

# On the Calibration of Image Semi-Supervised Learning Models

Mehrab Mustafy Rahman, Jayanth Mohan, Tiberiu Sosea, Cornelia Caragea

Computer Science

University of Illinois Chicago

mrahm@uic.edu, jmohal1@uic.edu, tsosea2@uic.edu, cornelia@uic.edu

## Abstract

Semi-supervised learning (SSL) has demonstrated high performance in image classification tasks by effectively utilizing both labeled and unlabeled data. However, existing SSL methods often suffer from poor calibration, with models yielding overconfident predictions that misrepresent actual prediction likelihoods. Recently, neural networks trained with `mixup` that linearly interpolates random examples from the training set have shown better calibration in supervised settings. However, calibration of neural models remains under-explored in semi-supervised settings. Although effective in supervised model calibration, random mixup of pseudolabels in SSL presents challenges due to the overconfidence and unreliability of pseudolabels. In this work, we introduce CalibrateMix, a targeted mixup-based approach that aims to improve the calibration of SSL models while maintaining or even improving their classification accuracy. Our method leverages training dynamics of labeled and unlabeled samples to identify “easy-to-learn” and “hard-to-learn” samples, which in turn are utilized in a targeted mixup of easy and hard samples. Experimental results across several benchmark image datasets show that our method achieves lower expected calibration error (ECE) and superior accuracy compared to existing SSL approaches.

**Code** — <https://github.com/mehrab-mustafy/CalibrateMix>

## Introduction

Deep neural networks (DNNs) (LeCun, Bengio, and Hinton 2015; Mathew, Amudha, and Sivakumari 2021) have achieved remarkable success across a wide range of computer vision tasks, including image classification (Lu and Weng 2007), object detection (Papageorgiou, Oren, and Poggio 1998), and semantic segmentation (Minaee et al. 2021). However, alongside accuracy, the predictive confidence of the models plays a vital role in real-world decision-making applications. For example, in critical applications such as autonomous driving, medical diagnosis, and disaster response, models must be accurate as well as reliably indicate when they are uncertain, so that additional safety measures can be triggered. For this reason, quantifying predictive uncertainty and calibration of the DNNs is a pivotal component toward building more reliable models. Despite strong performance,

DNNs often suffer from poor calibration (Guo et al. 2017), which means that the predictive confidence likely overestimates the model’s true accuracy. A key reason behind this is that modern DNNs are trained using one-hot encoded labels and the cross-entropy loss, which assumes that every training sample belongs with full certainty to a single class. This forces the model to assign the entire probability mass to a single class label, which in turn suppresses any expression of uncertainty even for “ambiguous” samples. To overcome the issues of overconfidence, label smoothing (Müller, Kornblith, and Hinton 2019) has been introduced. By softening the target distribution during training, label smoothing regularizes the model’s output probabilities, encouraging it to remain uncertain where appropriate. More recently, Thulasidasan et al. (2019) explored the use of mixup training (Zhang et al. 2018) for improving model calibration which creates augmented samples through convex combinations of input samples and their labels. Mixup distributes its probability into two classes, which introduces entropy, prevents overconfidence, and has proven to be an effective tool in model calibration.

However, this approach primarily targets the fully supervised setting, which requires a large amount of labeled data. With the emergence of AI in all domains, it is impractical to obtain large amounts of annotated data for every domain. To solve this issue, Semi-supervised learning (SSL) (Chapelle, Scholkopf, and Zien 2009) can be an effective strategy to leverage large amounts of unlabeled examples during training to boost model performance. Despite mixup being successful in supervised learning, mixing up unlabeled examples in an SSL setting poses certain challenges due to the uncertainty of pseudo-label correctness, especially at the early iterations of training. This makes it critical to ensure the quality of pseudo-labels before incorporating them into the learning process. To ensure the quality of pseudo-labels, a popular SSL method to learn from unlabeled examples is pseudo-labeling (Lee et al. 2013), which leverages a model to make predictions on unlabeled examples and assign them pseudo-labels, which are in turn used as (pseudo) ground truth during training. To ensure that correct pseudo-labels are used for model training, modern SSL frameworks such as FixMatch (Sohn et al. 2020) and FlexMatch (Zhang et al. 2021) utilize high confidence thresholds to maintain data quality and filter out potentially incorrect examples. However, the calibration of these SSL models is not well studied, and we found empirical evidence

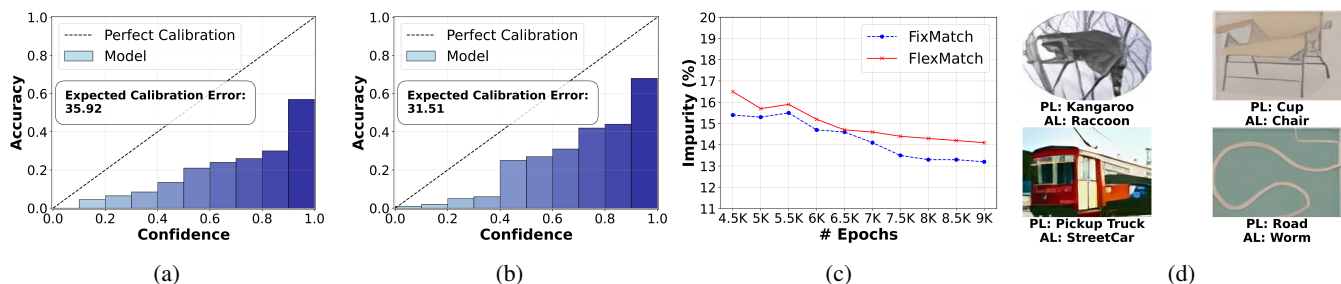


Figure 1: (a) FixMatch reliability diagram on CIFAR-100; (b) FlexMatch reliability diagram on CIFAR-100; (c) Impurity comparison on CIFAR-100; (d) Examples of incorrect pseudo-labels at the end of training. PL denotes the pseudo label predicted by the model, while AL denotes the ground-truth (actual) label. Mislabeling patterns for FlexMatch (left) and FixMatch (right) are shown.

that SSL models also suffer from poor calibration as shown in the reliability diagrams of FixMatch and FlexMatch on CIFAR-100 in Figures 1a and 1b, respectively. The diagrams, which plot accuracy as a function of confidence show that the confidence estimates of the models are not indicative of their correctness. Notably, FixMatch and FlexMatch predictions with confidences higher than 90% have less than 65% accuracy, contradicting the assumption that employing high confidence thresholds leads to high pseudo-label quality.

As shown in Figure 1c, the impurity of unlabeled data for both FixMatch and FlexMatch on CIFAR-100 (Krizhevsky, Hinton et al. 2009) is higher than 13%, indicating that more than 13% of the unlabeled data is utilized with incorrect pseudo-labels during training even at the later stages of training at 9000 epoch. Figure 1d shows examples of incorrect pseudo-labels introduced by both FixMatch and FlexMatch at the end of training. These incorrect pseudo-labels arise due to miscalibrated model predictions, manifested here as overconfidence in potentially incorrect predictions. This calibration gap in SSL can not be fully solved using random mixup. When incorrect pseudo-labels are used in mixup, the interpolation process can propagate label noise across training samples. This not only reinforces errors but also makes them harder to correct. For deep learning models that are highly over-parameterized and capable of achieving near-zero training error, such reinforcement of incorrect labels can lead to severe overfitting to mislabeled data. Thus, without addressing the underlying calibration issue and pseudo-label reliability, mixup can become a mechanism for error amplification rather than a regularization strategy. Hence, we investigate the calibration of SSL models under mixup and propose an enhanced mixup strategy.

In this paper, we propose CalibrateMix, a novel, targeted mixup-based (Zhang et al. 2018) framework that enhances the confidence calibration of SSL models. CalibrateMix monitors the training dynamics of unlabeled samples during training by keeping track of margins (Bartlett, Foster, and Telgarsky 2017; Elsayed et al. 2018; Sosea and Caragea 2023) of the outputs of the model. These margins are utilized to characterize each sample based on learning difficulty into two categories: “easy-to-learn” samples, which the model perceives to be correctly pseudo-labeled and “hard-to-learn”

samples, which the model perceives to be potentially incorrect or ambiguous. Then, our method performs targeted mixup, combining easy-to-learn labeled samples with hard-to-learn unlabeled samples, as well as hard-to-learn labeled samples with easy-to-learn unlabeled samples. This mixup strategy has two main advantages: First, the presence of easy-to-learn samples ensures that the resulting pseudo-label of the mixed sample is qualitative while the presence of the hard-to-learn samples ensures a more difficult environment for the model to learn from, which was shown to benefit SSL training (Xie et al. 2020). Second, CalibrateMix is easy-to-use in practice and can be applied on top of the most popular SSL frameworks such as FixMatch or FlexMatch.

To showcase the benefits of our method, we carry out extensive experiments on well-established small-scale SSL setups - CIFAR10, CIFAR100 (Krizhevsky, Hinton et al. 2009), SVHN (Netzer et al. 2011), STL10 (Coates, Ng, and Lee 2011); and on two large-scale benchmarks - ImageNet (Deng et al. 2009) and WebVision (Li et al. 2017), where we show that CalibrateMix outperforms strong baselines and previous works both in terms of accuracy and expected calibration error (ECE). Notably, CalibrateMix boosts the performance of FixMatch by 8.51% in ECE and by 1.54% in error rate on CIFAR-100 using only 25 labels per class. On the larger datasets, e.g., WebVision, we also see an improvement of 2% ECE on FixMatch and FlexMatch without loss in error rates.

Our contributions are as follows:

1. We introduce CalibrateMix, a novel targeted mixup framework that is primarily designed to improve the confidence calibration of semi-supervised learning (SSL) models.
2. Through extensive experiments on standard benchmark image datasets, we demonstrate that CalibrateMix outperforms existing methods in terms of Expected Calibration Error (ECE), while also often achieving lower error rates compared to multiple state-of-the-art SSL setups. Our method is compatible with any SSL frameworks, and we demonstrate its effectiveness by improving both calibration and, in multiple cases, accuracy when integrated with FixMatch, FlexMatch, and SoftMatch.
3. We carry out extensive ablation studies to understand the contribution of various components to the success of CalibrateMix, as well as an error analysis to analyze the

calibration and correctness of predictions in practice.

## Related Work

Despite achieving high accuracy, modern deep networks frequently remain miscalibrated, producing overly confident predictions (Guo et al. 2017; Minderer et al. 2021). Several effective methods have been proposed in supervised learning settings to mitigate calibration issues proactively. Dropout, initially introduced as a regularization technique, doubles as a Bayesian uncertainty estimator at inference via Monte Carlo sampling, providing improved calibration by averaging predictions across multiple stochastic forward passes (Srivastava et al. 2014). Another well established approach is label smoothing, which softens the one-hot encoded targets, thereby discouraging extreme predictions and significantly enhancing calibration and generalization performance, especially under limited or imbalanced data scenarios (Müller, Kornblith, and Hinton 2019). Mixup training, another promising approach, interpolates pairs of training samples and their labels, generating augmented data points that promote smoother predictive boundaries. Thulasidasan et al. (Thulasidasan et al. 2019) explicitly verified mixup’s effectiveness in reducing expected calibration error (ECE), highlighting its capability in generating well-calibrated confidence scores and robust predictive uncertainty. Further developments, such as RegMixup, have validated that Mixup based regularization not only enhances calibration but also robustness under distribution shifts, particularly beneficial in out-of-distribution (OOD) detection tasks (Pinto et al. 2022).

Other studies demonstrated that employing focal loss, which emphasizes learning from challenging predictions, also effectively mitigates model overconfidence, enhancing calibration and model reliability (Mukhoti et al. 2020). While supervised calibration methods have matured significantly, translating these advancements directly to semi-supervised learning (SSL) remains challenging. SSL leverages limited labeled data alongside abundant unlabeled data to enhance model learning. Prominent SSL frameworks, such as FixMatch (Sohn et al. 2020), FlexMatch (Zhang et al. 2021), and SoftMatch (Chen et al. 2023), primarily rely on pseudo-labeling and consistency regularization. FixMatch employs fixed high confidence thresholds to filter pseudo-labels, improving overall accuracy but inadvertently exacerbating calibration issues by neglecting calibration aware decision-making. Similarly, FlexMatch introduces a dynamic curriculum based threshold adjustment, enhancing class specific learning but still neglecting explicit calibration mechanisms (Zhang et al. 2021). SoftMatch further improves pseudo-label quality through adaptive weighting strategies; however, calibration remains an implicit, secondary consideration rather than a deliberate design goal (Chen et al. 2023). Consequently, SSL models often remain miscalibrated, assigning overly confident pseudo-labels to unlabeled samples, which compromises their accuracy and robustness. Recent SSL approaches, such as MarginMatch (Sosea and Caragea 2023), SequenceMatch (Nguyen 2024), and FineSSL (Gan and Wei 2024), indirectly target pseudo-label reliability and robustness through pseudo-margin analysis, consistency across augmentations, and robust finetuning strategies, respectively.

While these techniques implicitly reduce calibration errors, their main emphasis remains on enhancing pseudo-label quality or robustness rather than directly addressing calibration gaps. Thus, explicit calibration remains largely unaddressed. Motivated by this necessity, we propose CalibrateMix, a targeted mixup-based SSL approach explicitly designed to improve both model calibration and predictive accuracy.

## CalibrateMix

**Notation** Let  $D_L = \{(x_1, y_1), \dots, (x_{B_L}, y_{B_L})\}$  be a batch of labeled samples of size  $B_L$  and  $D_U = \{\hat{x}_1, \dots, \hat{x}_{B_U}\}$  be a batch of unlabeled samples of size  $B_U$ .

## Background

Most deep neural networks (DNNs) trained for classification are trained using one-hot encoded labels, where the entire probability mass is assigned to a single class. This leaves no room for uncertainty during training. As a result, models tend to become overconfident in their predictions. Hence, it is not surprising that modern DNNs are poorly calibrated (Guo et al. 2017). To prevent this overconfidence, Thulasidasan et al. (2019) explored the impacts of mixup training in supervised settings and found that mixup has proven to be effective in reducing miscalibration in supervised settings. Mixup training creates vicinity samples for training, effectively creating more samples for the model to learn from. The mixup augmented samples are generated by the following rule as mentioned in Zhang et al. (2018):

$$\tilde{x} = \gamma x_i + (1 - \gamma)x_j \quad (1)$$

$$\tilde{y} = \gamma y_i + (1 - \gamma)y_j \quad (2)$$

where,  $x_i$  and  $x_j$  are any two randomly selected samples from the train set and  $y_i$  and  $y_j$  are their corresponding one-hot labels. Mixup-augmented sample  $\tilde{x}$  is generated by interpolating between  $x_i$  and  $x_j$ . Similarly, the corresponding label  $\tilde{y}$  is obtained by mixing  $y_i$  and  $y_j$ . The mixup coefficient  $\gamma$ , comes from the Beta distribution, with the hyperparameter  $\alpha$  controlling the level of interpolation between the two samples. Unlike standard one-hot labels, when  $y_i \neq y_j$ , the mixup label  $\tilde{y}$  distributes the probability mass across two classes rather than a single class. This distribution introduces entropy and uncertainty into the training process, preventing the model from becoming overly confident in its predictions.

In SSL, selecting samples at random for mixup can be harmful because of the unreliability of pseudo-labels especially in the early stages of training when models are more prone to errors and because of the confirmation bias. This is particularly severe for difficult samples. From Figure 1, we observe that the pseudo-labels may remain incorrect at the end of training, and SSL models can produce high-confidence incorrect predictions. Such miscalibrated predictions not only hamper model learning progress but also risk propagating errors across training iterations. Given that pseudo-labeled data typically outnumbers labeled data in SSL settings, random pairing during mixup increases the chance of combining two erroneous pseudo-labels. This can yield misleading targets that further reinforce incorrect representations. Hence, to reduce the impact of noisy labels in training and ensure quality,

it is important to avoid mixup between two difficult samples or two pseudo-labeled samples.

### Proposed Approach: CalibrateMix

To address the limitations of standard random mixup in SSL and to deal with the propagation of overconfident or noisy pseudo-labels we propose **CalibrateMix**, a targeted mixup-based framework. CalibrateMix performs a controlled mixup between labeled and pseudo-labeled samples guided by their learning difficulty to generate higher-quality augmented training samples. The inclusion of labeled samples in the mixup ensures that no two pseudo-labeled samples take part in the mixup at the same time, while the inclusion of pseudo-labeled samples introduces entropy and uncertainty. This helps prevent the reinforcement of noise and encourages the model to express necessary uncertainty, leading to better calibrated confidence estimates. Hence, the core advantage of CalibrateMix lies in its structured pairing strategy that mixes labeled and unlabeled samples based on their learning difficulty. To quantify the learning difficulty of the samples of the labeled data we monitor the Area Under the Margin (AUM) (Pleiss et al. 2020) of the outputs of the model at each iteration. For the unlabeled data, we monitor the Average Pseudo Margin (APM) (Sosea and Caragea 2023) of the outputs of the model at each iteration. We do this because relying solely on the model’s current prediction confidence is insufficient. Confidence at a single iteration does not reliably reflect the sample’s margin and the correctness of its pseudo-label, as shown in Sosea and Caragea (2023). To calculate APM for an unlabeled sample  $\hat{x}_i$ , we use pseudo-margins defined as:

$$\text{PM}_c^t(\hat{x}_i) = z_c(\hat{x}_i) - \max_{j \neq c} z_j(\hat{x}_i) \quad (3)$$

where at iteration  $t$ ,  $z_c(\hat{x}_i)$  is the logit corresponding to sample  $\hat{x}_i$  for assigned pseudo-label  $c$  and  $\max_{j \neq c} z_j(\hat{x}_i)$  is the largest other logit corresponding to a label  $j$  other than  $c$  for the same sample  $\hat{x}_i$ . We use the pseudo-labels at the current iteration  $t$  as the “ground-truth”. Then, the average pseudo-margin (APM) for the unlabeled sample  $\hat{x}_i$  with pseudo-label  $c$  at iteration  $t$  is defined as follows:

$$\text{APM}_c^t(\hat{x}_i) = \frac{1}{t} \sum_{e=1}^t (\text{PM}_c^e(\hat{x}_i)) \quad (4)$$

An important consideration here is that for any previous iteration  $t'$ , if the pseudo-label was  $c'$  (where  $c \neq c'$ ) the pseudo-margin  $\text{PM}^{t'}$  is calculated with respect to  $c'$  and also the APM is averaged from 1 to  $t'$  with respect to  $c'$ . In practice, we maintain a vector of pseudo-margins for all classes accumulated over the training iterations and dynamically retrieve the accumulated pseudo-margin value of the argmax class  $c$  to obtain the  $\text{APM}_c^t$  at iteration  $t$ .

To handle old pseudo-margin deprecation across large number of iterations, we use an exponential moving average of pseudo-margins to place higher importance on recent iterations. Hence, APM follows:

$$\text{APM}_c^t(\hat{x}_i) = \text{PM}_c^t(\hat{x}_i) \cdot \frac{\delta}{1+t} + \text{APM}_c^{t-1}(\hat{x}_i) \cdot \left(1 - \frac{\delta}{1+t}\right) \quad (5)$$

Similar to Sosea and Caragea (2023), we set the smoothing parameter  $\delta = 0.997$ .

For both labeled and unlabeled data, samples with larger AUM and APM values, respectively, are considered easier for the model to learn, and those with smaller values are generally ambiguous, harder to learn, or mislabeled. These hard to learn or ambiguous samples are often the cause of overconfidence and errors. Hence, in the mixup, we include one easy sample and one hard sample. The easy-to-learn sample ensures that the resulting pseudo-label of the mixed sample is qualitative, while the presence of the hard-to-learn sample ensures more uncertainty for the model. To perform these difficulty-aware splits, we partition both labeled and pseudo-labeled data at each iteration into “easy-to-learn” and “hard-to-learn” subsets. This categorization is based on the batch medians of AUM for labeled samples ( $\tau_L$ ) and APM for pseudo-labeled samples ( $\tau_U$ ), following (Park and Caragea 2022) but in a batch-wise manner, effectively distinguishing samples by their learning difficulty.

After obtaining the difficulty-aware asymmetric splits, mixup is then performed as follows: for each easy labeled sample, we randomly select one from its *top-k* most dissimilar hard pseudo-labeled samples, and for each hard labeled sample, we do the same with the *top-k* most dissimilar easy pseudo-labeled samples. This dissimilarity based sample selection for mixup better simulates the out-of-domain (OOD) distribution by producing mixup augmented examples that deviate from the in-domain distribution. Furthermore, randomly selecting a sample from the *top-k* most dissimilar samples introduces diversity across iterations and avoids overfitting. As a result, the model becomes less likely to make overconfident predictions on ambiguous inputs, thereby improving calibration. For the mixup, we compute a convex combination of the two inputs using weights  $\gamma$  and  $1 - \gamma$  as follows:

$$\tilde{x}_1 = \gamma x_{le} + (1 - \gamma) \hat{x}_{uh} \quad (6)$$

$$\tilde{y}_1 = \gamma y_{le} + (1 - \gamma) \hat{y}_{uh} \quad (7)$$

$$\tilde{x}_2 = \gamma x_{lh} + (1 - \gamma) \hat{x}_{ue} \quad (8)$$

$$\tilde{y}_2 = \gamma y_{lh} + (1 - \gamma) \hat{y}_{ue} \quad (9)$$

where  $x_{le}, y_{le}$  are easy labeled samples and their ground-truth labels,  $x_{lh}, y_{lh}$  are hard labeled samples and their ground-truth labels,  $\hat{x}_{ue}, \hat{y}_{ue}$  are easy pseudo-labeled samples and their pseudo-labels, and  $\hat{x}_{uh}, \hat{y}_{uh}$  are hard pseudo-labeled samples and their pseudo-labels.

The final model input includes labeled, unlabeled, and mixup-augmented data. The total loss  $\mathcal{L}$  is computed as:

$$\mathcal{L} = \mathcal{L}_L + \lambda_U \cdot \mathcal{L}_U + \mathcal{L}_{\text{mixup}} \quad (10)$$

where  $\mathcal{L}_L$  is the supervised loss on labeled data,  $\mathcal{L}_U$  is the unsupervised loss on the unlabeled data and  $\lambda_U$  is a hyperparameter that controls the weight of  $\mathcal{L}_U$  in the total loss  $\mathcal{L}$ , and  $\mathcal{L}_{\text{mixup}}$  is the loss on mixup-augmented samples. In experiments, we set  $\lambda_U$  equals 1 consistent with previous approaches (Sohn et al. 2020; Chen et al. 2023).

The CalibrateMix algorithm is shown in Algorithm 1.

## Experimental Analysis

We test the performance of our approach across a variety of image benchmarks by following standard SSL settings (Chen

---

**Algorithm 1: CalibrateMix**

---

**Require:** Labeled batch  $D_L$ ; unlabeled batch  $D_U$ ; model  $\theta$ ; number of classes  $C$ ; weak augmentation  $\pi$ ; strong augmentation  $\Pi$ ; mixup coefficient  $\gamma$ ; number of dissimilar samples  $k$ ; confidence threshold  $\omega$ ; labeled-to-unlabeled loss ratio  $\lambda$ ; medians of labeled and unlabeled batches  $\tau_L, \tau_U$ ; current iteration  $t$ .

- 1: Compute pseudo-labels for the unlabeled batch:  
 $\hat{y}_i^t = \arg \max_{c \in \{1, \dots, C\}} p_\theta(\pi(\hat{x}_i))$  for  $\hat{x}_i \in D_U$  and build  $\hat{D}_U = \{(\hat{x}_i, \hat{y}_i^t) \mid \hat{x}_i \in D_U\}$
- 2: Calculate  $AUM^t$  for samples in  $D_L$  and  $APM^t$  for samples in  $D_U$
- 3: **Split** the labeled and unlabeled batch based on samples' learning difficulty:
- 4:  $D_{Leasy} \leftarrow \{(x_i, y_i) \in D_L \mid AUM^t(x_i, y_i) \geq \tau_L\}$
- 5:  $D_{Lhard} \leftarrow \{(x_i, y_i) \in D_L \mid AUM^t(x_i, y_i) < \tau_L\}$
- 6:  $\hat{D}_{Ueasy} \leftarrow \{(\hat{x}_i, \hat{y}_i^t) \in \hat{D}_U \mid APM^t(\hat{x}_i, \hat{y}_i^t) \geq \tau_U\}$
- 7:  $\hat{D}_{Uhard} \leftarrow \{(\hat{x}_i, \hat{y}_i^t) \in \hat{D}_U \mid APM^t(\hat{x}_i, \hat{y}_i^t) < \tau_U\}$
- 8: **Mix** easy labeled with hard unlabeled and hard labeled with easy unlabeled following Eqs (6), (7), (8), (9):
- 9:  $\hat{D}_{Mix1} \leftarrow \text{Mixup}(D_{Leasy}, \text{top-}k \text{ dissim. from } \hat{D}_{Uhard})$
- 10:  $\hat{D}_{Mix2} \leftarrow \text{Mixup}(D_{Lhard}, \text{top-}k \text{ dissim. from } \hat{D}_{Ueasy})$
- 11:  $D_M \leftarrow \hat{D}_{Mix1} + \hat{D}_{Mix2}$
- 12: **Optimize** total loss:  $\mathcal{L} \leftarrow \mathcal{L}_L + \lambda \cdot \mathcal{L}_U + \mathcal{L}_{Mixup}$ , where
- 13:  $\mathcal{L}_L = \frac{1}{|D_L|} \sum_{k=1}^{|D_L|} H(y_k, p_\theta(\pi(x_k)))$
- 14:  $\mathcal{L}_U = \frac{1}{|D_U|} \sum_{k=1}^{|D_U|} \mathbb{1}(\max(p_\theta(\pi(\hat{x}_k))) \geq \omega) \cdot H(\hat{y}_k^t, p_\theta(\Pi(\hat{x}_k)))$
- 15:  $\mathcal{L}_{Mixup} = \frac{1}{|D_M|} \sum_{k=1}^{|D_M|} H(\tilde{y}_k, p_\theta(\tilde{x}_k))$

---

et al. 2023; Wang et al. 2022). Specifically, we conduct experiments varying the amounts of labeled samples on CIFAR-10, CIFAR-100 (Krizhevsky, Hinton et al. 2009), SVHN (Netzer et al. 2011), STL-10 (Coates, Ng, and Lee 2011), ImageNet (Deng et al. 2009), and WebVision (Li et al. 2017). For the CIFAR-10, CIFAR-100, SVHN, and STL-10 datasets, we follow Sohn et al. (2020) and randomly select a small number of labeled samples ranging from 4 labels per class up to 400 labels per class and treat the remaining samples as unlabeled data, except for STL-10, which has its own set of unlabeled samples. For the ImageNet and WebVision datasets, we use 10% of the available labeled samples as labeled data and the remaining 90% as unlabeled data. In addition to the SSL setups, we also show results on the Fully Supervised setting.

We run all our experiments three times and report the average ECE and error rates, as well as their standard deviations. Similar to Sohn et al. (2020), we utilize Wide Residual Networks (Zagoruyko 2016) for small-scale datasets: WRN-28-2 for CIFAR-10 and SVHN, WRN-37-2 for STL-10 and WRN-28-8 for CIFAR-100; and use ResNet-50 (He et al. 2016) for the large-scale ImageNet and WebVision. Additionally, we utilize the SGD optimizer with a momentum of 0.9 and an initial learning rate of 0.03. We use a cosine annealing schedule to dynamically adjust the learning rate over a total of  $2^{20}$  training steps. For each dataset, the batch size for labeled data and mixup data is set to 64, while the batch size for unlabeled data is configured to be seven times larger than that of the labeled data. We also use weak (flip and shift) and strong (RandAugment (Cubuk et al. 2020)) augmentations for the

unlabeled data. We set mixup coefficient  $\gamma$  to be 0.4 similar to prior work (Thulasidasan et al. 2019). For CalibrateMix, we include warmup during training with the warmup phase consisting of approximately 100 epochs.

**Performance on CIFAR-10, CIFAR-100, STL-10, SVHN**

We compare CalibrateMix to relevant SSL approaches: FixMatch (Sohn et al. 2020), FlexMatch (Zhang et al. 2021), SoftMatch (Chen et al. 2023) and report ECE and top-1 error rates in Table 1. We further compare CalibrateMix against two calibration techniques: Random Mixup and Label Smoothing (LS) in Table 2. For Random Mixup, we randomly pair two samples from the combined pool of labeled and unlabeled data. Additionally, we apply LS following Müller, Kornblith, and Hinton (2019), with a factor of  $\gamma = 0.1$  to the labeled, unlabeled, and mixup samples for all methods: FixMatch, FlexMatch, SoftMatch, and CalibrateMix (including combining CalibrateMix with LS). The results for Label Smoothing are denoted as ‘LS’ and the results with Random Mixup are denoted as ‘RandomMixup’ in Table 2.

Across all datasets, CalibrateMix consistently improves model calibration and, in many cases, enhances classification accuracy as can be seen from Table 1. For example, when integrated with FixMatch, CalibrateMix yields better-calibrated models. On CIFAR-10 with only 4 labels per class, CalibrateMix achieves a reduction in ECE by 2.18%. On the more challenging CIFAR-100 dataset with 25 labels per class, the improvements are more significant, reducing ECE by 8.51% and error rate by 1.54%. On the more realistic STL-10 dataset, CalibrateMix reduces ECE by 2.73%, further demonstrating its robustness. CalibrateMix also significantly enhances FlexMatch, particularly under low-label settings. Furthermore, integrating CalibrateMix with SoftMatch provides significant gains. On CIFAR-100 with 25 labels per class, ECE is improved by 5.68%. On STL-10 with 4 labels per class, CalibrateMix achieves a 0.84% lower ECE and a 5.72% reduction in error rate compared to SoftMatch. We can also observe that CalibrateMix yields well calibrated models in the Fully Supervised setting.

CalibrateMix consistently demonstrates improvements in model calibration and often in error rates when compared to prior calibration strategies such as random mixup and label smoothing, as shown in Table 2. When applied to FixMatch, CalibrateMix outperforms random mixup across several settings. Notably, on CIFAR-100 and STL-10 with 4 labels per class, it achieves ECE reductions of 2.62% and 5.65%, along with error rate reductions of 6.48% and 3.39%, respectively. Furthermore, on CIFAR-100 with 25 labels per class, combining CalibrateMix with label smoothing results in a 4.48% reduction in ECE over using label smoothing alone with FixMatch. On FlexMatch, similar benefits are observed. When integrated with SoftMatch, CalibrateMix continues to enhance calibration by outperforming random mixup. For example, with 4 labels per class on CIFAR-100 and STL-10, CalibrateMix reduces ECE by 3.55% and 9.31%, and decreases error rates by 3.52% and 13.56%, respectively, compared to random mixup. Compared to label smoothing, ECE is reduced by 4.62% on CIFAR-100 and 3.21% on STL-10

Metric	Method	CIFAR-10			CIFAR-100			SVHN		STL-10	
		4	25	400	4	25	100	4	100	4	100
ECE	FixMatch	5.42 <sub>0.36</sub>	2.85 <sub>0.23</sub>	2.68 <sub>0.09</sub>	35.88 <sub>1.1</sub>	21.19 <sub>0.62</sub>	15.73 <sub>0.05</sub>	1.02 <sub>0.02</sub>	1.1 <sub>0.04</sub>	32.66 <sub>1.29</sub>	4.22 <sub>0.37</sub>
	FixMatch + CalibrateMix (Ours)	<b>3.24</b> <sub>0.12</sub>	<b>1.97</b> <sub>0.13</sub>	<b>1.25</b> <sub>0.27</sub>	<b>35.43</b> <sub>0.05</sub>	<b>12.68</b> <sub>0.42</sub>	<b>8.74</b> <sub>0.6</sub>	1.06 <sub>0.07</sub>	<b>0.82</b> <sub>0.13</sub>	<b>29.93</b> <sub>0.88</sub>	<b>3.9</b> <sub>0.25</sub>
ECE	FlexMatch	4.76 <sub>0.95</sub>	2.92 <sub>0.53</sub>	2.98 <sub>0.55</sub>	31.51 <sub>0.55</sub>	29.54 <sub>7.24</sub>	18.55 <sub>4.15</sub>	10.15 <sub>2.88</sub>	7.81 <sub>1.55</sub>	35.62 <sub>8.57</sub>	9.45 <sub>2.12</sub>
	FlexMatch + CalibrateMix (Ours)	<b>4.41</b> <sub>0.22</sub>	2.94 <sub>0.15</sub>	<b>2.95</b> <sub>0.24</sub>	<b>27.65</b> <sub>0.33</sub>	<b>28.23</b> <sub>0.27</sub>	<b>18.02</b> <sub>0.54</sub>	<b>10.01</b> <sub>0.08</sub>	<b>7.54</b> <sub>0.11</sub>	<b>32.45</b> <sub>0.23</sub>	<b>9.26</b> <sub>0.41</sub>
ECE	SoftMatch	3.02 <sub>0.21</sub>	2.95 <sub>0.53</sub>	2.21 <sub>0.03</sub>	26.4 <sub>1.2</sub>	19.58 <sub>0.1</sub>	15.34 <sub>0.32</sub>	1.44 <sub>0.03</sub>	1.17 <sub>0.07</sub>	13.24 <sub>0.06</sub>	4.31 <sub>0.4</sub>
	SoftMatch + CalibrateMix (Ours)	<b>1.69</b> <sub>0.15</sub>	<b>2.09</b> <sub>0.18</sub>	<b>1.40</b> <sub>0.26</sub>	<b>23.1</b> <sub>0.3</sub>	<b>13.9</b> <sub>0.1</sub>	<b>11.9</b> <sub>0.33</sub>	<b>0.70</b> <sub>0.05</sub>	<b>0.80</b> <sub>0.03</sub>	<b>12.4</b> <sub>0.08</sub>	<b>3.72</b> <sub>0.6</sub>
ECE	FullySupervised		2.22 <sub>0.13</sub>			6.77 <sub>0.67</sub>		0.4 <sub>0.04</sub>		24.14 <sub>0.38</sub>	
	FullySupervised + CalibrateMix (Ours)		<b>1.84</b> <sub>0.22</sub>			<b>4.74</b> <sub>0.13</sub>		<b>0.29</b> <sub>0.11</sub>		<b>13.77</b> <sub>0.23</sub>	
Error Rate	FixMatch	7.29 <sub>0.05</sub>	<b>4.91</b> <sub>0.02</sub>	4.3 <sub>0.02</sub>	44.45 <sub>0.15</sub>	29.88 <sub>0.17</sub>	22.88 <sub>0.03</sub>	3.65 <sub>0.23</sub>	<b>2.04</b> <sub>0.1</sub>	<b>36.34</b> <sub>0.4</sub>	6.2 <sub>0.07</sub>
	FixMatch + CalibrateMix (Ours)	<b>7.13</b> <sub>0.02</sub>	4.97 <sub>0.01</sub>	<b>4.26</b> <sub>0.08</sub>	<b>44.05</b> <sub>0.11</sub>	<b>28.34</b> <sub>0.41</sub>	<b>22.3</b> <sub>0.3</sub>	<b>3.63</b> <sub>0.11</sub>	2.17 <sub>0.03</sub>	37.3 <sub>0.12</sub>	<b>5.78</b> <sub>0.06</sub>
Error Rate	FlexMatch	<b>5.03</b> <sub>0.08</sub>	<b>4.98</b> <sub>0.07</sub>	4.28 <sub>0.02</sub>	40.15 <sub>1.87</sub>	27.73 <sub>0.35</sub>	21.93 <sub>0.33</sub>	7.88 <sub>1.23</sub>	6.78 <sub>1.15</sub>	29.78 <sub>4.01</sub>	6.29 <sub>0.54</sub>
	FlexMatch + CalibrateMix (Ours)	5.04 <sub>0.05</sub>	4.99 <sub>0.02</sub>	<b>4.23</b> <sub>0.05</sub>	<b>40.02</b> <sub>0.16</sub>	<b>26.33</b> <sub>0.31</sub>	<b>21.55</b> <sub>0.16</sub>	<b>7.56</b> <sub>0.14</sub>	<b>6.71</b> <sub>0.08</sub>	<b>29.03</b> <sub>0.43</sub>	<b>6.25</b> <sub>0.09</sub>
Error Rate	SoftMatch	5.09 <sub>0.07</sub>	<b>4.90</b> <sub>0.02</sub>	4.16 <sub>0.07</sub>	37.11 <sub>0.3</sub>	<b>26.76</b> <sub>0.1</sub>	<b>22.11</b> <sub>0.3</sub>	2.59 <sub>0.8</sub>	2.09 <sub>0.02</sub>	20.90 <sub>2.47</sub>	6.10 <sub>0.06</sub>
	SoftMatch + CalibrateMix (Ours)	<b>5.03</b> <sub>0.06</sub>	4.93 <sub>0.03</sub>	<b>4.05</b> <sub>0.04</sub>	<b>36.68</b> <sub>0.18</sub>	26.86 <sub>0.13</sub>	22.21 <sub>0.03</sub>	<b>2.58</b> <sub>0.01</sub>	<b>2.08</b> <sub>0.2</sub>	<b>15.18</b> <sub>0.63</sub>	<b>6.07</b> <sub>0.04</sub>
Error Rate	FullySupervised		4.63 <sub>0.03</sub>			19.42 <sub>0.17</sub>		2.19 <sub>0.01</sub>		34.41 <sub>0.17</sub>	
	FullySupervised + CalibrateMix (Ours)		<b>4.51</b> <sub>0.08</sub>			<b>18.97</b> <sub>0.26</sub>		<b>2.17</b> <sub>0.04</sub>		<b>27.81</b> <sub>0.21</sub>	

Table 1: Expected Calibration Errors (ECE, %) and Top-1 error rates (%) on CIFAR-10, CIFAR-100, SVHN, and STL-10 datasets by FixMatch, FlexMatch, SoftMatch and CalibrateMix (the lower the better). Values are reported in the format  $X_Y$  where  $X$  is the mean and  $Y$  is the Standard Deviation across 3 runs. Better scores in comparison are shown in **bold**.

on the 4 labels per class setting. These consistent reductions across different models and datasets highlight CalibrateMix as a strong solution to existing calibration techniques, making it an effective and versatile regularization method for SSL.

### Performance on ImageNet and WebVision

We report the performance of CalibrateMix on two large-scale datasets: ImageNet (Deng et al. 2009) and WebVision (Li et al. 2017). We randomly sample 10% examples from the training set to be used as labeled samples and use the rest of the examples as unlabeled data. We show the results of FixMatch, FlexMatch, SoftMatch and FullySupervised in terms of ECE and Error rate in Table 3. We note that CalibrateMix considerably boosts the calibration of our models in all setups and yields small error rate improvements in most settings. Specifically, CalibrateMix outperforms SoftMatch by 1.79% ECE on WebVision and pushes the performance over FixMatch by 0.99% ECE on ImageNet. Additionally, CalibrateMix reduces ECE by 1.98% over FlexMatch on the WebVision dataset. When applied to the FullySupervised setting, it also reduces ECE by 1.50% compared to the base supervised model on WebVision. These results highlight the effectiveness of our approach in large-scale, real-world scenarios, leading to better-calibrated models.

### Ablation Study

We conduct experiments to capture: (1) the effect of the targeted mixup on different setups such as how to mix labeled and unlabeled samples as well as with and without warmup, and (2) mixup with and without considering the dissimilar samples (without considering cosine similarity).

### Effect of Different Mixup Strategies With and Without Warmup

We run this ablation on CIFAR-100 and STL-10 (4 labels per class) using SoftMatch, with results reported in Table 4. We denote “easy-to-learn” labeled samples as  $LE$ , “hard-to-learn”

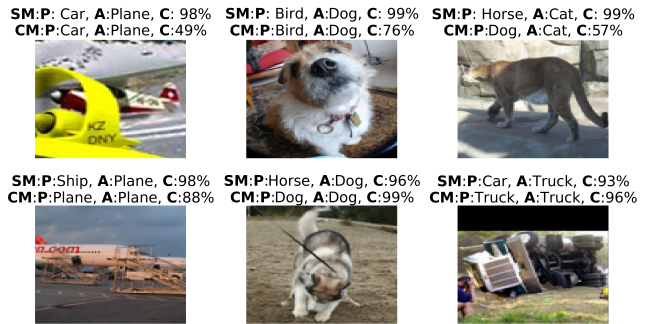


Figure 2: Confidence comparison on images from the STL-10 dataset for **SoftMatch (SM)** and **CalibrateMix (CM)**. For each image, **P** denotes the predicted class, **A** the ground-truth (actual) class, and **C** the associated prediction confidence.

labeled samples as  $LH$ , and the corresponding unlabeled categories as  $UE$  and  $UH$ . Our first ablation compares CalibrateMix, which mixes only across labeled and unlabeled easy/hard sets ( $LE+UH$ ,  $LH+UE$ ), to a mixup strategy that mixes both across and within labeled/unlabeled sets ( $LE+LH$ ,  $LE+UH$ ,  $LH+UE$ ,  $UE+UH$ ). We evaluate this “mixup all” variant with and without warmup.

As shown in the table, CalibrateMix yields lower ECE and error rates than mixup all +warmup. Moreover, removing warmup significantly increases error rates on both CIFAR-100 and STL-10, despite producing slightly lower ECE on CIFAR-100.

### Effect of Selecting Dissimilar Samples

We evaluate our choice of using the  $top-k$  most dissimilar samples for mixup by experimenting SoftMatch + CalibrateMix with different values of  $k$  (5%, 10%, 15%). Results on CIFAR-100 and STL-10 (Table 4) show that  $k = 5\%$  with cosine similarity yields the best performance. Removing cosine sim-

Metric	Method	CIFAR-10			CIFAR-100			SVHN		STL-10	
		4	25	400	4	25	100	4	100	4	100
ECE	FixMatch + RandomMixup	5.26 <sub>0.61</sub>	3.19 <sub>1.42</sub>	1.72 <sub>0.11</sub>	38.05 <sub>2.51</sub>	14.78 <sub>0.21</sub>	8.72 <sub>1.91</sub>	2.87 <sub>0.25</sub>	0.87 <sub>0.06</sub>	35.58 <sub>3.61</sub>	4.16 <sub>0.30</sub>
	FixMatch + CalibrateMix (Ours)	<b>3.24</b> <sub>0.12</sub>	<b>1.97</b> <sub>0.13</sub>	<b>1.25</b> <sub>0.27</sub>	<b>35.43</b> <sub>0.05</sub>	<b>12.68</b> <sub>0.42</sub>	8.74 <sub>0.6</sub>	<b>1.06</b> <sub>0.07</sub>	<b>0.82</b> <sub>0.13</sub>	<b>29.93</b> <sub>0.88</sub>	<b>3.9</b> <sub>0.25</sub>
	FixMatch + LS	4.07 <sub>0.48</sub>	2.52 <sub>0.22</sub>	2.69 <sub>0.11</sub>	34.11 <sub>2.01</sub>	17.37 <sub>1.11</sub>	7.56 <sub>0.95</sub>	1.23 <sub>0.05</sub>	1.19 <sub>0.06</sub>	32.96 <sub>1.73</sub>	4.02 <sub>0.32</sub>
	FixMatch + CalibrateMix (Ours) + LS	<b>4.07</b> <sub>0.41</sub>	<b>2.06</b> <sub>0.23</sub>	<b>1.22</b> <sub>0.45</sub>	<b>31.71</b> <sub>1.33</sub>	<b>12.89</b> <sub>1.1</sub>	<b>7.11</b> <sub>0.8</sub>	<b>1.22</b> <sub>0.10</sub>	<b>0.99</b> <sub>0.13</sub>	<b>29.56</b> <sub>2.34</sub>	<b>3.9</b> <sub>0.28</sub>
ECE	FlexMatch + RandomMixup	4.52 <sub>0.18</sub>	3.57 <sub>0.31</sub>	3.24 <sub>0.02</sub>	27.80 <sub>0.08</sub>	<b>26.72</b> <sub>1.05</sub>	21.79 <sub>0.19</sub>	10.45 <sub>0.34</sub>	7.80 <sub>0.76</sub>	<b>31.80</b> <sub>1.38</sub>	9.35 <sub>0.23</sub>
	FlexMatch + CalibrateMix (Ours)	<b>4.41</b> <sub>0.22</sub>	<b>2.94</b> <sub>0.15</sub>	<b>2.95</b> <sub>0.24</sub>	<b>27.65</b> <sub>0.33</sub>	28.23 <sub>0.27</sub>	<b>18.02</b> <sub>0.54</sub>	<b>10.01</b> <sub>0.08</sub>	<b>7.54</b> <sub>0.11</sub>	32.45 <sub>0.23</sub>	<b>9.26</b> <sub>0.41</sub>
	FlexMatch + LS	4.01 <sub>0.32</sub>	<b>2.23</b> <sub>0.21</sub>	2.96 <sub>0.31</sub>	29.41 <sub>5.71</sub>	29.03 <sub>3.17</sub>	18.17 <sub>3.61</sub>	<b>9.65</b> <sub>1.16</sub>	7.23 <sub>1.34</sub>	31.87 <sub>4.5</sub>	7.86 <sub>1.54</sub>
	FlexMatch + CalibrateMix (Ours) + LS	<b>3.99</b> <sub>0.32</sub>	2.24 <sub>0.19</sub>	<b>2.91</b> <sub>0.14</sub>	<b>26.84</b> <sub>1.34</sub>	<b>27.53</b> <sub>0.105</sub>	<b>17.03</b> <sub>0.61</sub>	9.67 <sub>0.15</sub>	<b>7.01</b> <sub>0.16</sub>	<b>30.44</b> <sub>1.73</sub>	<b>7.1</b> <sub>0.34</sub>
ECE	SoftMatch + RandomMixup	2.75 <sub>0.64</sub>	2.17 <sub>0.10</sub>	1.84 <sub>0.04</sub>	26.65 <sub>1.12</sub>	14.12 <sub>0.32</sub>	<b>11.62</b> <sub>0.58</sub>	2.44 <sub>0.28</sub>	0.84 <sub>0.05</sub>	21.71 <sub>3.47</sub>	4.07 <sub>0.05</sub>
	SoftMatch + CalibrateMix (Ours)	<b>1.69</b> <sub>0.15</sub>	<b>2.09</b> <sub>0.18</sub>	<b>1.4</b> <sub>0.26</sub>	<b>23.1</b> <sub>0.3</sub>	<b>13.9</b> <sub>0.1</sub>	11.9 <sub>0.33</sub>	<b>0.7</b> <sub>0.05</sub>	<b>0.8</b> <sub>0.03</sub>	<b>12.4</b> <sub>0.08</sub>	<b>3.72</b> <sub>0.6</sub>
	SoftMatch + LS	6.62 <sub>0.32</sub>	6.83 <sub>0.3</sub>	7.02 <sub>0.08</sub>	20.89 <sub>0.40</sub>	7.88 <sub>0.37</sub>	6.11 <sub>0.35</sub>	8.31 <sub>0.34</sub>	7.80 <sub>0.23</sub>	11.26 <sub>3.44</sub>	7.46 <sub>0.5</sub>
	SoftMatch + CalibrateMix (Ours) + LS	<b>4.76</b> <sub>0.39</sub>	<b>5.70</b> <sub>0.37</sub>	<b>6.49</b> <sub>0.11</sub>	<b>16.27</b> <sub>0.44</sub>	<b>6.25</b> <sub>0.59</sub>	<b>4.49</b> <sub>0.19</sub>	<b>4.64</b> <sub>1.49</sub>	<b>7.32</b> <sub>0.06</sub>	<b>8.05</b> <sub>2.18</sub>	<b>7.01</b> <sub>0.25</sub>
Error Rate	FixMatch + RandomMixup	8.98 <sub>1.85</sub>	6.89 <sub>1.33</sub>	4.35 <sub>0.04</sub>	50.53 <sub>2.19</sub>	<b>27.17</b> <sub>0.33</sub>	<b>20.76</b> <sub>0.23</sub>	<b>2.48</b> <sub>0.04</sub>	2.40 <sub>0.13</sub>	40.69 <sub>5.30</sub>	6.44 <sub>0.11</sub>
	FixMatch + CalibrateMix (Ours)	<b>7.13</b> <sub>0.02</sub>	<b>4.97</b> <sub>0.01</sub>	<b>4.26</b> <sub>0.08</sub>	<b>44.05</b> <sub>0.11</sub>	28.34 <sub>0.41</sub>	22.3 <sub>0.3</sub>	3.63 <sub>0.11</sub>	<b>2.17</b> <sub>0.03</sub>	<b>37.3</b> <sub>0.12</sub>	<b>5.78</b> <sub>0.06</sub>
	FixMatch + LS	7.41 <sub>0.61</sub>	5.06 <sub>0.98</sub>	4.32 <sub>0.17</sub>	44.57 <sub>4.16</sub>	29.82 <sub>2.64</sub>	23.62 <sub>1.44</sub>	3.71 <sub>0.46</sub>	2.11 <sub>0.24</sub>	<b>36.43</b> <sub>2.32</sub>	6.56 <sub>0.28</sub>
	FixMatch + CalibrateMix (Ours) + LS	<b>7.13</b> <sub>0.61</sub>	<b>4.99</b> <sub>0.23</sub>	<b>4.27</b> <sub>0.25</sub>	<b>44.11</b> <sub>1.42</sub>	<b>28.64</b> <sub>1.24</sub>	<b>22.52</b> <sub>1.05</sub>	<b>3.64</b> <sub>0.31</sub>	<b>2.11</b> <sub>0.25</sub>	37.54 <sub>1.76</sub>	<b>5.72</b> <sub>0.42</sub>
Error Rate	FlexMatch + RandomMixup	5.84 <sub>0.34</sub>	5.45 <sub>0.08</sub>	4.69 <sub>0.05</sub>	40.62 <sub>0.88</sub>	26.53 <sub>0.43</sub>	22.12 <sub>0.13</sub>	7.78 <sub>0.45</sub>	6.96 <sub>0.36</sub>	31.60 <sub>6.21</sub>	6.66 <sub>0.13</sub>
	FlexMatch + CalibrateMix (Ours)	<b>5.04</b> <sub>0.05</sub>	<b>4.99</b> <sub>0.02</sub>	<b>4.23</b> <sub>0.05</sub>	<b>40.02</b> <sub>0.16</sub>	<b>26.33</b> <sub>0.31</sub>	<b>21.55</b> <sub>0.16</sub>	<b>7.56</b> <sub>0.14</sub>	<b>6.71</b> <sub>0.08</sub>	<b>29.03</b> <sub>0.43</sub>	<b>6.25</b> <sub>0.09</sub>
	FlexMatch + LS	<b>4.81</b> <sub>0.51</sub>	4.72 <sub>0.47</sub>	<b>4.01</b> <sub>0.34</sub>	40.11 <sub>1.64</sub>	27.56 <sub>1.15</sub>	21.53 <sub>1.04</sub>	7.51 <sub>0.31</sub>	6.91 <sub>0.37</sub>	29.44 <sub>1.27</sub>	<b>5.88</b> <sub>0.49</sub>
	FlexMatch + CalibrateMix (Ours) + LS	4.83 <sub>0.41</sub>	<b>4.71</b> <sub>0.37</sub>	4.05 <sub>0.22</sub>	<b>40.01</b> <sub>1.87</sub>	<b>26.02</b> <sub>1.46</sub>	<b>21.35</b> <sub>1.21</sub>	<b>7.45</b> <sub>0.48</sub>	<b>6.63</b> <sub>0.32</sub>	<b>28.66</b> <sub>0.28</sub>	6.04 <sub>0.16</sub>
Error Rate	SoftMatch + RandomMixup	7.44 <sub>0.96</sub>	6.02 <sub>0.36</sub>	4.42 <sub>0.08</sub>	40.20 <sub>0.6</sub>	<b>25.42</b> <sub>0.15</sub>	<b>21.58</b> <sub>0.86</sub>	<b>2.51</b> <sub>0.07</sub>	2.64 <sub>0.18</sub>	28.74 <sub>6.26</sub>	<b>5.78</b> <sub>0.14</sub>
	SoftMatch + CalibrateMix (Ours)	<b>5.03</b> <sub>0.05</sub>	<b>4.93</b> <sub>0.03</sub>	<b>4.05</b> <sub>0.04</sub>	<b>36.68</b> <sub>0.18</sub>	26.86 <sub>0.13</sub>	22.21 <sub>0.03</sub>	2.58 <sub>0.01</sub>	<b>2.08</b> <sub>0.2</sub>	<b>15.18</b> <sub>0.63</sub>	6.07 <sub>0.04</sub>
	SoftMatch + LS	5.56 <sub>0.46</sub>	5.19 <sub>0.12</sub>	4.46 <sub>0.09</sub>	39.10 <sub>0.46</sub>	<b>26.50</b> <sub>0.25</sub>	22.13 <sub>0.11</sub>	3.14 <sub>0.27</sub>	3.07 <sub>0.20</sub>	24.57 <sub>3.93</sub>	<b>5.63</b> <sub>0.08</sub>
	SoftMatch + CalibrateMix (Ours) + LS	<b>5.49</b> <sub>0.34</sub>	<b>5.09</b> <sub>0.14</sub>	<b>4.43</b> <sub>0.11</sub>	<b>37.61</b> <sub>1</sub>	27.41 <sub>1.18</sub>	<b>21.81</b> <sub>0.28</sub>	<b>2.9</b> <sub>0.18</sub>	<b>2.84</b> <sub>0.07</sub>	<b>19.24</b> <sub>5.76</sub>	5.97 <sub>0.12</sub>

Table 2: Expected Calibration Errors (ECE, %) and Top-1 error rates (%) on CIFAR-10, CIFAR-100, SVHN, and STL-10 datasets by Random Mixup, Label Smoothing (LS) and CalibrateMix (the lower the better). Values are reported in the format  $X_Y$  where  $X$  is the mean and  $Y$  is the Standard Deviation across 3 runs. Better scores in comparison are shown in **bold**.

Dataset	ImageNet		WebVision	
	ECE	Error	ECE	Error
Fixmatch	9.44	<b>43.54</b>	12.5	44.51
Fixmatch + Ours	<b>8.45</b>	43.61	<b>10.4</b>	<b>44.47</b>
Flexmatch	10.45	<b>42.9</b>	13.32	43.7
Flexmatch + Ours	<b>9.76</b>	43.01	<b>11.34</b>	<b>43.66</b>
SoftMatch	9.88	40.05	10.8	42.01
SoftMatch + Ours	<b>7.56</b>	<b>39.88</b>	<b>9.01</b>	<b>41.77</b>
FullySupervised	6.55	<b>21.06</b>	7.57	27.54
FullySupervised + Ours	<b>5.76</b>	21.87	<b>6.07</b>	<b>27.22</b>

Table 3: ECE and Top-1 error rates on ImageNet and WebVision (the lower the better). Better results in comparison are shown in **bold**.

ilarity consistently degrades calibration and accuracy, with large increases in ECE and error rate on STL-10.

## Error Analysis

Figure 2 presents the confidence, predicted labels, and ground-truth labels for six STL-10 test images under the 4-labels-per-class setting. In the top row, both CalibrateMix and SoftMatch misclassify the images. However, SoftMatch does so with extremely high confidence (over 98%), whereas CalibrateMix assigns noticeably lower confidence, reflecting appropriate uncertainty and thus better calibration. In the bottom row, SoftMatch again produces overconfident errors, while CalibrateMix makes correct predictions with high confidence, indicating better reliability. Overall, these cases show that CalibrateMix reduces overconfident mistakes while still maintaining strong confidence on correct predictions.

Method	CIFAR-100		STL-10	
	ECE	Err	ECE	Err
SoftMatch	26.4 <sub>1.2</sub>	37.11 <sub>0.3</sub>	13.24 <sub>0.06</sub>	20.90 <sub>2.47</sub>
+ mixup all -warmup	16.5 <sub>3.35</sub>	86.95 <sub>1.05</sub>	18.77 <sub>2.21</sub>	23.41 <sub>2.42</sub>
+ mixup all +warmup	30.63 <sub>0.64</sub>	41.5 <sub>0.76</sub>	16.66 <sub>3.35</sub>	19.70 <sub>5.07</sub>
+ CM k=0 (no cosine)	23.25 <sub>0.32</sub>	37.13 <sub>0.35</sub>	16.81 <sub>4.01</sub>	15.69 <sub>2.78</sub>
+ CM k=10	24.42 <sub>1.04</sub>	37.8 <sub>1.3</sub>	12.71 <sub>1.02</sub>	14.92 <sub>1.03</sub>
+ CM k=15	23.14 <sub>0.04</sub>	36.92 <sub>0.02</sub>	14.66 <sub>0.71</sub>	15.38 <sub>0.85</sub>
+ CM k=5 (Ours)	23.1 <sub>0.3</sub>	36.68 <sub>0.18</sub>	12.4 <sub>0.08</sub>	15.18 <sub>0.63</sub>

Table 4: Ablation on CIFAR-100 and STL-10 (4 labels/class). CalibrateMix has warmup and cosine similarity. CM stands for CalibrateMix.

## Conclusion

We proposed CalibrateMix, a targeted mixup strategy that improves the calibration of SSL models by utilizing training dynamics and dissimilarity-aware pairing of easy and hard samples. Specifically, CalibrateMix leverages Area Under the Margin (AUM) and Average Pseudo Margin (APM) to identify sample difficulty, and then performs mixup between labeled and pseudo-labeled samples based on their difficulty and feature dissimilarity. Our method consistently reduces ECE and generally improves accuracy (yielding lower error rates), especially in low-label settings. This enhances the reliability of AI systems in high-stakes applications while also lowering data annotation costs. CalibrateMix contributes to positive societal impact, enabling safer deployments through better-calibrated predictions. For future work, we aim to extend CalibrateMix to other tasks such as object detection, semantic segmentation, and out-of-domain (OOD) scenarios.

## Acknowledgements

This research is supported in part by the NSF IIS award 2107518, a UIC Discovery Partners Institute (DPI) award, and a Google CAHSI award. Any opinions, findings, and conclusions expressed here are those of the authors and do not necessarily reflect the views of NSF, DPI, or Google CAHSI. We thank our anonymous reviewers for their constructive feedback, which helped improve the quality of our paper.

## References

- Bartlett, P. L.; Foster, D. J.; and Telgarsky, M. J. 2017. Spectrally-normalized margin bounds for neural networks. *Advances in neural information processing systems*, 30.
- Chapelle, O.; Scholkopf, B.; and Zien, A. 2009. Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews]. *IEEE Transactions on Neural Networks*, 20(3): 542–542.
- Chen, H.; Tao, R.; Fan, Y.; Wang, Y.; Wang, J.; Schiele, B.; Xie, X.; Raj, B.; and Savvides, M. 2023. Softmatch: Addressing the quantity-quality trade-off in semi-supervised learning. *arXiv preprint arXiv:2301.10921*.
- Coates, A.; Ng, A.; and Lee, H. 2011. An Analysis of Single-Layer Networks in Unsupervised Feature Learning. In Gordon, G.; Dunson, D.; and Dudík, M., eds., *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, 215–223. Fort Lauderdale, FL, USA: PMLR.
- Cubuk, E. D.; Zoph, B.; Shlens, J.; and Le, Q. V. 2020. Randaugment: Practical automated data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 702–703.
- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.
- Elsayed, G.; Krishnan, D.; Mobahi, H.; Regan, K.; and Bengio, S. 2018. Large margin deep networks for classification. *Advances in neural information processing systems*, 31.
- Gan, K.; and Wei, T. 2024. Erasing the Bias: Fine-Tuning Foundation Models for Semi-Supervised Learning. *arXiv preprint arXiv:2405.11756*.
- Guo, C.; Pleiss, G.; Sun, Y.; and Weinberger, K. Q. 2017. On calibration of modern neural networks. In *International conference on machine learning*, 1321–1330. PMLR.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- Krizhevsky, A.; Hinton, G.; et al. 2009. Learning multiple layers of features from tiny images.
- LeCun, Y.; Bengio, Y.; and Hinton, G. 2015. Deep learning. *nature*, 521(7553): 436–444.
- Lee, D.-H.; et al. 2013. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, 896. Atlanta.
- Li, W.; Wang, L.; Li, W.; Agustsson, E.; and Van Gool, L. 2017. Webvision database: Visual learning and understanding from web data. *arXiv preprint arXiv:1708.02862*.
- Lu, D.; and Weng, Q. 2007. A survey of image classification methods and techniques for improving classification performance. *International journal of Remote sensing*, 28(5): 823–870.
- Mathew, A.; Amudha, P.; and Sivakumari, S. 2021. Deep learning techniques: an overview. *Advanced Machine Learning Technologies and Applications: Proceedings of AMLTA 2020*, 599–608.
- Minaee, S.; Boykov, Y.; Porikli, F.; Plaza, A.; Kehtarnavaz, N.; and Terzopoulos, D. 2021. Image segmentation using deep learning: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(7): 3523–3542.
- Minderer, M.; Djolonga, J.; Romijnders, R.; Hubis, F.; Zhai, X.; Houlsby, N.; Tran, D.; and Lucic, M. 2021. Revisiting the calibration of modern neural networks. *Advances in neural information processing systems*, 34: 15682–15694.
- Mukhoti, J.; Kulharia, V.; Sanyal, A.; Golodetz, S.; Torr, P.; and Dokania, P. 2020. Calibrating deep neural networks using focal loss. *Advances in neural information processing systems*, 33: 15288–15299.
- Müller, R.; Kornblith, S.; and Hinton, G. E. 2019. When does label smoothing help? *Advances in neural information processing systems*, 32.
- Netzer, Y.; Wang, T.; Coates, A.; Bissacco, A.; Wu, B.; Ng, A. Y.; et al. 2011. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, 4. Granada.
- Nguyen, K.-B. 2024. Sequencematch: Revisiting the design of weak-strong augmentations for semi-supervised learning. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 96–106.
- Papageorgiou, C. P.; Oren, M.; and Poggio, T. 1998. A general framework for object detection. In *Sixth international conference on computer vision (IEEE Cat. No. 98CH36271)*, 555–562. IEEE.
- Park, S. Y.; and Caragea, C. 2022. On the calibration of pre-trained language models using mixup guided by area under the margin and saliency. *arXiv preprint arXiv:2203.07559*.
- Pinto, F.; Yang, H.; Lim, S. N.; Torr, P.; and Dokania, P. 2022. Using mixup as a regularizer can surprisingly improve accuracy & out-of-distribution robustness. *Advances in Neural Information Processing Systems*, 35: 14608–14622.
- Pleiss, G.; Zhang, T.; Elenberg, E.; and Weinberger, K. Q. 2020. Identifying mislabeled data using the area under the margin ranking. *Advances in Neural Information Processing Systems*, 33: 17044–17056.
- Sohn, K.; Berthelot, D.; Carlini, N.; Zhang, Z.; Zhang, H.; Raffel, C. A.; Cubuk, E. D.; Kurakin, A.; and Li, C.-L. 2020.

Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33: 596–608.

Sosea, T.; and Caragea, C. 2023. MarginMatch: Improving semi-supervised learning with pseudo-margins. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15773–15782.

Srivastava, N.; Hinton, G.; Krizhevsky, A.; Sutskever, I.; and Salakhutdinov, R. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1): 1929–1958.

Thulasidasan, S.; Chennupati, G.; Bilmes, J. A.; Bhattacharya, T.; and Michalak, S. 2019. On mixup training: Improved calibration and predictive uncertainty for deep neural networks. *Advances in neural information processing systems*, 32.

Wang, Y.; Chen, H.; Fan, Y.; SUN, W.; Tao, R.; Hou, W.; Wang, R.; Yang, L.; Zhou, Z.; Guo, L.-Z.; Qi, H.; Wu, Z.; Li, Y.-F.; Nakamura, S.; Ye, W.; Savvides, M.; Raj, B.; Shinnozaki, T.; Schiele, B.; Wang, J.; Xie, X.; and Zhang, Y. 2022. USB: A Unified Semi-supervised Learning Benchmark for Classification. In Koyejo, S.; Mohamed, S.; Agarwal, A.; Belgrave, D.; Cho, K.; and Oh, A., eds., *Advances in Neural Information Processing Systems*, volume 35, 3938–3961. Curran Associates, Inc.

Xie, Q.; Luong, M.-T.; Hovy, E.; and Le, Q. V. 2020. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10687–10698.

Zagoruyko, S. 2016. Wide residual networks. *arXiv preprint arXiv:1605.07146*.

Zhang, B.; Wang, Y.; Hou, W.; Wu, H.; Wang, J.; Okumura, M.; and Shinnozaki, T. 2021. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information Processing Systems*, 34: 18408–18419.

Zhang, H.; Cissé, M.; Dauphin, Y. N.; and Lopez-Paz, D. 2018. mixup: Beyond Empirical Risk Minimization. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.