Uncertainty-based methods focus on finding which samples would be the most uncertain for the model, whereas diversity-based methods aim to construct a subset of representative samples of the input distribution. Our proposed method, TiDAL, lies in uncertainty-based methods. The significant difference between TiDAL and previous uncertainty-based methods is that TiDAL estimates data uncertainty using TD that contains additional hints generated during model training. In contrast, the previous methods leverage only static information (e.g., loss [52, 20] and predicted probabilities [13, 45, 25]) obtained by a model snapshot at the data selection phase.

6.2. Training Dynamics

TD focuses on how deep neural networks are optimized under back-propagation-based stepwise weight updates. Many studies try to understand how gradient descent can effectively obtain the global minimum by analyzing the loss landscape of neural networks [24, 32] or its loss trajectory [3]. Some also import alternative models that are more mathematically approachable to analyze, such as neural tangent kernels [21], deep Gaussian processes [30], or stochastic differential equations [54]. On the other hand, the phenomenological and practical viewpoint of TD also exists. [48] coin the term Forgetting Dynamics to assert that unforgettable samples are often less helpful, and [9] show that the model could prefer samples that are often wrongly predicted throughout model training. TD is also commonly used in noisy label literature to find potential noisy labels as they tend to fit later on model training [2, 39] or locate samples that can be relabeled correctly [46]. Furthermore, [57] calculate the Dynamic Instance Hardness score by monitoring losses of each sample or whether the prediction gets flipped so that higher scored samples can be prioritized for curriculum learning, and [23] feed the loss history to the auxiliary neural network to mediate the curriculum for training. [29] also introduce temporal ensembling for semi-supervised learning, where the model fits towards averaged probability outputs. [47, 37] devise Data Maps to inspect datasets with two TD measures; confidence and its variability across epochs on the true class prediction. [55] further extend the Data Maps for AL, whether the target classifier was consistently correct or not during training. The proposed method splits the labeled samples by applying a heuristic threshold on the level of consistency to train a binary classifier that is trained to discern uncertain samples. Even though the work, similar to ours, also utilizes TD, it relies on empirical observations and heuristic choices to separate the certain and uncertain samples. In this study, we link the concept of TD to AL with both empirical and theoretical results to estimate the uncertainty of unlabeled samples, which is often neglected in previous TD studies.

7. Conclusion

We propose a novel active learning method, Training Dynamics for Active Learning (TiDAL), by linking the concept of training dynamics to active learning. We provide motivating observations and theoretical evidence for using training dynamics to estimate the uncertainty of unlabeled data. Since tracking the training dynamics of large-scale unlabeled data is infeasible, TiDAL utilizes a training dynamics prediction module to efficiently predict the training dynamics of the unlabeled data. Based on the predicted training dynamics, TiDAL quantifies data uncertainty using the common uncertainty estimators: entropy and mar-

gin. Extensive experiments on multiple benchmark datasets demonstrate the effectiveness of our method, surpassing the existing state-of-the-art active learning methods. We further analyze that using our training dynamics prediction module is effective and the module successfully predicts the TD of unlabeled data.

References

- [1] Dana Angluin. Queries and concept learning. *Machine learning*, 2(4):319–342, 1988. [8](#)
- [2] Eric Arazo, Diego Ortego, Paul Albert, Noel O’Connor, and Kevin McGuinness. Unsupervised label noise modeling and loss correction. In *International conference on machine learning*, pages 312–321. PMLR, 2019. [8](#)
- [3] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. In *International Conference on Learning Representations*, 2018. [8](#)
- [4] Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep batch active learning by diverse, uncertain gradient lower bounds. *arXiv preprint arXiv:1906.03671*, 2019. [8](#)
- [5] Les E Atlas, David A Cohn, and Richard E Ladner. Training connectionist networks with queries and selective sampling. In *Advances in neural information processing systems*, pages 566–573. Citeseer, 1990. [1](#)
- [6] William H Beluch, Tim Genewein, Andreas Nürnberger, and Jan M Köhler. The power of ensembles for active learning in image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9368–9377, 2018. [1](#), [5](#)
- [7] Javad Zolfaghari Bengar, Joost van de Weijer, Laura Lopez Fuentes, and Bogdan Raducanu. Class-balanced active learning for image classification. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1536–1545, 2022. [15](#)
- [8] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. *Advances in Neural Information Processing Systems*, 32:1567–1578, 2019. [3](#), [5](#), [15](#)
- [9] Haw-Shiuan Chang, Erik Learned-Miller, and Andrew McCallum. Active bias: Training more accurate neural networks by emphasizing high variance samples. *Advances in Neural Information Processing Systems*, 30, 2017. [1](#), [8](#)
- [10] David Cohn, Les Atlas, and Richard Ladner. Improving generalization with active learning. *Machine learning*, 15(2):201–221, 1994. [8](#)
- [11] Charles Corbière, Nicolas Thome, Avner Bar-Hen, Matthieu Cord, and Patrick Pérez. Addressing failure prediction by learning model confidence. *Advances in Neural Information Processing Systems*, 32, 2019. [4](#)
- [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. [5](#)
- [13] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *International Conference on Machine Learning*, pages 1183–1192. PMLR, 2017. [1](#), [5](#), [8](#)
- [14] Daniel Gissin and Shai Shalev-Shwartz. Discriminative active learning. *arXiv preprint arXiv:1907.06347*, 2019. [1](#), [8](#)
- [15] Jia Gong, Zhipeng Fan, QiuHong Ke, Hossein Rahmani, and Jun Liu. Meta agent teaming active learning for pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11079–11089, 2022. [8](#)
- [16] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017. [15](#)
- [17] Hangfeng He and Weijie Su. The local elasticity of neural networks. In *International Conference on Learning Representations*, 2019. [3](#), [12](#)
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [5](#), [15](#)
- [19] Youngkyu Hong, Seungju Han, Kwanghee Choi, Seokjun Seo, Beomsu Kim, and Buru Chang. Disentangling label distribution for long-tailed visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6626–6636, 2021. [3](#), [5](#), [15](#)
- [20] Siyu Huang, Tianyang Wang, Haoyi Xiong, Jun Huan, and Dejing Dou. Semi-supervised active learning with temporal output discrepancy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3447–3456, 2021. [8](#)
- [21] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *Advances in neural information processing systems*, 31, 2018. [8](#)
- [22] Heinrich Jiang, Been Kim, Melody Guan, and Maya Gupta. To trust or not to trust a classifier. *Advances in neural information processing systems*, 31, 2018. [3](#), [12](#)
- [23] Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *International Conference on Machine Learning*, pages 2304–2313. PMLR, 2018. [8](#)
- [24] Kenji Kawaguchi. Deep learning without poor local minima. *Advances in neural information processing systems*, 29, 2016. [8](#)
- [25] Kwanyoung Kim, Dongwon Park, Kwang In Kim, and Se Young Chun. Task-aware variational adversarial active learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8166–8175, 2021. [4](#), [5](#), [8](#), [16](#)
- [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. [15](#)
- [27] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009. [5](#)

- [28] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. The cifar-10 dataset. *online: <http://www.cs.toronto.edu/kriz/cifar.html>*, 55(5), 2014. 15
- [29] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. In *International Conference on Learning Representations (ICLR)*. OpenReview.net, 2017. 1, 8
- [30] Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as gaussian processes. In *International Conference on Learning Representations*, 2018. 8
- [31] David D Lewis and William A Gale. A sequential algorithm for training text classifiers. In *Proceedings of the 17th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 3–12, 1994. 1, 8
- [32] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31, 2018. 8
- [33] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu. Large-scale long-tailed recognition in an open world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2537–2546, 2019. 3, 5, 15
- [34] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisaccho, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011. 5
- [35] Curtis Northcutt, Lu Jiang, and Isaac Chuang. Confident learning: Estimating uncertainty in dataset labels. *Journal of Artificial Intelligence Research*, 70:1373–1411, 2021. 15
- [36] Nicolas Papernot and Patrick McDaniel. Deep k-nearest neighbors: Towards confident, interpretable and robust deep learning. *arXiv preprint arXiv:1803.04765*, 2018. 3, 12
- [37] Seo Yeon Park and Cornelia Caragea. A data cartography based mixup for pre-trained language models. *arXiv preprint arXiv:2205.03403*, 2022. 8
- [38] Amin Parvaneh, Ehsan Abbasnejad, Damien Teney, Gholamreza Reza Haffari, Anton van den Hengel, and Javen Qinfeng Shi. Active learning by feature mixing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12237–12246, 2022. 8
- [39] Geoff Pleiss, Tianyi Zhang, Ethan Elenberg, and Kilian Q Weinberger. Identifying mislabeled data using the area under the margin ranking. *Advances in Neural Information Processing Systems*, 33:17044–17056, 2020. 2, 8, 15
- [40] Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang Chen, and Xin Wang. A survey of deep active learning. *ACM Computing Surveys (CSUR)*, 54(9):1–40, 2021. 8
- [41] Dan Roth and Kevin Small. Margin-based active learning for structured output spaces. In *European Conference on Machine Learning*, pages 413–424. Springer, 2006. 2
- [42] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In *International Conference on Learning Representations*, 2018. 1, 5, 8
- [43] Claude Elwood Shannon. A mathematical theory of communication. *The Bell system technical journal*, 27(3):379–423, 1948. 2, 5
- [44] Changjian Shui, Fan Zhou, Christian Gagné, and Boyu Wang. Deep active learning: Unified and principled method for query and training. In *International Conference on Artificial Intelligence and Statistics*, pages 1308–1318. PMLR, 2020. 8
- [45] Samarth Sinha, Sayna Ebrahimi, and Trevor Darrell. Variational adversarial active learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5972–5981, 2019. 2, 3, 4, 5, 8, 16
- [46] Hwanjun Song, Minseok Kim, and Jae-Gil Lee. Selfie: Refurbishing unclean samples for robust deep learning. In *International Conference on Machine Learning*, pages 5907–5915. PMLR, 2019. 3, 8
- [47] Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A Smith, and Yejin Choi. Dataset cartography: Mapping and diagnosing datasets with training dynamics. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9275–9293, 2020. 1, 2, 3, 8
- [48] Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. In *International Conference on Learning Representations*, 2018. 1, 8, 15
- [49] Toan Tran, Thanh-Toan Do, Ian Reid, and Gustavo Carneiro. Bayesian generative active deep learning. In *International Conference on Machine Learning*, pages 6295–6304. PMLR, 2019. 8
- [50] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8769–8778, 2018. 5
- [51] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017. 5
- [52] Donggeun Yoo and In So Kweon. Learning loss for active learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 93–102, 2019. 1, 2, 3, 4, 5, 8, 15, 16
- [53] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021. 12
- [54] Jiayao Zhang, Hua Wang, and Weijie Su. Imitating deep learning dynamics via locally elastic stochastic differential equations. *Advances in Neural Information Processing Systems*, 34, 2021. 3, 8, 12
- [55] Mike Zhang and Barbara Plank. Cartography active learning. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 395–406, 2021. 5, 8

- [56] Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhao-Min Chen. Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9719–9728, 2020. [5](#), [15](#)
- [57] Tianyi Zhou, Shengjie Wang, and Jeffrey Bilmes. Curriculum learning by dynamic instance hardness. *Advances in Neural Information Processing Systems*, 33:8602–8613, 2020. [8](#)
- [58] Jia-Jie Zhu and José Bento. Generative adversarial active learning. *arXiv preprint arXiv:1702.07956*, 2017. [8](#)

Appendix

A. Details on the Theoretical Evidence

A.1. Proof of Theorem 1

We adopt the settings of [54], with a slight modification of the assumption that sample-level local elasticity affects the training dynamics instead of class-level local elasticity.

Consider the binary classification problem with two classes $k = 1, 2$ where class 1 consists of both certain (easy) and uncertain (hard) samples and class 2 only consists of samples with the same certainty (easiness). Let $\mathcal{S}_{1,e}$, $\mathcal{S}_{1,h}$ and \mathcal{S}_2 denote the easy samples from class 1, hard samples from class 1, and samples from class 2 respectively, which constitutes the partition of the whole set of training samples \mathcal{S} : $\mathcal{S} = \mathcal{S}_{1,e} \cup \mathcal{S}_{1,h} \cup \mathcal{S}_2$. Let the corresponding sample sizes be $n_{1,e} = |\mathcal{S}_{1,e}|$, $n_{1,h} = |\mathcal{S}_{1,h}|$, $n_2 = |\mathcal{S}_2|$ and $n = |\mathcal{S}| = n_{1,e} + n_{1,h} + n_2$, respectively.

At each iteration m , a training candidate sample $J_m \in \mathcal{S}$ with class L_m is sampled uniformly from the whole training set \mathcal{S} with replacement. Training using this sample J_m via SGD affects the training dynamics of other samples $s \in \mathcal{S}$ of class k as:

$$X_s^k(m) = X_s^k(m-1) + hE_{s,J_m}X_{J_m}^{L_m}(m-1) + \sqrt{h}\zeta_s^k(m-1), \quad (12)$$

where $X > 0$ is logit of the true label, $h > 0$ is the step size, $\zeta \sim \mathcal{N}(0, \sigma^2)$ denotes the noise term arises during training, and $E \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ refers to the sample-level local elasticity [17] where each entry $E_{s,s'}$ measures the strength of the local elasticity of s' by s . For simplicity, we assume this local elasticity does not depend on the time step m . Furthermore, we consider that the sample-level local elasticity only depends on the set $\mathcal{S}_{1,e}$, $\mathcal{S}_{1,h}$ and \mathcal{S}_2 in which each samples are in.

Let

$$\bar{X}^{1,e}(t) = \frac{1}{n_{1,e}} \sum_{s \in \mathcal{S}_{1,e}} X_s^1(t), \bar{X}^{1,h}(t) = \frac{1}{n_{1,h}} \sum_{s \in \mathcal{S}_{1,h}} X_s^1(t), \bar{X}^2(t) = \frac{1}{n_2} \sum_{s \in \mathcal{S}_2} X_s^2(t) \quad (13)$$

be the averaged logits for certain samples in class 1, uncertain samples in class 1, and class 2 respectively.

Regarding the strength of local elasticity between ‘‘class’’ of samples, for some constants α_e , α_h and β , we set the value of $E_{s,s'}$ to model sample-level local elasticity for (1) between easy and hard samples in the class 1 and (2) between classes 1 and 2. We use the values $\alpha_e > \alpha_h > \beta > 0$ to define the easiness such that the power exerted by sample-level local elasticity between easy samples are stronger for the pair of easy samples than for the pair consists of one or more hard sample.

- $E_{s,s'} = \alpha_e$ if $(s, s') \in (\mathcal{S}_{1,e} \times \mathcal{S}_{1,e}) \cup (\mathcal{S}_2 \times \mathcal{S}_2)$,
- $E_{s,s'} = \alpha_h$ if $(s, s') \in (\mathcal{S}_{1,e} \times \mathcal{S}_{1,h}) \cup (\mathcal{S}_{1,h} \times \mathcal{S}_{1,e}) \cup (\mathcal{S}_{1,h} \times \mathcal{S}_{1,h})$ (either $s \in \mathcal{S}_{1,h}$ or $s' \in \mathcal{S}_{1,h}$),
- $E_{s,s'} = \beta$ otherwise.

Intuitively, one can interpret the above assumption as easy samples being clustered with each other [22, 36], hence having a stronger influence on each other due to the local elasticity. On the contrary, hard samples are often distant from other same-class samples. Their influence is often limited, as memorizing is easy for the neural nets due to their large capacity [53]. Finally, we ignore the influence of other class samples in this proof for simplicity, as we are only considering the logits of the true label.

Theorem 1. (Formal) *On average, the convergence speed of logit is faster for easy samples than hard samples. Formally:*

$$\frac{d\bar{X}^{1,e}(t)}{dt} > \frac{d\bar{X}^{1,h}(t)}{dt}. \quad (14)$$

Proof. Fix a target sample $s \in \mathcal{S}$, and execute the dynamics (12) r times since step m . Accumulated change for feature X becomes

$$X_s^k(m+r) - X_s^k(m) = h \sum_{q=1}^r E_{k,L_{m+q}} X_{J_{m+q}}^{L_{m+q}}(m+q-1) + \epsilon_{s,k,r,h}, \quad (15)$$

where $\epsilon = \sqrt{h} \sum_{q=1}^r \zeta_s^k(m+q-1) \sim \mathcal{N}(0, \sigma^2 r h)$ is the accumulated noise terms during r updates. Regarding terms inside the summation, we can divide cases based on which sample J_r (with corresponding class L_r) is actually selected as a training candidate at iteration $\nu (= m+q)$:

$$E_{k,J_\nu} X_{J_\nu}^{L_\nu}(\nu-1) = \mathbf{1}_{J_\nu \in \mathcal{S}_{1,e}} E_{k,J_\nu} X_{J_\nu}^1(\nu-1) + \mathbf{1}_{J_\nu \in \mathcal{S}_{1,h}} E_{k,J_\nu} X_{J_\nu}^1(\nu-1) + \mathbf{1}_{J_\nu \in \mathcal{S}_2} E_{k,J_\nu} X_{J_\nu}^2(\nu-1), \quad (16)$$

hence the summand from (15) becomes (omitting time index for X for simplicity)

$$h \sum_{q=1}^r \left(\mathbf{1}_{J_{m+q} \in \mathcal{S}_{1,e}} E_{k,J_{m+q}} X_{J_{m+q}}^1 + \mathbf{1}_{J_{m+q} \in \mathcal{S}_{1,h}} E_{k,J_{m+q}} X_{J_{m+q}}^1 + \mathbf{1}_{J_{m+q} \in \mathcal{S}_2} E_{k,J_{m+q}} X_{J_{m+q}}^2 \right),$$

and for sufficiently large r we can approximate the summations as the sample-average dynamics:

$$\begin{aligned} & h \sum_{q=1}^r \left(\mathbf{1}_{J_{m+q} \in \mathcal{S}_{1,e}} E_{k,J_{m+q}} X_{J_{m+q}}^1 + \mathbf{1}_{J_{m+q} \in \mathcal{S}_{1,h}} E_{k,J_{m+q}} X_{J_{m+q}}^1 + \mathbf{1}_{J_{m+q} \in \mathcal{S}_2} E_{k,J_{m+q}} X_{J_{m+q}}^2 \right) \\ & \approx hr \left(\mathbb{P}(J \in \mathcal{S}_{1,e}) \frac{\sum_{s \in \mathcal{S}_{1,e}} E_{k,s} X_s^1}{n_{1,e}} + \mathbb{P}(J \in \mathcal{S}_{1,h}) \frac{\sum_{s \in \mathcal{S}_{1,h}} E_{k,s} X_s^1}{n_{1,h}} + \mathbb{P}(J \in \mathcal{S}_2) \frac{\sum_{s \in \mathcal{S}_2} E_{k,s} X_s^2}{n_2} \right) \\ & \approx hr \left(\frac{n_{1,e}}{n} \frac{\sum_{s \in \mathcal{S}_{1,e}} E_{k,s} X_s^1}{n_{1,e}} + \frac{n_{1,h}}{n} \frac{\sum_{s \in \mathcal{S}_{1,h}} E_{k,s} X_s^1}{n_{1,h}} + \frac{n_2}{n} \frac{\sum_{s \in \mathcal{S}_2} E_{k,s} X_s^2}{n_2} \right) \end{aligned} \quad (17)$$

As the components of E only depend on the subset sample relies, we can rewrite accumulated dynamics of logits (15) for three cases separately, utilizing the notation of averaged logit (13):

$$\begin{aligned} X_s^{1,e}(m+r) - X_s^{1,e}(m) &= hr \left(\frac{n_{1,e}}{n} \alpha_e \bar{X}^{1,e}(m) + \frac{n_{1,h}}{n} \alpha_h \bar{X}^{1,h}(m) + \frac{n_2}{n} \beta \bar{X}^2(m) \right) + \epsilon_{s,k,r,h} \\ X_s^{1,h}(m+r) - X_s^{1,h}(m) &= hr \left(\frac{n_{1,e}}{n} \alpha_h \bar{X}^{1,e}(m) + \frac{n_{1,h}}{n} \alpha_h \bar{X}^{1,h}(m) + \frac{n_2}{n} \beta \bar{X}^2(m) \right) + \epsilon_{s,k,r,h} \\ X_s^2(m+r) - X_s^2(m) &= hr \left(\frac{n_{1,e}}{n} \beta \bar{X}^{1,e}(m) + \frac{n_{1,h}}{n} \beta \bar{X}^{1,h}(m) + \frac{n_2}{n} \alpha_e \bar{X}^2(m) \right) + \epsilon_{s,k,r,h}, \end{aligned} \quad (18)$$

with a little bit of abbreviated notation for class 1: $X_s^{1,e} = X_s^1$ for easy sample s , and similarly for hard samples. The differential counterpart of the above difference equation is

$$\begin{aligned} dX_s^{1,e}(t) &= \left(\frac{n_{1,e}}{n} \alpha_e \bar{X}^{1,e}(t) + \frac{n_{1,h}}{n} \alpha_h \bar{X}^{1,h}(t) + \frac{n_2}{n} \beta \bar{X}^2(t) \right) dt + \sigma dW^s(t) \\ dX_s^{1,h}(t) &= \left(\frac{n_{1,e}}{n} \alpha_h \bar{X}^{1,e}(t) + \frac{n_{1,h}}{n} \alpha_h \bar{X}^{1,h}(t) + \frac{n_2}{n} \beta \bar{X}^2(t) \right) dt + \sigma dW^s(t) \\ dX_s^2(t) &= \left(\frac{n_{1,e}}{n} \beta \bar{X}^{1,e}(t) + \frac{n_{1,h}}{n} \beta \bar{X}^{1,h}(t) + \frac{n_2}{n} \alpha_e \bar{X}^2(t) \right) dt + \sigma dW^s(t), \end{aligned} \quad (19)$$

where $W^s(t)$ is standard Wiener process per sample. Averaging each differential equation with respect to each set of samples and ignoring error terms yield a set of simultaneous deterministic differential equations for averaged logits:

$$\begin{aligned} d\bar{X}^{1,e}(t) &= \left(\frac{n_{1,e}}{n} \alpha_e \bar{X}^{1,e}(t) + \frac{n_{1,h}}{n} \alpha_h \bar{X}^{1,h}(t) + \frac{n_2}{n} \beta \bar{X}^2(t) \right) dt \\ d\bar{X}^{1,h}(t) &= \left(\frac{n_{1,e}}{n} \alpha_h \bar{X}^{1,e}(t) + \frac{n_{1,h}}{n} \alpha_h \bar{X}^{1,h}(t) + \frac{n_2}{n} \beta \bar{X}^2(t) \right) dt \\ d\bar{X}^2(t) &= \left(\frac{n_{1,e}}{n} \beta \bar{X}^{1,e}(t) + \frac{n_{1,h}}{n} \beta \bar{X}^{1,h}(t) + \frac{n_2}{n} \alpha_e \bar{X}^2(t) \right) dt, \end{aligned} \quad (20)$$

To compare the convergence speed of average logit between certain and uncertain samples in the same class 1, observe that

$$\frac{d\bar{X}^{1,e}(t)}{dt} - \frac{d\bar{X}^{1,h}(t)}{dt} = \frac{n_{1,e}}{n} (\alpha_e - \alpha_h) \bar{X}^{1,e}(t) > 0. \quad (21)$$

□

With additional assumptions on the other class logits being the same, one can also conclude that the estimated probability of the true label will increase steeply during training for the easy samples. After increasing to some extent, the probability will saturate to one; hence the snapshot model predictions will contain less useful information than monitoring its training dynamics. However, future work on extending the above theorem is needed. Starting from the basic idea above, that sample proximity and its amount influence the training dynamics, one can further relax the above assumptions, such as concentrating on the individuality of each sample or considering the changing elasticities during training. We hope our work ignites the theoretical research on uncertainty from the viewpoint of training dynamics.

A.2. Proof of Theorem 2

We aim to show the effectiveness of the proposed estimators, entropy (Equation 2) and margin (Equation 3), especially in the case where the probabilities converge. After training, it is commonly observed that the probabilities of the true label of all the samples tend to converge to one, whereas the speed of the convergence differs (Theorem 1). Hence, we show that the estimators can effectively discern the differences during training.

For each time step t during training, we have a sequence of predicted probabilities $p^{(t)}(y = c|x)$ corresponds to t , for each target class $c = 1, 2, \dots, C$. In our paper, we regard the area under the predicted probability $\bar{p}^{(T)}(y = c|x)$ of the sample x as the training dynamics (Equation 5), which is indeed a well-known metric of area under the curve, except that it is normalized properly to have value between 0 and 1. For convenience, let

$$\mathbf{s}(x) = \begin{bmatrix} s_1(x) \\ s_2(x) \\ \vdots \\ s_C(x) \end{bmatrix} = \begin{bmatrix} \bar{p}^{(T)}(y = 1|x) \\ \bar{p}^{(T)}(y = 2|x) \\ \vdots \\ \bar{p}^{(T)}(y = C|x) \end{bmatrix}$$

be the vector consisting the area under the prediction curve for each class up to final epoch T . By construction, the components in $\mathbf{s}(x)$ are nonnegative and sum to 1.

Theorem 2. (Formal) Assume that all target classes have the same area under the prediction curve except for the true class y . Suppose two training samples $(x_1, y_1), (x_2, y_2) \in \mathcal{D}$ satisfies

- a. $p^{(T)}(y_1|x_1) = p^{(T)}(y_2|x_2)$ (same predicted probability at the end of training)
- b. $\frac{1}{2} < s_{y_1}(x_1) < s_{y_2}(x_2)$ (but different TD, in terms of the area under the curve)

Then, the following inequalities hold:

1. $H(\mathbf{s}(x_1)) > H(\mathbf{s}(x_2))$;
2. $M(\mathbf{s}(x_1)) < M(\mathbf{s}(x_2))$.

Proof. By the assumption, for all target class c except the true class y , the area under the prediction curve is given by

$$s_c(x) = \frac{1 - s_y(x)}{C - 1}, \quad (22)$$

and the corresponding entropy can be calculated as

$$\begin{aligned} H(\mathbf{s}(x)) &= \sum_{c=1}^C (-s_c(x) \log(s_c(x))) \\ &= -s_y(x) \log s_y(x) - (C - 1) \cdot \left(\frac{1 - s_y(x)}{C - 1} \right) \log \left(\frac{1 - s_y(x)}{C - 1} \right) \\ &= -\{s_y(x) \log s_y(x) + (1 - s_y(x)) \log(1 - s_y(x))\} + (1 - s_y(x)) \log(C - 1) \\ &= H_2(s_y(x)) + (1 - s_y(x)) \log(C - 1). \end{aligned} \quad (23)$$

where $H_2(p) = -p \log p - (1 - p) \log(1 - p)$ stands for the binary entropy function. Since $H_2(p)$ is a decreasing function for $p > \frac{1}{2}$,

$$H(\mathbf{s}(x_1)) - H(\mathbf{s}(x_2)) = \{H_2(s_{y_1}(x_1)) - H_2(s_{y_2}(x_2))\} + \{s_{y_2}(x_2) - s_{y_1}(x_1)\} \log(C - 1) > 0,$$

which proves the first inequality stated.

The first assumption also gives the simplified formulation for the margin

$$M(\mathbf{s}(x)) = s_y(x) - \frac{1 - s_y(x)}{C - 1} = \frac{C}{C - 1} s_y(x) - \frac{1}{C - 1}, \quad (24)$$

in which the second inequality directly follows:

$$M(\mathbf{s}(x_1)) - M(\mathbf{s}(x_2)) = \frac{C}{C - 1} (s_{y_1}(x_1) - s_{y_2}(x_2)) < 0. \quad (25)$$

□

While the final predicted probabilities $p^{(T)}(y|x)$ of the training samples tend to converge to 1 for the true class y , otherwise 0, their TD (in this case $s(x) = \bar{p}^{(T)}$) may be different depending on the easiness of the samples. Thus, the degree of the easiness of the samples (i.e. uncertainty) could be captured from TD \bar{p} , whereas the predictions \mathbf{p} from a model snapshot cannot.

B. Details on the Motivating Observation

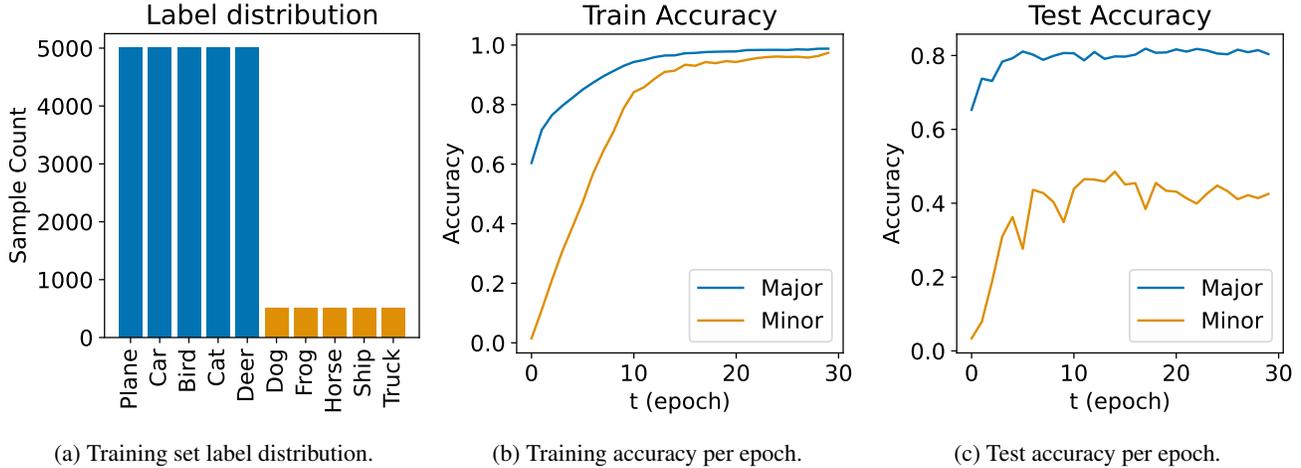


Figure 8: Training label distribution and accuracy curves for the motivating experiment in §3.1.

§3.1 empirically show that using TD is effective in separating uncertain samples from certain samples. Before diving into the experimental details, we want to emphasize that it is difficult to control the level of data difficulty (or uncertainty). First and foremost, human perception of data difficulty will be highly subjective and potentially different from its model counterpart. This limitation hinders the quantitative analysis, and thus some previous works had to rely on qualitative substitutes or analyze mislabeled samples which are impossible to control its difficulty [39, 35, 48]. Also, even if we could obtain sample-wise difficulty, it is often nontrivial to analyze the overall trend during training due to sheer data size.

To avoid the two challenges above, we borrow the settings from studies on long-tail visual recognition [33, 7]. [8] show that generalization error is bounded by the inverse square root of the dataset size. Further, many long-tail literature [33, 56, 19] have also empirically shown that it is hard for the deep neural network-based model to train with fewer samples, showing lower accuracy. Hence, we consider the major and minor classes as certain and uncertain classes, as the binned classification error is often used as the definition of confidence [16].

We train ResNet-18 [18] on the CIFAR10 dataset [28, 8] with an imbalance ratio of 10 for 30 epochs using the Adam optimizer [26]. Figure 8a shows the label distribution of the training dataset. Similar to [7], we choose classes 0, 1, 2, 3 and 4 as the major class and the rest as the minor class, randomly removing 90% of the training samples for the minor class. We reduce the inter-class differences of CIFAR10 by merging five classes into one, and demonstrate both the overall distribution and samplewise scores in Figure 2. We conclude that TD successfully captures data uncertainties, where its characteristics are more helpful in separating uncertain samples from certain samples than the information obtained from a model snapshot. Also, we empirically reaffirm that the major classes being more advantageous than minor classes in terms of accuracy during model training (Figure 8b, 8c).

C. Details on the TD Prediction Module

One can offer numerous alternatives on the design of the TD prediction module m , but we adopt the architecture of the loss prediction module [52] except for the last layer. By adopting the architecture used in the previous study, it is intended to show that the performance improvement of TiDAL does not come from adopting an advanced prediction module architecture, but from using TD. The TD prediction module takes several hidden feature maps extracted between the mid-level blocks of the target classifier f as inputs. Through a global average pooling layer and a fully-connected layer, each feature map is reduced to a fixed dimensional feature vector. All the reduced feature vectors are concatenated to take multi-level knowledge

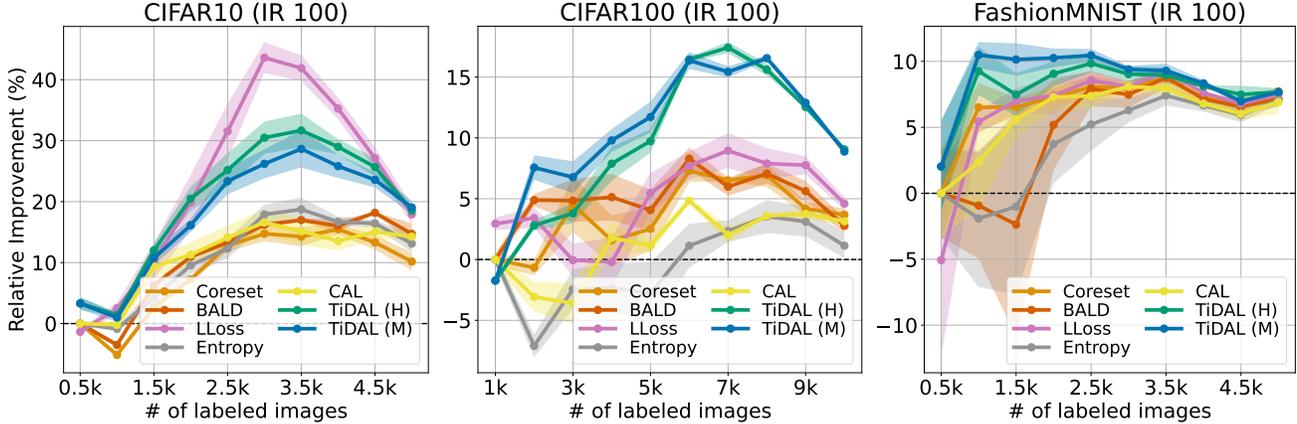


Figure 9: Averaged relative accuracy improvement curves and their 95% confidence interval (shaded) of AL methods over the number of labeled samples on synthetically imbalanced datasets. We use the imbalance ratio (IR) of 100 on CIFAR10, CIFAR100, and FashionMNIST.

of the target classifier into consideration for TD prediction. Using a single Softmax layer, the TD prediction module outputs a C -dimensional prediction $\hat{y}^{(t)} \in [0, 1]^C$, which are used as the predicted TD.

For a better understanding of the architecture of our TD prediction module m , please refer to [52].

D. Additional Experiments

We conduct additional experiments to further demonstrate the effectiveness of our method, TiDAL. We provide the dataset statistics in Table 1.

Table 1: Statistics of the dataset used for experiments.

Dataset	# of classes	# of samples	Imbalance ratio
CIFAR10	10	50k	{1, 10, 100}
CIFAR100	100	50k	{1, 10, 100}
FashionMNIST	10	60k	{1, 10, 100}
SVHN	10	73k	2.98
iNaturalist2018	8k	437k	500

D.1. Additional Results on Imbalanced Datasets

Figure 9 shows the experimental results on the imbalance ratio 100. Except for CIFAR10, our methods show superiority over other state-of-the-art methods.

D.2. Additional Results on Absolute Accuracy

Figure 10 and 11 provides the absolute accuracy plots for the completeness of the evaluation for real and synthetic data, respectively. We can observe the superiority of our method further on many of the settings.

D.3. Additional Baselines

Figure 12 compares our TiDAL with VAAL [45] and TA-VAAL [25]. Except for the case of CIFAR10 with the imbalance ratio of 100, both TiDAL strategies excel in performance. Note that both VAAL and TA-VAAL use a semi-supervised approach to train the selection module and further leverage the unlabeled data for training.

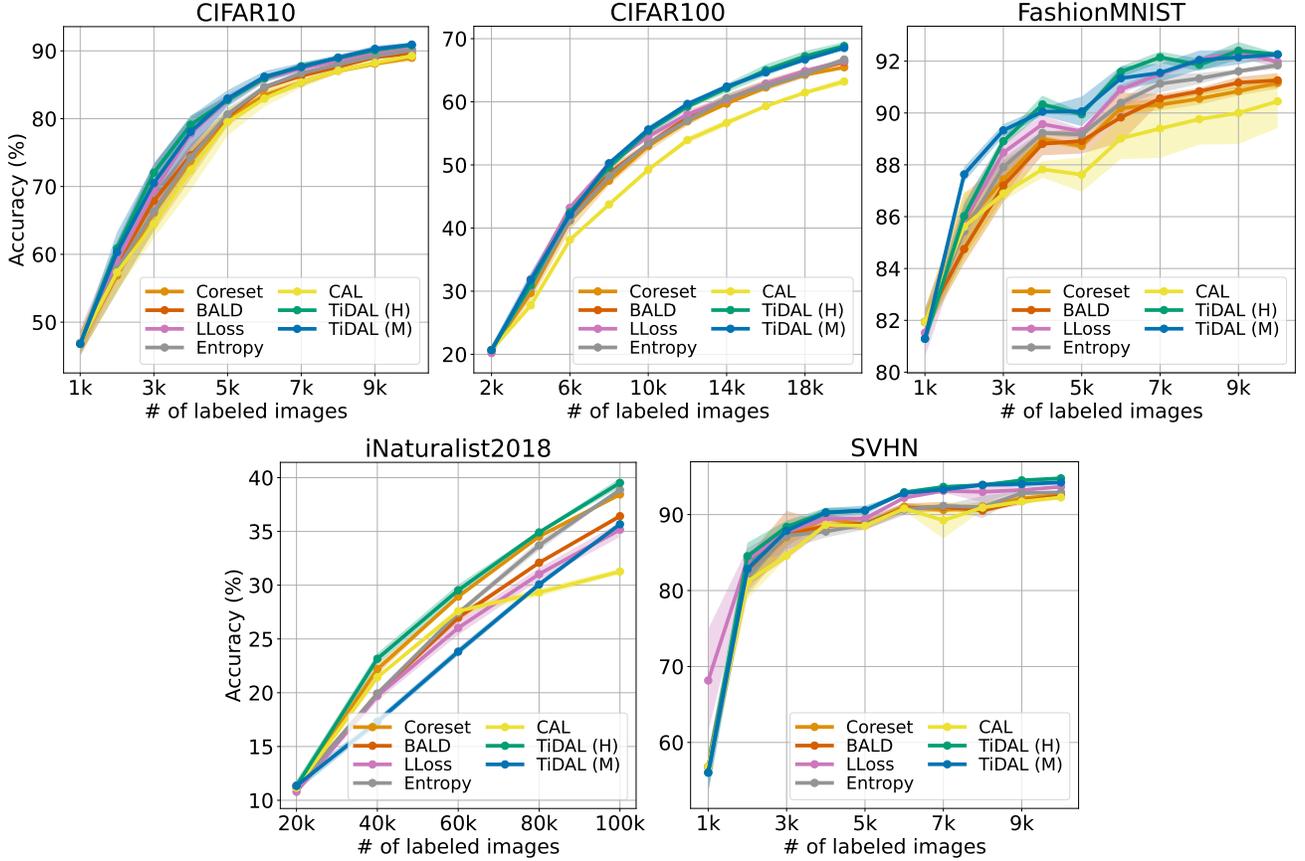


Figure 10: Averaged absolute accuracy improvement curves and its 95% confidence interval (shaded) of AL methods over the number of labeled samples on balanced and imbalanced datasets.

D.4. Variants of Training Dynamics-Aware Margin

We introduced two TD-aware strategies: entropy \bar{H} and margin \bar{M} , in §2. We further demonstrate various uncertainty estimation strategies as follows:

$$\bar{M}_0(\tilde{\mathbf{p}}_m) = \tilde{p}_m(\tilde{y}|x) - \max_{c \neq \tilde{y}} \tilde{p}_m(c|x), \quad (26)$$

$$\bar{P}(\tilde{\mathbf{p}}_m) = \tilde{p}_m(\hat{y}|x), \quad (27)$$

$$\bar{P}_0(\tilde{\mathbf{p}}_m) = \tilde{p}_m(\tilde{y}|x), \quad (28)$$

where $\tilde{y} = \operatorname{argmax}_c \tilde{p}_m(c|x)$ is the class of the maximum module output.

\hat{M}_0 is the naive variant of the margin \hat{M} where it does not utilize the predicted label \hat{y} of the target classifier f . It calculates the margin between the biggest and the second biggest outputs of the module m . \bar{P} uses the module output on the predicted label \hat{y} from the target classifier f and \bar{P}_0 is the naive variant of \bar{P} that uses the maximum output of the module m .

Figure 13 shows the average accuracy of three runs for the entropy \bar{H} and margin \bar{M} , where we show the accuracy of a single run for other strategies. We can observe that the naive variant of the margin \bar{M}_0 generally underperforms compared to the margin \bar{M} except CIFAR100 with the imbalance ratio of 100. There seems to be no clear dominance between \bar{P} and its naive variant \bar{P}_0 . However, both \bar{P} and \bar{P}_0 perform moderately well on both CIFAR100 and FashionMNIST despite its simplicity. Future studies may concentrate on broader query strategies based on various training dynamics and its module predictions.

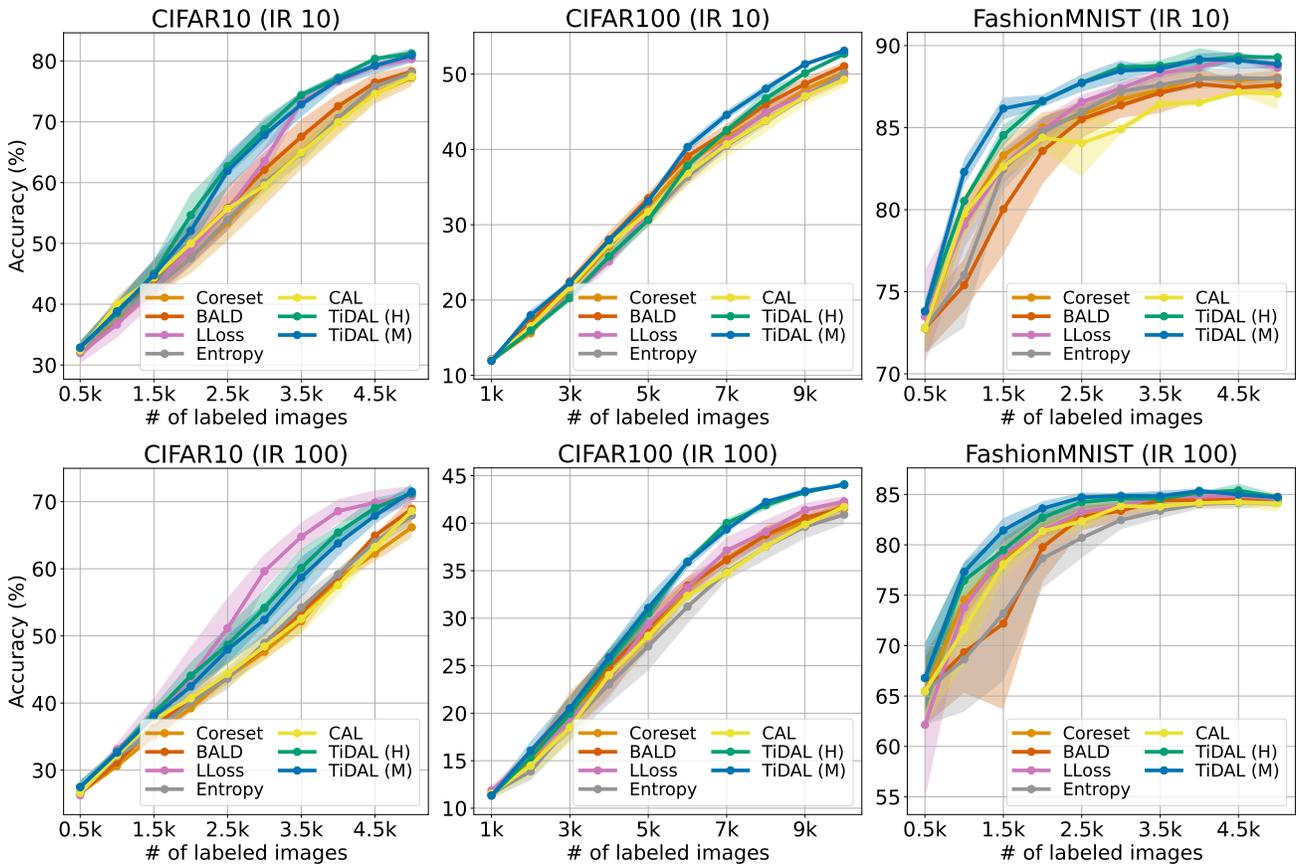


Figure 11: Averaged absolute accuracy improvement curves and their 95% confidence interval (shaded) of AL methods over the number of labeled samples on synthetically imbalanced datasets. We use the imbalance ratio (IR) of 10 and 100 on CIFAR10, CIFAR100, and FashionMNIST.

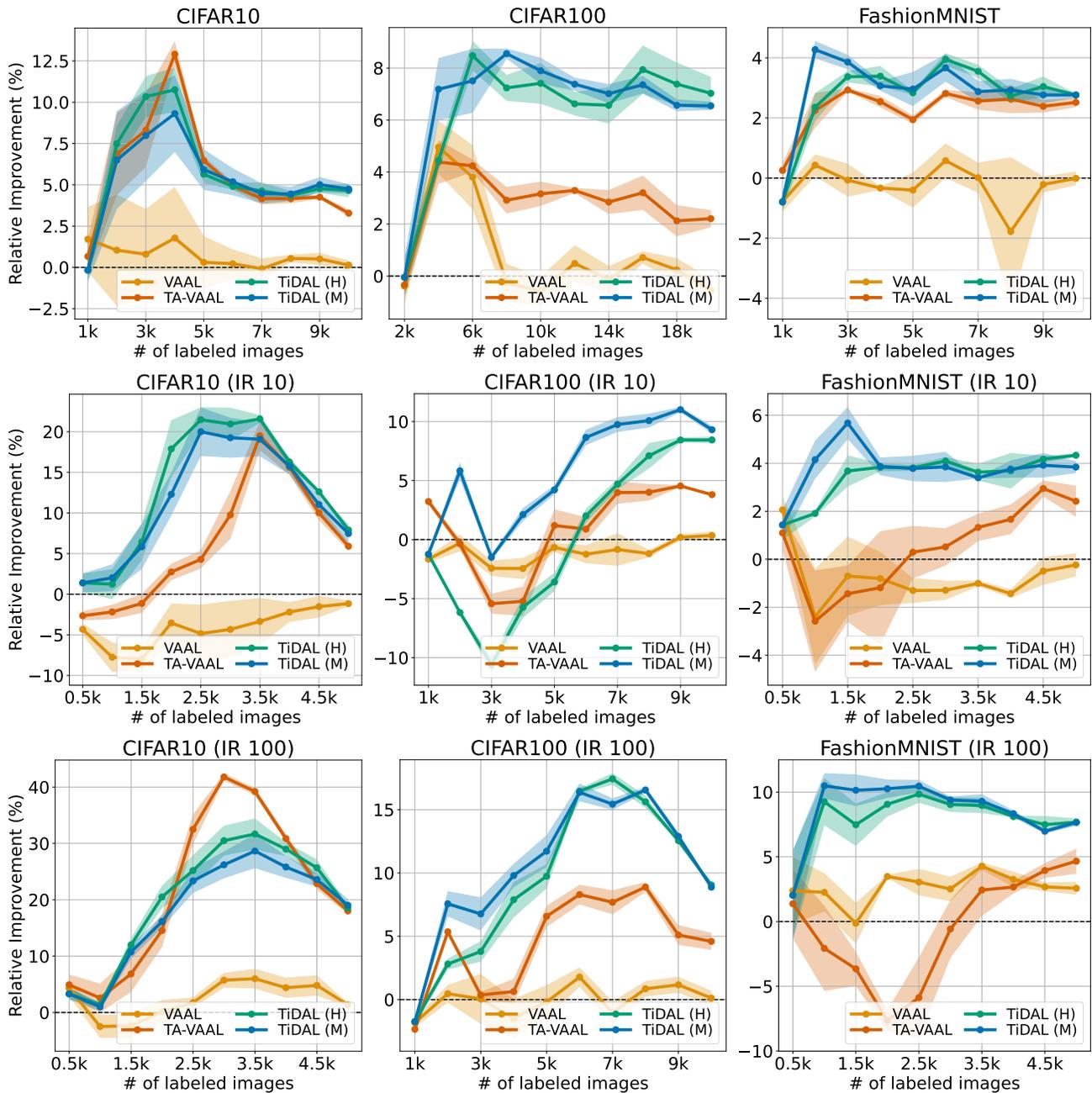


Figure 12: Averaged relative accuracy improvement curves and their 95% confidence interval (shaded) of AL methods over the number of labeled samples on balanced and synthetically imbalanced datasets. We use the imbalance ratio (IR) of 10 and 100 on CIFAR10, CIFAR100, and FashionMNIST to synthetically imbalance the dataset.

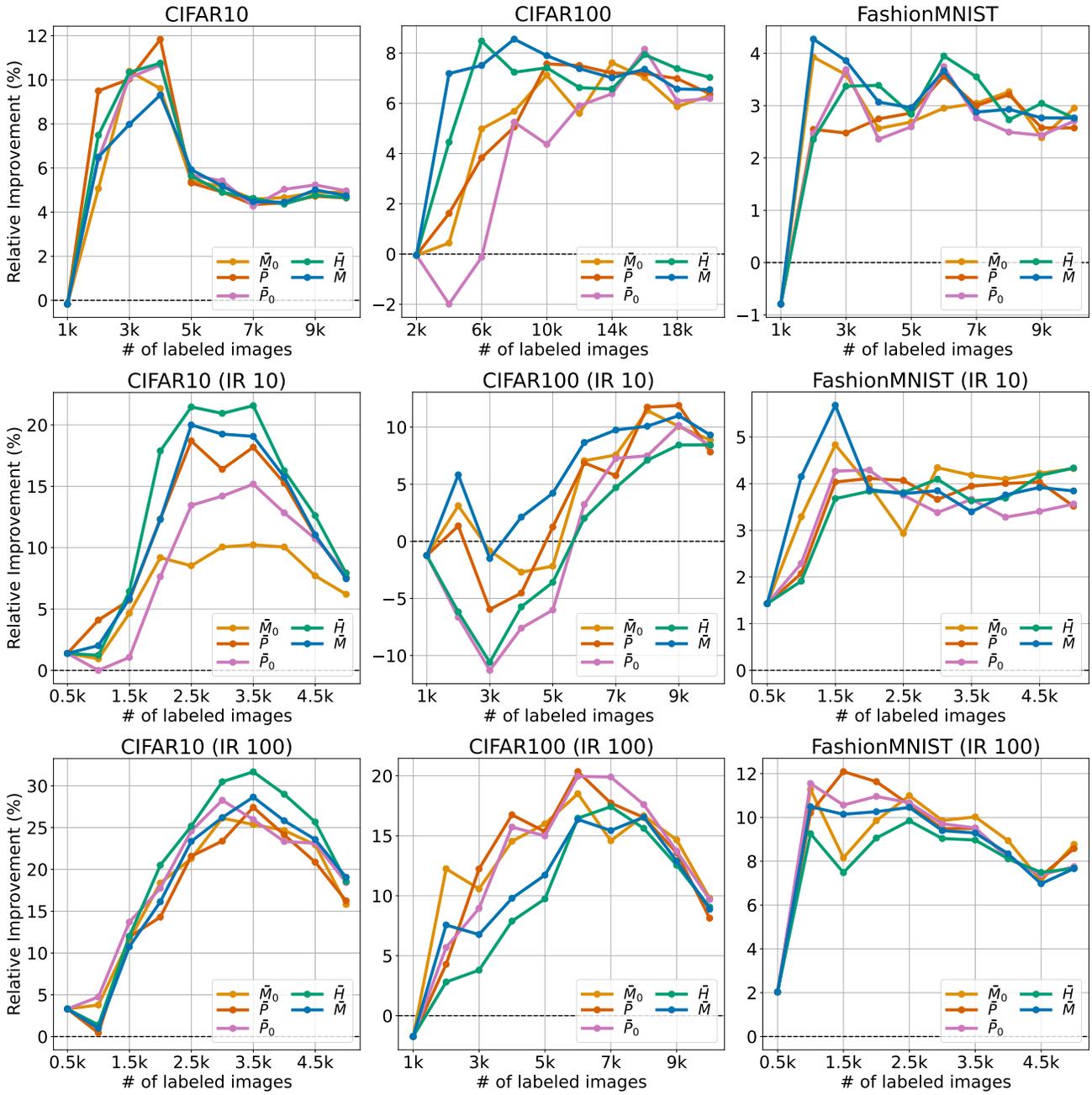


Figure 13: Averaged relative accuracy improvement curves of different uncertainty estimation strategies over the number of labeled samples on balanced and synthetically imbalanced datasets. We use the imbalance ratio (IR) of 10 and 100 on CIFAR10, CIFAR100, and FashionMNIST to synthetically imbalance the dataset.