

# PabLO: Improving Semi-Supervised Learning with Pseudolabeling Optimization

Harit Vishwakarma

HVISHWAKARMA@CS.WISC.EDU

Yi Chen\*

YI.CHEN@WISC.EDU

Satya Sai Srinath Namburi\*

SGNAMBURI@CS.WISC.EDU

Sui Jiet Tay

SSTAY2@WISC.EDU

Ramya Korlakai Vinayak

RAMYA@ECE.WISC.EDU

Frederic Sala

FREDSALA@CS.WISC.EDU

## Abstract

Modern semi-supervised learning (SSL) methods frequently rely on pseudolabeling and consistency regularization. The main technical challenge in pseudolabeling is identifying the points that can reliably be labeled. Existing methods use ad-hoc or hand-crafted notions of confidence and threshold selection functions to choose points. Though such hand-designed strategies shine on benchmark datasets, they may not fare well in specialized settings. To address this challenge we propose a framework to learn confidence functions and thresholds explicitly aligned with the SSL task, obviating the need for manual designs. Our approach formulates an optimization problem over a flexible space of confidence functions and thresholds, allowing us to obtain optimal scoring functions—while remaining compatible with the most popular and performant SSL techniques today. Extensive empirical evaluation of our method shows up to 11% improvement in test accuracy over the standard baselines while requiring substantially fewer training iterations.

## 1 Introduction

Obtaining high-quality labeled data is a major bottleneck in machine learning. The semi-supervised learning (SSL) paradigm tackles this problem by training models on a small amount of labeled data and a large quantity of unlabeled data (Chapelle et al., 2006; Zhu, 2005; van Engelen and Hoos, 2019). While SSL as a field dates back decades and includes a wide variety of approaches, modern SSL methods frequently rely on a pair of ideas: pseudolabeling (McLachlan, 1975; Amini et al., 2023; Rosenberg et al., 2005; Lee, 2013; Rizve et al., 2021) and consistency regularization (Laine and Aila, 2017; Bachman et al., 2014; Sajjadi et al., 2016; Fan et al., 2021; Kukačka et al., 2017). SSL techniques marrying these ideas have delivered strong performance on a number of benchmark datasets.

---

\*. Authors contributed equally to this work

The main challenge with pseudolabeling is balancing accurate point selection with efficient model training. Methods often use confidence thresholds to choose points, aiming to avoid cascading errors while ensuring timely convergence. Numerous ad-hoc strategies exist for both confidence calculation and threshold setting. These include maximum softmax probability (MSP) scores for confidence as well as using fixed thresholds (Sohn et al., 2020), class-wise and adaptive thresholds (Zhang et al., 2021; Wang et al., 2023; Xu et al., 2021; Chen et al., 2023), and instance-dependent thresholds (Li et al., 2023) for thresholding. While such hand-designed strategies have been tuned to produce good performance on popular benchmarks, *practitioners may struggle using them for SSL in specialized settings.*

A promising solution to the pseudolabeling challenge is a framework that learns confidence functions and thresholds *explicitly aligned* with the SSL task, eliminating the need for manual experimentation. Inspired by threshold-based auto-labeling (TBAL) (Vishwakarma et al., 2024), a data development technique, we propose a framework that adapts TBAL principles to learn confidence functions and thresholds specifically for pseudolabeling-based SSL.

Our approach involves two aspects. First, we formulate an optimization problem over a flexible space of confidence functions and thresholds to optimize the quantity/quality tradeoff in pseudolabeling. The space we optimize over is broad enough to subsume many existing manually-designed approaches. That is, we *learn confidence functions and thresholds*. Second, we develop strategies to make the framework compatible with SSL approaches.

Experimentally, we couple our framework to some of the most prominent SSL techniques in use today, including Fixmatch (Sohn et al., 2020) and Freematch (Wang et al., 2023). We observed accuracy lifts of up to 11%, 6%, and 3% on popular benchmarks like SVHN, CIFAR-10, and CIFAR-100 respectively, along with substantial improvements in convergence speed. Our main contributions are,

1. We formulate the goal of maximizing pseudolabel quality and quantity as an optimization problem over a flexible space of scoring functions and thresholds, allowing us to characterize the optimal scoring functions.
2. From this framework we derive a practical method to learn scoring functions and thresholds, optimizing for the quality and quantity of pseudolabels. This method can be flexibly integrated with existing SSL methods that rely on pseudolabeling.
3. We provide extensive empirical evaluation and obtain strong results: the combination of our method with baseline SSL approaches outperforms the baselines, achieving **substantial improvements in test accuracy, by up to 11%, and requires substantially fewer iterations.**

## 2 Background and Problem Setup

We begin with notation, then provide useful background and a statement of our goal.

**Notation.** Consider a feature space  $\mathcal{X}$  and label space  $\mathcal{Y} = \{1, \dots, k\}$  in a  $k$ -class classification task. As usual in semi-supervised learning, we have access to a set  $X_u = \{\mathbf{x}_u\}_{u=1}^{n_u}$  of unlabeled data drawn from the distribution  $P_x$  over  $\mathcal{X}$ . We also have access to  $D_l = \{(\mathbf{x}_l, y_l)\}_{l=1}^{N_l}$ ,

a set of labeled data points drawn from the joint distribution  $P_{xy}$ , with  $n_l \ll N_u$ . Let  $h : \mathcal{X} \rightarrow \mathcal{Y}$  denote a model and  $g : \mathcal{X} \rightarrow T^k \subseteq \mathbb{R}^k$  be an associated confidence function giving a score  $g(\mathbf{x})$  indicating the confidence of  $h$  on its prediction for any data point  $\mathbf{x}$ . For any  $\mathbf{x}$  the hard label prediction is  $\hat{y} := h(\mathbf{x})$ . When the prediction  $\hat{y}$  is used as a pseudolabel we denote it as  $\tilde{y}$ . In general, for a vector  $\mathbf{v} \in \mathbb{R}^d$ ,  $\mathbf{v}[i]$  denotes its  $i$ -th component. The vector  $\mathbf{t}$  denotes thresholds over the scores  $k$ -classes, and  $\mathbf{t}[y]$  is its  $y$ -th entry, i.e., the score for class  $y$ .

## 2.1 Pseudolabeling-based Semi-Supervised Learning

Given, as above, a large collection of unlabeled data  $X_u$  and a small set of labeled points  $D_l$ , inductive semi-supervised learning (SSL) seeks to learn a classifier  $\hat{h}_{\text{ssl}}$  from the model class  $\mathcal{H}$ . The promise of SSL is that by effectively using  $X_u$  in the learning process it can learn a better classifier than its supervised counterpart, which learns only from  $D_l$ .

In many recent pseudolabeling-based SSL techniques, in each iteration of training, a batch of labeled and unlabeled data is obtained, then the sum of the losses  $\hat{\mathcal{L}} = \hat{\mathcal{L}}_s + \lambda_u \hat{\mathcal{L}}_u + \lambda_r \hat{\mathcal{L}}_r$  is minimized w.r.t to the model  $h$ . Here  $\hat{\mathcal{L}}_s$  is the supervised loss,  $\hat{\mathcal{L}}_u$  unsupervised loss, and  $\hat{\mathcal{L}}_r$  is (the sum of) regularization term(s). The constants  $\lambda_u, \lambda_r$  are hyperparameters controlling the relative importance of the corresponding terms.

**Supervised loss.** Given a batch of labeled data  $D_l^b$  the supervised loss is computed as follows,  $\hat{\mathcal{L}}_s(h|D_l^b) = \frac{1}{|D_l^b|} \sum_{(x,y) \in D_l^b} H(y, h, \mathbf{x})$ . Here  $H(y, h, \mathbf{x})$  is the standard cross-entropy loss between the 1-hot representation of  $y$  and the softmax output of  $h$  on input  $\mathbf{x}$ .

**Unsupervised loss and consistency regularization.** For the unlabeled batch  $X_u^b$ , pseudolabels  $\tilde{y} = h(\mathbf{x})$  are computed for each  $\mathbf{x} \in X_u^b$ . Then, a pseudolabeling mask  $S(\mathbf{x}, g, \mathbf{t} | h) = \mathbb{1}(g(\mathbf{x})[\tilde{y}] \geq \mathbf{t}[\tilde{y}])$ , is 1 for points having confidence score bigger than predetermined threshold corresponding to the predicted class. Recent methods, couple this loss and consistency regularization together by doing pseudolabeling on weakly augmented data using weak transform  $\omega$  and then defining the cross-entropy loss on the strongly augmented data using strong transformation  $\Omega$ . The loss is

$$\hat{\mathcal{L}}_u(h | g, \mathbf{t}, \tilde{D}_u^b) = \frac{1}{|\tilde{D}_u^b|} \sum_{(x, \tilde{y}) \in \tilde{D}_u^b} S(\omega(\mathbf{x}), g, \mathbf{t} | h) \cdot H(\tilde{y}, h, \Omega(\mathbf{x})).$$

**Regularization.** A regularization term (or a sum of multiple regularizers) is often included along with the above two losses to encourage desired behavior. For instance, Freematch (Wang et al., 2023) adds a self-adaptive class fairness regularizer to encourage diverse predictions during the initial training phase. Similarly, a regularizer is introduced in (Mishra et al., 2024) to encourage calibration in model’s confidence scores. Including such regularizers has been fruitful in pseudolabeling-based SSL.

## 2.2 Problem Statement

The success of pseudolabeling-based SSL hinges heavily on maximizing the quality and quantity of the pseudolabels. These are defined as follows:

**Pseudolabeling coverage (quantity).** Given a set of points  $X$ , the pseudolabeling coverage is the fraction of points that were pseudolabeled using  $h, g$  and  $\mathbf{t}$ . This measurement captures the quantity of pseudolabels and is defined as

$$\hat{\mathcal{P}}(g, \mathbf{t} \mid h, X) := \frac{1}{|X|} \sum_{(\mathbf{x}) \in X} S(\mathbf{x}, g, \mathbf{t} \mid h), \quad \mathcal{P}(g, \mathbf{t} \mid h) := \mathbb{E}_{\mathbf{x}}[S(\mathbf{x}, g, \mathbf{t} \mid h)]. \quad (1)$$

**Pseudolabeling error (quality).** This is the fraction of pseudolabeled points that received wrong labels. This metric captures the quality of pseudolabels:

$$\hat{\mathcal{E}}(g, \mathbf{t} \mid h, D) := \frac{\sum_{(\mathbf{x}, y, \tilde{y}) \in D} S(\mathbf{x}, g, \mathbf{t} \mid h) \cdot \mathbf{1}(h(\mathbf{x}) \neq y)}{\sum_{(\mathbf{x}, y, \tilde{y}) \in D} S(\mathbf{x}, g, \mathbf{t} \mid h)}, \quad (2)$$

$$\mathcal{E}(g, \mathbf{t} \mid h) = \frac{\mathbb{E}_{\mathbf{x}}[S(\mathbf{x}, g, \mathbf{t} \mid h) \cdot \mathbf{1}(h(\mathbf{x}) \neq y)]}{\mathcal{P}(g, \mathbf{t} \mid h)}. \quad (3)$$

**Goal.** In pseudolabeling-based SSL, we want to learn a classifier  $\hat{h}_{\text{ssl}}$  that generalizes well on the unseen data. It is common wisdom in the literature that maximizing the quality and quantity of the pseudolabeled points during the training procedure will lead to a better model. Departing from the hand-crafting strategies to achieve this objective, we seek learnable solutions to confidence scores and thresholding to achieve high coverage at low error, that will eventually lead to a classification model  $\hat{h}_{\text{ssl}}$  with much higher test accuracy.

### 3 Methodology

Our approach is to integrate learnable confidence functions and thresholds into existing pseudolabeling-based SSL pipelines. To do so, we build on a recently-developed technique (Vishwakarma et al., 2024) to improve the performance of threshold-based auto-labeling (TBAL) (SGT, 2022; Vishwakarma et al., 2023; Qiu et al., 2023) systems. In order to make such an approach compatible with SSL, we apply a simple notion—*accumulating pseudolabels*—that may also be useful for other methods.

#### 3.1 Pseudolabeling Optimization Framework

The fundamental problem in pseudolabeling is, given a classifier  $\hat{h}_i$ , to correctly identify the points in the pool of unlabeled data  $X_u$  where the predictions of  $\hat{h}_i$  are correct. Since the classifier is frequently undertrained during the SSL process, it may not have high accuracy. That is, it might only be accurate in some small part of the feature space, which we hope to identify via the confidence scores and appropriate thresholds. As discussed earlier, existing solutions (Lee, 2013; Sohn et al., 2020; Wang et al., 2023) use maximum softmax probability (MSP) from the model  $\hat{h}_i$  in concert with heuristics for thresholds that are either fixed or vary dynamically based on the learning status of the model. Some recent works have observed that MSP scores tend to be miscalibrated and proposed solutions to obtain more calibrated scores (Mishra et al., 2024; Loh et al., 2023), which also led to performance gains.

**Theoretical Framework.** Instead of trying to improve calibration or heuristics for thresholding, we propose to express the objective of pseudolabeling as an optimization problem

over the space of confidence functions and thresholds. The objective is to maximize the quantity i.e. the pseudolabeling coverage (eq. (1)) while keeping the pseudolabeling error low (eq. (3)) i.e. have high quality.

More specifically, one approach to formalizing this optimization problem is to seek to maximize the pseudolabeling coverage while ensuring pseudolabeling error is at most  $\epsilon \in (0, 1)$ , for some hyperparameter  $\epsilon$ . In other words, given the classifier  $\hat{h}_i$  in any iteration  $i$  of SSL, then,

$$g_i^*, \mathbf{t}_i^* \in \arg \max_{g \in \mathcal{G}, \mathbf{t} \in T^k} \mathcal{P}(g, \mathbf{t} | \hat{h}_i) \quad \text{s.t.} \quad \mathcal{E}(g, \mathbf{t} | \hat{h}_i) \leq \epsilon,$$

are the optimal confidence functions and thresholds for pseudolabeling using  $\hat{h}_i$ 's predictions. The *quality* of the pseudolabels can be controlled using  $\epsilon$ . This follows the recipe for TBAL (Vishwakarma et al., 2024), with one additional complication: for SSL, it is not clear what value of  $\epsilon$  is suitable, while in TBAL  $\epsilon$  is a system-level constant provided as input.

The most attractive property of this framework is that, irrespective of the choice of  $\epsilon$ , it provides the scores and threshold that yield maximum pseudolabeling coverage at that error level, freeing us from making arbitrary choices of confidence scores, calibration techniques, and thresholding heuristics. Instead, we solve the optimization problem over a flexible enough space will subsume specific strategies. Next, we discuss how to make the framework tractable.

**Practical Version.** The optimization problem discussed earlier involves population-level quantities which are usually not accessible in practice. Thus we have to fall back to using their finite sample estimates and smooth variations to make the optimization problem tractable. We adapt the steps from (Vishwakarma et al., 2024) to obtain such a practical version of the optimization problem. There, the authors first estimate the coverage and error using a small amount held-out labeled data (called calibration data  $D_{\text{cal}}$ ) curated from the validation data. They then introduce differentiable surrogates for the 0-1 variables. Let  $\sigma(\alpha, z) := 1/(1 + \exp(-\alpha z))$  denote the sigmoid function on  $\mathbb{R}$  with scale parameter  $\alpha \in \mathbb{R}$ . The surrogates are as follows,

$$\tilde{\mathcal{P}}(g, \mathbf{t} | h, D_{\text{cal}}) := \frac{1}{|D_{\text{cal}}|} \sum_{(\mathbf{x}, y, \tilde{y}) \in D_{\text{cal}}} \sigma(\alpha, g(\mathbf{x})[\tilde{y}] - \mathbf{t}[\tilde{y}]), \quad (4)$$

$$\tilde{\mathcal{E}}(g, \mathbf{t} | h, D_{\text{cal}}) := \frac{\sum_{(\mathbf{x}, y, \tilde{y}) \in D_{\text{cal}}} \mathbf{1}(y \neq \tilde{y}) \sigma(\alpha, g(\mathbf{x})[\tilde{y}] - \mathbf{t}[\tilde{y}])}{\sum_{(\mathbf{x}, y, \tilde{y}) \in D_{\text{cal}}} \sigma(\alpha, g(\mathbf{x})[\tilde{y}] - \mathbf{t}[\tilde{y}])}. \quad (5)$$

Using these surrogates the following practical optimization problem is obtained. It is also converted into unconstrained formulation by introducing the penalty term  $\lambda \in \mathbb{R}^+$  controlling the relative importance of the pseudolabeling error and coverage.

$$\hat{g}_i, \hat{\mathbf{t}}_i \in \arg \min_{g \in \mathcal{G}, \mathbf{t} \in T^k} -\tilde{\mathcal{P}}(g, \mathbf{t} | \hat{h}_i, D_{\text{cal}}) + \lambda \tilde{\mathcal{E}}(g, \mathbf{t} | \hat{h}_i, D_{\text{cal}}) \quad (\text{P1})$$

We use 2-layer neural nets as a choice of  $\mathcal{G}$ . The optimization problem (P1) is nonconvex, but differentiable and we solve it using Stochastic Gradient Descent (SGD). See Appendix C for more details on our choice of  $\mathcal{G}$  and training details and hyperparameters.

### 3.2 Threshold Estimation

While we can obtain both the confidence scores and thresholds by solving (P1), we propose to adapt the threshold estimation procedure from (Vishwakarma et al., 2024) as it avoids potential generalization issues due to learning them simultaneously from the same data  $D_{\text{cal}}$  and ensures stricter control over the pseudolabeling errors. It is also decoupled from any particular choice of scoring function, hence it can replace the thresholding procedure in the existing SSL pipelines as well.

Our procedure is simple. It takes in a confidence function  $\tilde{g}_i$  and another part of the held-out validation data referred to as  $D_{\text{th}}$ . It estimates thresholds for each class separately and estimates the pseudolabeling errors  $\hat{\mathcal{E}}(\tilde{g}_i, t \mid h, D_{\text{th}}, \tilde{y})$  on the super level sets of  $\tilde{g}_i$ . Here we slightly abuse notation: instead of  $\mathbf{t} \in T^k$ , we use  $t \in T$ , to indicate the estimate of pseudolabeling error at threshold  $t$  for class  $y$ . To obtain a threshold  $\tilde{\mathbf{t}}[y]$  for class  $y$ , the procedure finds the smallest  $t \in T$  such that  $\hat{\mathcal{E}}(\tilde{g}_i, t \mid h, D_{\text{th}}, \tilde{y}) + C_1 \hat{\sigma}(\hat{\mathcal{E}}) \leq \epsilon$ . Here  $C_1$  is a constant and  $\hat{\sigma}(z) = \sqrt{z \cdot (1 - z)}$  and  $\hat{\mathcal{E}}$  is used for brevity in place of  $\hat{\mathcal{E}}(\tilde{g}_i, t \mid h, D_{\text{th}}, \tilde{y})$ . Using the thresholds found using this procedure ensures pseudolabeling error remains below (or close to) the a tolerance level  $\epsilon$ .

*Remarks.* Departing from fixed thresholds as in (Sohn et al., 2020), prior works have proposed adaptive and class-wise heuristic thresholding schemes based on the model’s learning status, such as in (Djurisic et al., 2023; Zhang et al., 2021; Wang et al., 2023) and others. In contrast, our approach is a principled way to estimate adaptive and class-wise pseudolabeling thresholds while providing strict control over the quality of pseudolabels. This has been used in prior works in TBAL for the problem of creating reliable datasets and is backed by theoretical guarantees for the quality of pseudolabels produced (Vishwakarma et al., 2023).

### 3.3 Pseudolabeling and Accumulation

In the usual pseudolabeling-based SSL setups, the pseudolabels inferred by the model for a mini-batch will be discarded after each iteration. Moreover, it is not guaranteed that a previously pseudolabeled point will get pseudolabeled in the current iteration as well. Given the quality of pseudolabels is high, it is appealing to reuse the past pseudolabel for a point that did not get pseudolabeled in the current iteration. We propose to do so for techniques where the quality of pseudolabels is assured. We refer to this trick as “pseudolabel accumulation”.

Mathematically, if  $s_{i-1}, \tilde{y}_{i-1}$  are the previous mask and psuedolabels and  $s_i, \tilde{y}_{i-1}$  are the current mask and pseudolabels for the same point, then with accumulation,

$$\tilde{y}_i \leftarrow s_i \tilde{y}_i + (1 - s_i) \tilde{y}_{i-1}; \quad \text{and} \quad s_i \leftarrow s_i \vee s_{i-1}.$$

Here  $\vee$  is the boolean `or` operation and the steps are executed in the order. In words, if a point is pseudo-labeled in the current iteration, we use that label. If not, but it was pseudo-labeled previously, we reuse that label. If it’s never been pseudo-labeled, it remains unlabeled.

We refer to our method as `PabL0`. A more formal listing of the steps is detailed in Algorithm 1, deferred to Appendix B due to space constraints.

Table 1: Details of the dataset we use in experiments.  $k$  is the no. of classes.  $N_l$  is the no. of labeled data points used for training the backbone model  $h$ .  $N_u$  is the no. of unlabelled data points used for consistency regularization and pseudolabeling for all the methods.  $N_{\text{val}}$  is the no. of points used for model selection in all methods.  $N_{\text{test}}$  is the no. of test data points.  $N_{\text{cal}}$  is the number of points used for learning the  $g$  function.  $N_{\text{th}}$  is the no. of points used for threshold estimation.

Dataset	Backbone Model	$h$	$k$	$N_u$	$N_{\text{val}}$	$N_{\text{test}}$	$N_l$	$N_{\text{cal}}$	$N_{\text{th}}$	Augmentation
CIFAR-10	WRN-28-2		10	50K	6K	4K	250	1K	1K	Weak, Strong
CIFAR-100	WRN-28-2		100	50K	6K	4K	2500	3K	3K	Weak, Strong
SVHN	WRN-28-2		10	604,388	15,620	10,412	250	3K	3K	Weak, Strong

## 4 Experiments

We evaluate our method empirically to verify the following claims: **C1**. Our method produces models with improved test accuracy while taking fewer iterations. **C2**. In certain cases, we may wish to produce a high-quality dataset using pseudolabeling (rather than a single high-quality model). For such scenarios, PabLO achieves much higher dataset coverage and accuracy. Additionally, we conduct ablation studies, deferred to the Appendix C.

### 4.1 Experimental Setup

First, we provide a brief description of the experimental setup. Details are in Appendix C.

**Methods.** We use two simple base methods capturing the core ideas of pseudolabeling (PL) and consistency regularization (CR). The first is *Fixmatch* (Sohn et al., 2020) which uses fixed thresholds on MSP scores for PL along with CR. *Freematch* (Wang et al., 2023) improves upon it by using adaptive, class-wise thresholds and class fairness regularization (CFR) along with CR, and is a promising method among others using dynamic thresholds for PL. We include their combinations with recently proposed *Bayesian Model Averaging (BAM)* (Loh et al., 2023) and *Margin Regularization (MR)*<sup>1</sup> (Mishra et al., 2024) to improve calibration in SSL. We replace the pseudolabeling component by our method PabLO to obtain *Fixmatch + Ours* (a combination of PabLO and CR) and *Freematch + Ours* (a combination of PabLO, CR, and CFR).

**Datasets.** We experiment with 3 datasets: *CIFAR-10* (Krizhevsky et al., 2009), *CIFAR-100* (Krizhevsky et al., 2009) and *SVHN* (Netzer et al., 2011). More details are summarized in Table 1. We use a portion of the validation data ( $N_{\text{val}}$ ) for our method, split into  $N_{\text{cal}}$ , used to calibrate the function  $g$ , and  $N_{\text{th}}$ , used to estimate the threshold.

**Models and Training.** The backbone encoder is a Wide ResNet-28-2 for all the datasets. We use the default hyperparameters and dataset-specific settings (learning rates, batch size, optimizers and schedulers) following previous baseline recommendations (Wang et al., 2022). We run till 25K iterations—in contrast to the extremely large number of iterations ( $2^{20}$ ) in prior works—which may be unrealistic in practice due to resource constraints. For confidence functions class  $\mathcal{G}$  we use a class of 2-layer neural nets and provide it last two layers representations from  $h$  as input, as in (Vishwakarma et al., 2023). We train it using SGD, the hyperparameters are deferred to Appendix C. We use  $\epsilon = 5\%$  across all settings.

1. We assign this name for convenience.

Table 2: Top-1 Accuracy for CIFAR-10, CIFAR-100 and SVHN averaged across 3 random seeds. The best accuracy is **bolded**

Dataset	CIFAR-10	CIFAR-100	SVHN
# Labels	250	2500	250
Fixmatch	88.15 $\pm$ 1.27	50.07 $\pm$ 1.12	96.54 $\pm$ 0.05
Fixmatch + MR	87.85 $\pm$ 1.10	44.75 $\pm$ 1.36	96.58 $\pm$ 0.04
Fixmatch + BaM	86.44 $\pm$ 1.47	44.58 $\pm$ 0.41	95.99 $\pm$ 0.06
<b>Fixmatch + Ours</b>	<b>93.03 <math>\pm</math> 0.44</b>	<b>53.17 <math>\pm</math> 1.27</b>	<b>96.61 <math>\pm</math> 0.16</b>
Freematch	90.17 $\pm$ 0.13	57.21 $\pm$ 0.78	85.25 $\pm$ 1.70
Freematch + MR	90.17 $\pm$ 0.45	57.23 $\pm$ 1.18	84.65 $\pm$ 1.03
Freematch + BaM	88.34 $\pm$ 0.99	51.98 $\pm$ 1.74	86.28 $\pm$ 1.75
<b>Freematch + Ours</b>	<b>93.08 <math>\pm</math> 0.05</b>	<b>60.96 <math>\pm</math> 0.53</b>	<b>96.48 <math>\pm</math> 0.33</b>

## 4.2 Results and Discussion

To verify our main claims, we compare the baselines, their combinations with our method, and methods that induce calibrated scores in SSL. We run all methods with three random seeds and report the mean and std. deviation of accuracy across three runs in Table 2.

**C1. Test accuracy improvements.** Our method maximizes pseudolabeling coverage and accuracy, producing more accurate pseudolabels. As Table 2 shows, integrating our method into Fixmatch and Freematch significantly improves test accuracy on CIFAR-10, CIFAR-100, and SVHN. Notably, we see a 6% improvement on CIFAR-10 with Fixmatch, a 3% improvement on the harder CIFAR-100 with Fixmatch, and an 11% improvement on SVHN with Freematch.

**C2. Improved pseudolabeling coverage and accuracy.** As our method is designed to maximize coverage and accuracy of pseudolabels, we expect it to maintain high pseudolabeling accuracy and coverage from the beginning. To test this, we log the pseudolabeling coverage and accuracy in each iteration on the batch of unlabeled data used in that iteration. We refer to these as batch pseudolabeling coverage (batch-pl-cov) and batch pseudolabeling accuracy (batch-pl-acc). We show these for CIFAR-10 and CIFAR-100 settings in Figure 1 and 2 in the Appendix. As expected, the batch-pl-acc is high right from the beginning and it is close to the desired level of 95% (with  $\epsilon = 5\%$ ) throughout for CIFAR-10. However, for CIFAR-100 possibly due to high class cardinality it drops to around 70%, This is similar to the baselines but yields much higher coverage. Similar results hold for SVHN ( Figure 3).

## 5 Conclusion

We built a framework, inspired by ideas from autolabeling, that learns confidence functions and thresholds explicitly aligned with the SSL task. This approach eliminates the need for manual designs and hand-crafted notions of confidence, which can be limited in specialized data settings. By formulating an optimization problem over a flexible space of confidence functions and thresholds, we characterized optimal scoring functions. We derived our practical method to learn the scores and evaluated it empirically, where it achieved up to 11% improvement in test accuracy over standard baselines, while also reducing training iterations.



## References

- Ryan Prescott Adams and Zoubin Ghahramani. Archipelago: nonparametric bayesian semi-supervised learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 1–8, 2009.
- Massih-Reza Amini, Vasilii Feofanov, Loic Pauletto, Lies Hadjadj, Emilie Devijver, and Yury Maximov. Self-training: A survey, 2023.
- Eric Arazo, Diego Ortego, Paul Albert, Noel E O’Connor, and Kevin McGuinness. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In *2020 International joint conference on neural networks (IJCNN)*, pages 1–8. IEEE, 2020.
- Philip Bachman, Ouais Alsharif, and Doina Precup. Learning with pseudo-ensembles. In *Advances in Neural Information Processing Systems*, volume 27, 2014.
- Avrim Blum and Shuchi Chawla. Learning from labeled and unlabeled data using graph mincuts. 2001.
- Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *Proc. of the eleventh annual conference on Computational learning theory*, pages 92–100. ACM, 1998.
- Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien, editors. *Semi-Supervised Learning*. The MIT Press, 2006. ISBN 9780262033589.
- Hao Chen, Ran Tao, Yue Fan, Yidong Wang, Jindong Wang, Bernt Schiele, Xing Xie, Bhiksha Raj, and Marios Savvides. Softmatch: Addressing the quantity-quality tradeoff in semi-supervised learning. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=ymt1zQXBdIF>.
- Charles Corbière, Nicolas THOME, Avner Bar-Hen, Matthieu Cord, and Patrick Pérez. Addressing failure prediction by learning model confidence. In *Advances in Neural Information Processing Systems 32*, pages 2902–2913. 2019.
- Andrija Djuricic, Nebojsa Bozanic, Arjun Ashok, and Rosanne Liu. Extremely simple activation shaping for out-of-distribution detection. In *The Eleventh International Conference on Learning Representations*, 2023.
- Yasser El-Manzalawy, Elyse E Munoz, Scott E Lindner, and Vasant Honavar. PlasmoSep: Predicting surface-exposed proteins on the malaria parasite using semisupervised self-training and expert-annotated data. *Proteomics*, 16(23):2967–2976, 2016.
- Yue Fan, Anna Kukleva, and Bernt Schiele. Revisiting consistency regularization for semi-supervised learning, 2021.
- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021.

- 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.
- Chirag Gupta and Aaditya Ramdas. Top-label calibration and multiclass-to-binary reductions. In *International Conference on Learning Representations, 2022*. URL <https://openreview.net/forum?id=WqoBaaPHS->.
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations, 2017*.
- Like Hui, Mikhail Belkin, and Stephen Wright. Cut your losses with squentropy. In *Proceedings of the 40th International Conference on Machine Learning*, pages 14114–14131, 2023.
- Thorsten Joachims. Transductive inference for text classification using support vector machines. In Ivan Bratko and Saso Dzeroski, editors, *Proceedings of ICML-99, 16th International Conference on Machine Learning*, pages 200–209, 1999.
- Jacob Kahn, Ann Lee, and Awni Hannun. Self-training for end-to-end speech recognition. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7084–7088. IEEE, 2020.
- Giannis Karamanolakis, Subhabrata Mukherjee, Guoqing Zheng, and Ahmed Hassan Awadallah. Self-training with weak supervision. *arXiv preprint arXiv:2104.05514*, 2021.
- Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised learning with deep generative models. In *Advances in Neural Information Processing Systems*, volume 27, 2014.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Jan Kukačka, Vladimir Golkov, and Daniel Cremers. Regularization for deep learning: A taxonomy, 2017.
- Meelis Kull, Miquel Perello Nieto, Markus Kängsepp, Telmo Silva Filho, Hao Song, and Peter Flach. Beyond temperature scaling: Obtaining well-calibrated multi-class probabilities with dirichlet calibration. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Ananya Kumar, Percy S Liang, and Tengyu Ma. Verified uncertainty calibration. *Advances in Neural Information Processing Systems*, 32, 2019.
- Aviral Kumar, Sunita Sarawagi, and Ujjwal Jain. Trainable calibration measures for neural networks from kernel mean embeddings. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2805–2814. PMLR, 10–15 Jul 2018.
- Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *Fifth International Conference on Learning Representations, 2017*.

- Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *ICML Workshop on Challenges in Representation Learning*, 2013.
- Muyang Li, Runze Wu, Haoyu Liu, Jun Yu, Xun Yang, Bo Han, and Tongliang Liu. Instant: Semi-supervised learning with instance-dependent thresholds. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Charlotte Loh, Rumen Dangovski, Shivchander Sudalairaj, Seungwook Han, Ligong Han, Leonid Karlinsky, Marin Soljacic, and Akash Srivastava. Mitigating confirmation bias in semi-supervised learning via efficient bayesian model averaging. *Transactions on Machine Learning Research*, 2023.
- Geoffrey J McLachlan. Iterative reclassification procedure for constructing an asymptotically optimal rule of allocation in discriminant analysis. *Journal of the American Statistical Association*, 70(350):365–369, 1975.
- Shambhavi Mishra, Balamurali Murugesan, Ismail Ben Ayed, Marco Pedersoli, and Jose Dolz. Do not trust what you trust: Miscalibration in semi-supervised learning, 2024.
- Jooyoung Moon, Jihyo Kim, Younghak Shin, and Sangheum Hwang. Confidence-aware learning for deep neural networks. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 7034–7044, 2020.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Baolin Wu, Andrew Y Ng, et al. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, page 7. Granada, Spain, 2011.
- Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 427–436, 2015.
- Kamal Nigam, Andrew Kachites McCallum, Sebastian Thrun, and Tom Mitchell. Text classification from labeled and unlabeled documents using em. *Machine learning*, 39: 103–134, 2000.
- Partha Niyogi. Manifold regularization and semi-supervised learning: Some theoretical analyses. *Journal of Machine Learning Research*, 14(5), 2013.
- Avital Oliver, Augustus Odena, Colin A Raffel, Ekin Dogus Cubuk, and Ian Goodfellow. Realistic evaluation of deep semi-supervised learning algorithms. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Samet Oymak and Talha Cihad Gulcu. Statistical and algorithmic insights for semi-supervised learning with self-training. *arXiv preprint arXiv:2006.11006*, 2020.
- Hang Qiu, Krishna Chintalapudi, and Ramesh Govindan. MCAL: Minimum cost human-machine active labeling. In *The Eleventh International Conference on Learning Representations*, 2023.

- Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat, and Mubarak Shah. In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning. In *International Conference on Learning Representations*, 2021.
- Chuck Rosenberg, Martial Hebert, and Henry Schneiderman. Semi-supervised self-training of object detection models. In *Seventh IEEE Workshops on Applications of Computer Vision (WACV/MOTION'05) - Volume 1*, volume 1, pages 29–36, 2005. doi: 10.1109/ACVMOT.2005.107.
- Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, page 1171–1179, 2016.
- Henry Scudder. Probability of error of some adaptive pattern-recognition machines. *IEEE Transactions on Information Theory*, 11(3):363–371, 1965.
- SGT. Aws sagemaker ground truth. <https://aws.amazon.com/sagemaker/data-labeling/>, 2022. Accessed: 2022-11-18.
- Aarti Singh, Robert Nowak, and Jerry Zhu. Unlabeled data: Now it helps, now it doesn't. In *Advances in Neural Information Processing Systems*, volume 21. Curran Associates, Inc., 2008.
- Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608, 2020.
- Amarnag Subramanya and Partha Pratim Talukdar. *Graph-based semi-supervised learning*. Springer Nature, 2022.
- Jesper E. van Engelen and Holger H. Hoos. A survey on semi-supervised learning. *Machine Learning*, 109:373 – 440, 2019.
- Vladimir Naumovich Vapnik, Vlamimir Vapnik, et al. *Statistical learning theory*. 1998.
- Harit Vishwakarma, Huguang Lin, Frederic Sala, and Ramya Korlakai Vinayak. Promises and pitfalls of threshold-based auto-labeling. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Harit Vishwakarma, Sui Jiet Tay, Satya Sai Srinath Namburi, Frederic Sala, Ramya Korlakai Vinayak, et al. Pearls from pebbles: Improved confidence functions for auto-labeling. *arXiv preprint arXiv:2404.16188*, 2024.
- Yidong Wang, Hao Chen, Yue Fan, Wang Sun, Ran Tao, Wenxin Hou, Renjie Wang, Linyi Yang, Zhi Zhou, Lan-Zhe Guo, Heli Qi, Zhen Wu, Yu-Feng Li, Satoshi Nakamura, Wei Ye, Marios Savvides, Bhiksha Raj, Takahiro Shinozaki, Bernt Schiele, Jindong Wang, Xing Xie, and Yue Zhang. Usb: A unified semi-supervised learning benchmark for classification.

In *Thirty-sixth Conference on Neural Information Processing Systems, Datasets and Benchmarks Track*, 2022.

Yidong Wang, Hao Chen, Qiang Heng, Wenxin Hou, Yue Fan, Zhen Wu, Jindong Wang, Marios Savvides, Takahiro Shinozaki, Bhiksha Raj, Bernt Schiele, and Xing Xie. Freematch: Self-adaptive thresholding for semi-supervised learning. In *The Eleventh International Conference on Learning Representations*, 2023. URL [https://openreview.net/forum?id=PDrUPTXJI\\_A](https://openreview.net/forum?id=PDrUPTXJI_A).

Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and Quoc Le. Unsupervised data augmentation for consistency training. In *Advances in Neural Information Processing Systems*, volume 33, 2020.

Yi Xu, Lei Shang, Jinxing Ye, Qi Qian, Yu-Feng Li, Baigui Sun, Hao Li, and Rong Jin. Dash: Semi-supervised learning with dynamic thresholding. In *International Conference on Machine Learning*, pages 11525–11536. PMLR, 2021.

Bianca Zadrozny and Charles Elkan. Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 694–699, 2002.

Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jindong Wang, Manabu Okumura, and Takahiro Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information Processing Systems*, 34:18408–18419, 2021.

Xiaojin Zhu. Semi-supervised learning literature survey. In *University of Wisconsin-Madison, Department of Computer Sciences*, 2005.

## Supplementary Material

We discuss related works in Appendix A formal algorithm in Appendix B. Additional experimental results and details are in Appendix C.

### Appendix A. Related Work

**Semi-supervised learning (SSL).** There is a rich literature on SSL spanning multiple decades (Zhu, 2005; Chapelle et al., 2006; Singh et al., 2008; Oliver et al., 2018). This literature comprises of a wide variety of approaches. Among these significant focus has been placed on self-training (also called pseudolabeling) (Scudder, 1965; Blum and Mitchell, 1998; Rosenberg et al., 2005; Lee, 2013; Oymak and Gulcu, 2020; Amini et al., 2023), generative models Nigam et al. (2000); Adams and Ghahramani (2009); Kingma et al. (2014), graph-based strategies (Blum and Chawla, 2001; Niyogi, 2013; Subramanya and Talukdar, 2022), and transductive approaches (Vapnik et al., 1998; Joachims, 1999). Due to their simplicity, pseudolabeling-based approaches have gained prominence and are widely used in application areas such as NLP (Karamanolakis et al., 2021), speech recognition (Kahn et al., 2020), and protein prediction (El-Manzalawy et al., 2016). Our paper focuses on recent variants of this, discussed next.

**Pseudolabeling based SSL.** These methods generate artificial labels for unlabeled data and use them for training the model. A crucial challenge here is the issue of confirmation bias (Arazo et al., 2020) i.e., when a model starts to reinforce its own mistakes. To overcome this and to maintain high quality of pseudolabels, confidence-based thresholding is applied. Here only the unlabeled data where confidence is higher than a particular threshold is used (Sohn et al., 2020). Due to the limitations of fixed thresholds, adaptive thresholds based on the classifier’s learning status have been introduced to improve performance (Xu et al., 2021; Zhang et al., 2021; Wang et al., 2023). Nearly all of these methods also use some form of consistency regularization (Laine and Aila, 2017; Bachman et al., 2014; Sajjadi et al., 2016; Fan et al., 2021; Kukačka et al., 2017) where the core idea is that the model should produce similar prediction when presented with different versions (perturbations) of inputs and all the present SSL methods (Xie et al., 2020; Wang et al., 2023; Sohn et al., 2020; Zhang et al., 2021; Chen et al., 2023; Xu et al., 2021).

**Confidence functions and calibration.** Miscalibration (overconfidence) in neural networks plagues various applications (Nguyen et al., 2015; Hendrycks and Gimpel, 2017; Guo et al., 2017), including SSL. To mitigate this in general, a range of solutions have been proposed, including training-time methods (Moon et al., 2020; Kumar et al., 2018; Hui et al., 2023; Corbière et al., 2019; Foret et al., 2021) and post-hoc methods (Guo et al., 2017; Kumar et al., 2019; Gupta and Ramdas, 2022; Kull et al., 2019; Zadrozny and Elkan, 2002). In pseudolabeling based SSL, recent works (Rizve et al., 2021; Loh et al., 2023; Mishra et al., 2024) noted the issue of miscalibration. To promote calibration, Loh et al. (2023) use Bayesian neural nets by replacing the model’s final layer with a Bayesian layer. Rizve et al. (2021) improve pseudolabeling with negative labels and an uncertainty-aware pseudolabel selection technique. Mishra et al. (2024) incorporate a regularizer in pseudolabeling to encourage calibration.

While calibration is generally desirable, it may not be enough to solve the overconfidence issue in SSL and other applications. Pseudolabeling requires scores that effectively distinguish correct from incorrect predictions, aligning with the ordinal ranking criterion (Hendrycks and Gimpel, 2017; Moon et al., 2020; Foret et al., 2021; Corbière et al., 2019). Instead of trial-and-error with various options, we propose a flexible framework that learns confidence functions directly optimized for pseudolabeling objectives. This builds upon principles used in threshold-based auto-labeling (TBAL) (Vishwakarma et al., 2024), a technique for creating labeled datasets.

## Appendix B. Appendix to the Method Section

The full algorithm we use is:

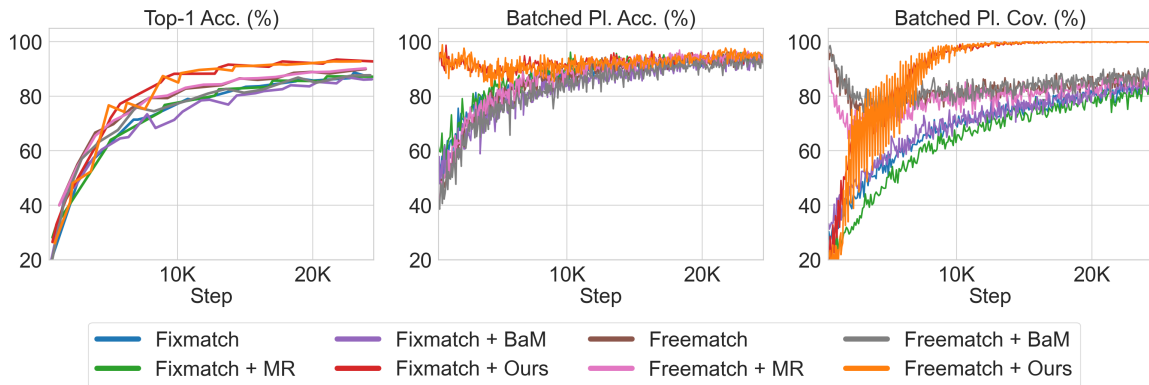


Figure 1: Left to Right: Top-1 accuracy, Batched pseudolabeling accuracy and Batched pseudolabeling coverage of our method and baselines on CIFAR-10. We plot the values for every 200 steps.

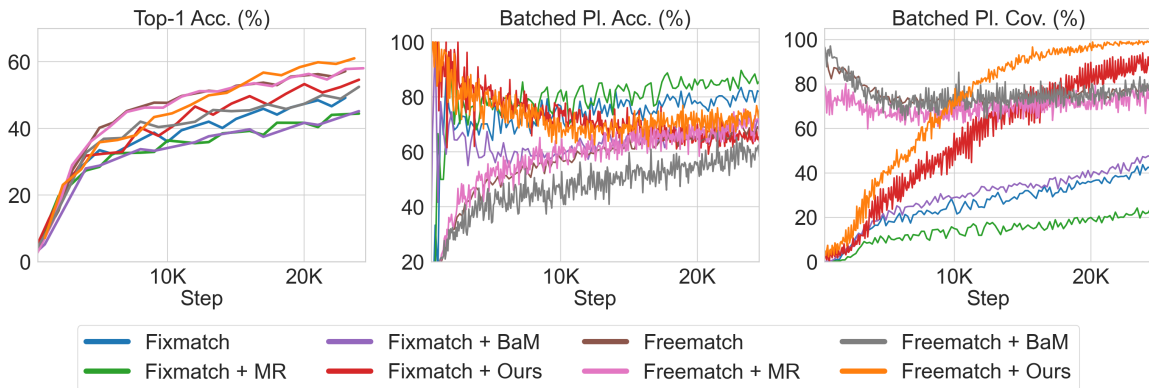


Figure 2: Left to Right: Top-1 accuracy, Batched pseudolabeling accuracy and Batched pseudolabeling coverage of our method and baselines on CIFAR-100. We plot the values for every 200 steps.

### Appendix C. Additional Experiments and Details

**Compute.** For all our experiments, we used an NVIDIA RTX A6000 which has 48GB of VRAM and an NVIDIA RTX 4090 with 24GB of VRAM. The runtime depends on several factors including CPU I/O and GPU load, but on average, the baselines took around 8 hours, while our method took around 15 hours for 25K iterations.

**Hyperparameters.** For the baselines, we have used their default settings. To maintain consistency and experiment the efficiency of method, we used WRN-28-2 which is 1.4M parameter model for all the datasets. We summarize the main hyperparameters we have used in our method in Table 3.



Table 3: Hyperparameters used for our method.

Method	Hyperparameter	Values
Learning $g$ function	optimizer	SGD
	learning rate	0.01
	batch size	64
	max epoch	500
	weight decay	0.01
	momentum	0.9
Estimating $t$	optimizer	SGD
	learning rate	0.01
	batch size	64
	max epoch	500
	weight decay	0.01
	momentum	0.9

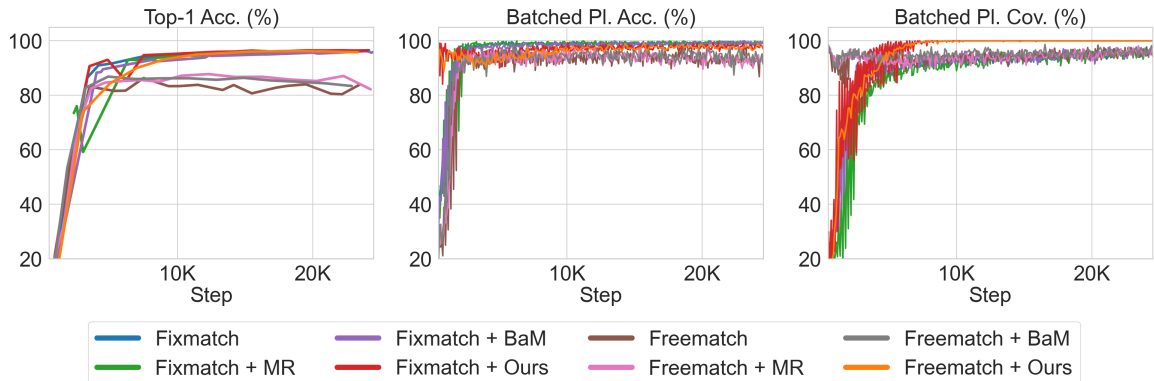


Figure 3: Left to Right: Top-1 Accuracy, Batched pseudolabeling Accuracy and batched pseudolabeling coverage of our method and various baselines on SVHN. We plots the values for every 200 steps.

### C.1 Ablation Studies

We perform ablations that give insights into the role of various parts of it. We run all the ablation experiments on the CIFAR-10 data setting.

**A1. Is pseudolabel accumulation helpful?** Accumulation allows the methods to use old pseudolabel for points that couldn't get pseudolabeled in the current iteration. Thus we expect accumulation could help in improving the utilization of unlabeled data and could lead to better test accuracy in cases where the pseudolabel quality is assured to be high in all iterations. We run two variations of our method and baselines — with accumulation and without it and report the results in Table 4. We observe that our method has similar test accuracy irrespective of accumulation. However, with accumulation it achieves better coverage

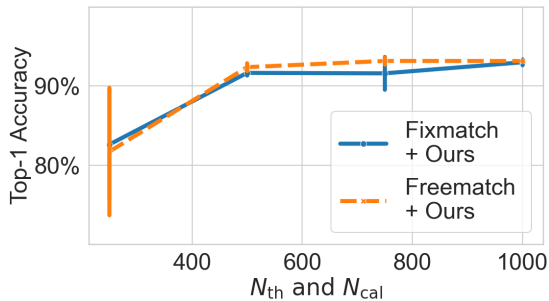


Figure 4: Top-1 accuracy of our method with different  $N_{th}$  and  $N_{cal}$ .

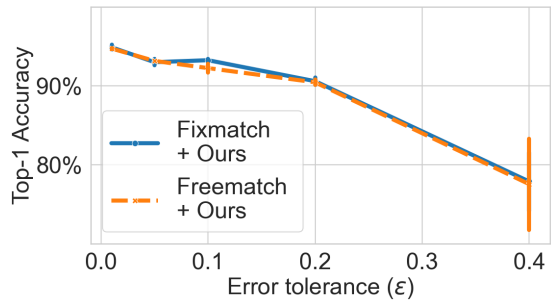


Figure 5: Top-1 accuracy of our method with different error tolerance  $\epsilon$ .

in early iterations as observed in Figure 6. These results are not surprising, since our method ensures high quality of pseudolabels while maximizing coverage, it is able to eventually catch up with the version using accumulation, leading to similar final test accuracies. On the other hand, having accumulation hurts the performance of baseline models. This might be because the pseudo labels generated by the baseline models are not accurate especially in the earlier iterations, thus degrading the overall performance. Overall, we believe accumulation is going to be helpful when we have pseudolabels with high accuracy. The plots for coverage and accuracy over the entire run are in Figures 7, 8 in the Appendix C.

## A2. Does error tolerance affect performance?

In our method, the error tolerance parameter  $\epsilon$  is a knob to control the amount of noise in pseudolabels. A common wisdom in pseudolabeling is higher noise will lead to worse performance, which is our expectation too. To see this, we run our method with  $\epsilon \in \{0.01, 0.05, 0.1, 0.2, 0.4\}$  in the CIFAR-10 setting. We run each setting with 3 random seeds and report the results in Figure 5. The results are as expected — higher values of  $\epsilon$  lead to degraded test accuracy due to high noise in the pseudolabels and with decreasing  $\epsilon$  leads to improved accuracy. These results also suggest that prioritizing the quality (accuracy) of pseudolabels over quantity is a better choice in pseudolabeling. The results are also summarized in Table 6 and Figure 10.

Table 4: Results on CIFAR-10 with and without pseudolabel accumulation (Acc) for all the methods.

Method	Acc—True	Acc—False
Fixmatch	66.30 $\pm$ 1.68	88.15 $\pm$ 1.27
Fixmatch + MR	64.24 $\pm$ 1.93	87.85 $\pm$ 1.10
Fixmatch + BaM	84.50 $\pm$ 2.60	86.44 $\pm$ 1.47
Freematch	85.17 $\pm$ 4.74	90.17 $\pm$ 0.13
Freematch + MR	80.67 $\pm$ 2.39	90.17 $\pm$ 0.45
Freematch + BaM	88.92 $\pm$ 0.49	88.34 $\pm$ 0.99
<b>Fixmatch + Ours</b>	<b>93.03 <math>\pm</math> 0.44</b>	<b>93.34 <math>\pm</math> 0.50</b>
<b>Freematch + Ours</b>	<b>93.08 <math>\pm</math> 0.05</b>	<b>93.01 <math>\pm</math> 0.24</b>

## A3. How much data is needed to learn the $g$ and $\mathbf{t}$ ?

We take  $N_{cal}$  and  $N_{th}$  from the validation data to learn the confidence function  $g$  and estimate the thresholds  $\mathbf{t}$  respectively. Intuitively larger values of these should lead to good  $g$  and  $\mathbf{t}$  that can extract the expected level of pseudolabeling coverage and accuracy from the classifier at hand. However, the task of learning good  $g$  and estimating thresholds is not super hard and we expect it will take fewer samples to be successful. To understand this better we run our method with  $N_{cal}$  and

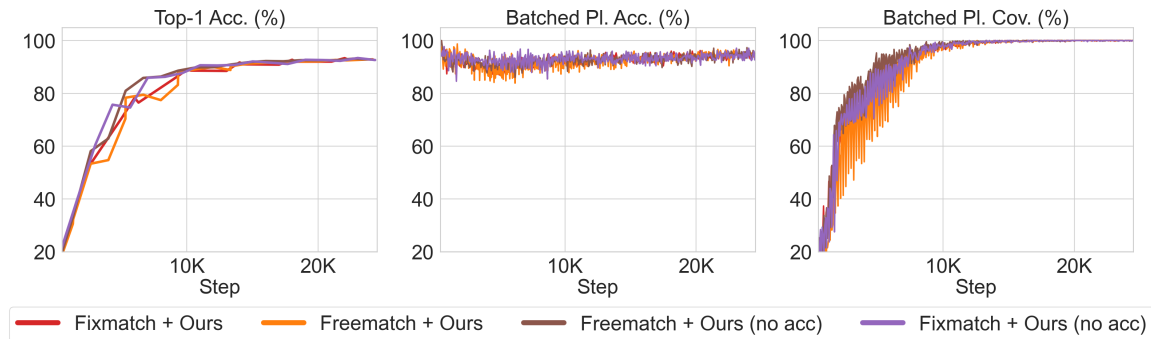


Figure 6: Left to Right: Top-1 Accuracy, Batched pseudolabeling Accuracy and Batched pseudolabeling coverage of our method with and without pseudolabeling accumulation enabled.

Table 5: Results on CIFAR-10 with varying  $N_{\text{cal}}$  and  $N_{\text{th}}$ .

Method	$N_{\text{cal}} = N_{\text{th}} = 250$	$N_{\text{cal}} = N_{\text{th}} = 500$	$N_{\text{cal}} = N_{\text{th}} = 750$
Fixmatch + Ours	$82.67 \pm 7.08$	$91.74 \pm 0.41$	$91.66 \pm 2.11$
Freematch + Ours	$82.13 \pm 7.93$	$92.33 \pm 0.49$	$93.20 \pm 0.53$

$N_{\text{th}}$  in  $\{250, 500, 750, 1000\}$  on CIFAR-10 setting for 3 random seeds and report the result in Fig 4. We observe that our method can achieve desired performance with just 500 labeled points (i.e 50 labels per class). This is interesting because we can achieve 90% accuracy by just using 250 points ( $N_l$ ) for training  $h$  and a total of 1K for learning  $g$ . Refer Table 5 and Figure 9 for more details.

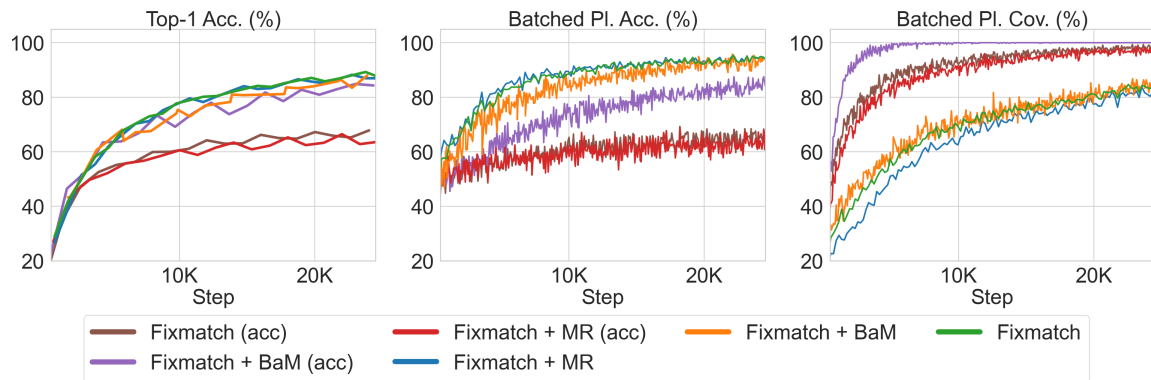


Figure 7: **(A1.)** Left to Right: Top-1 Accuracy, Batched pseudolabeling Accuracy and batched pseudolabeling coverage of Fixmatch with and without pseudolabeling accumulation enabled on CIFAR-10. It can be seen that enabling pseudolabeling accumulation worsen the performance of baseline methods in terms of accuracy and coverage.

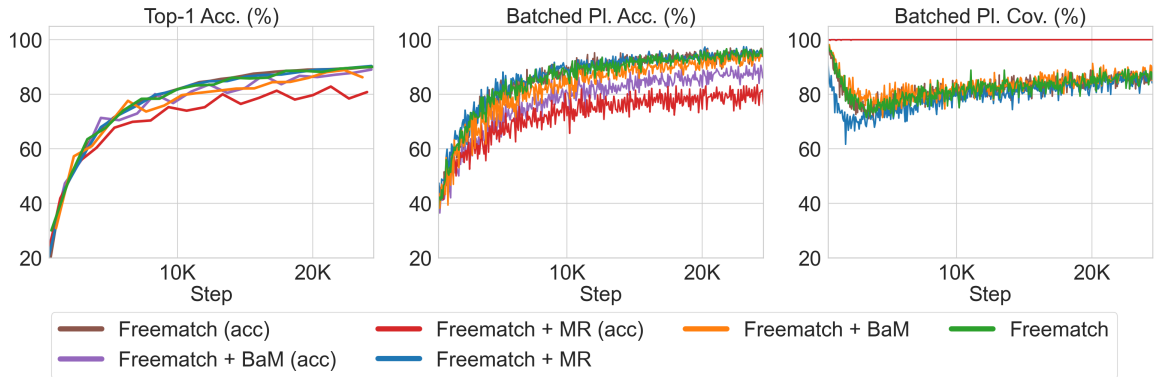


Figure 8: **(A1.)** Left to Right: Top-1 Accuracy, Batched pseudolabeling Accuracy and batched pseudolabeling coverage of Freematch with and without pseudolabeling accumulation enabled on CIFAR-10. It can be seen that enabling pseudolabeling accumulation worsen the performance of baseline methods in terms of accuracy and coverage.

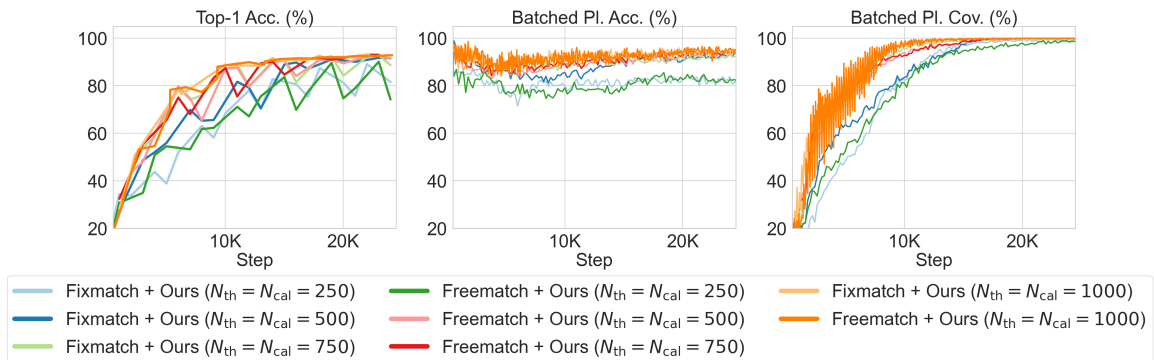


Figure 9: **(A3.)** Left to Right: Top-1 Accuracy, Batched pseudolabeling Accuracy and batched pseudolabeling coverage of our method with  $N_{th} = N_{cal} \in \{250, 500, 750, 1000\}$  on CIFAR-10. We observe that having more calibration and threshold estimation points benefits the performance of our method.

Table 6: Results on CIFAR-10 with varying  $\epsilon$ .

Method	$\epsilon = 0.01$	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.4$
Fixmatch + Ours	<b>94.85 ± 0.28</b>	93.24 ± 0.18	90.52 ± 0.43	80.62 ± 1.22
Freematch + Ours	<b>94.67 ± 0.09</b>	92.11 ± 0.84	90.20 ± 0.65	82.23 ± 1.31

---

**Algorithm 1** Pseudolabeling Based SSL with PabL0
 

---

**Input:** Labeled data for training  $D_l$ , Validation data  $D_{\text{val}}$ , unlabeled pool  $X_u$ , error tolerance  $\epsilon$ , use-accumulation flag, num\_iters, batch size  $B$ , replication factor  $\mu$ , weak  $\omega$  and strong  $\Omega$  augmentations.

**Output:**  $\hat{h}_{\text{ssl}}$ , model with the best validation accuracy.

```

1:  $\tilde{Y} \leftarrow [0] \times n_u, S \leftarrow [0] \times n_u, i \leftarrow 1.$ 
2:  $D_{\text{cal}}, D_{\text{th}} \leftarrow \text{draw\_randomly}(D_{\text{val}}, N_{\text{cal}}, N_{\text{th}})$ 
3: while  $i \leq \text{num\_iters}$  do
4:    $D_l^b, X_u^b, I_u^b \leftarrow \text{draw\_random\_batch}(\mu D_l, \mu X_u, B)$ 
5:    $X_{u,w}^b, X_{u,s}^b \leftarrow \omega(X_u^b), \Omega(X_u^b)$ 
6:   if use-PabL0 then
7:     if  $i \% F = 0$  then
8:        $\hat{g}_i \leftarrow \text{solve\_opt\_problem\_P1}(\hat{h}_i, D_{\text{cal}})$ 
9:        $\hat{t}_i \leftarrow \text{estimate\_thresholds}(\hat{h}_i, \hat{g}_i, D_{\text{th}})$ 
10:       $\tilde{Y}^f \leftarrow \hat{h}_i(\omega(X_u)), S^f \leftarrow \mathbf{1}(\hat{g}_i(\omega(X_u)) \geq \hat{t})$ 
11:      if use-accumulation then
12:         $\tilde{Y}, S \leftarrow S^f \tilde{Y}^f + (1 - S^f) \tilde{Y}; S \leftarrow S \vee S^f$ 
13:      else
14:         $\tilde{Y}, S \leftarrow \tilde{Y}^f, S^f$ 
15:      end if
16:    end if
17:     $\tilde{Y}^b, S^b \leftarrow \tilde{Y}[I_u^b], S[I_u^b]$ 
18:  else
19:     $\tilde{Y}^b, S^b \leftarrow \text{baseline\_pseudo\_labeling}(\hat{h}_i, X_{u,w}^b)$ 
20:    if use-accumulation then
21:      for  $j \in I_u^b$  do
22:         $\tilde{Y}[j] \leftarrow S^b[j] \tilde{Y}^b[j] + (1 - S^b[j]) \tilde{Y}[j]$ 
23:         $S[j] \leftarrow S[j] \vee S^b[j]$ 
24:      end for
25:    end if
26:  end if
27:   $\hat{\mathcal{L}}_s(\hat{h}_i) \leftarrow \text{supervised\_loss}(h, D_l^b)$ 
28:   $\hat{\mathcal{L}}_u(\hat{h}_i) \leftarrow \text{unsupervised\_loss}(h, X_{u,w}^b, X_{u,s}^b, \tilde{Y}^b, S^b)$ 
29:   $\hat{\mathcal{L}}_r(\hat{h}_i) \leftarrow \text{baseline\_regularizers}()$ 
30:   $\hat{\mathcal{L}}(\hat{h}_i) \leftarrow \hat{\mathcal{L}}_s(\hat{h}_i) + \lambda_u \hat{\mathcal{L}}_u(\hat{h}_i) + \lambda_r \hat{\mathcal{L}}_r(\hat{h}_i)$ 
31:   $\hat{h}_{i+1} \leftarrow \text{SGD\_update}(\hat{\mathcal{L}}(\hat{h}_i)); i \leftarrow i + 1$ 
32:  if  $i \% \text{eval\_freq} = 0$  then
33:    eval_acc  $\leftarrow \text{evaluate\_model}(\hat{h}_i, D_{\text{val}})$ 
34:    If eval_acc is best so far then  $\hat{h}_{\text{ssl}} = \hat{h}_i.$ 
35:  end if
36: end while

```

---

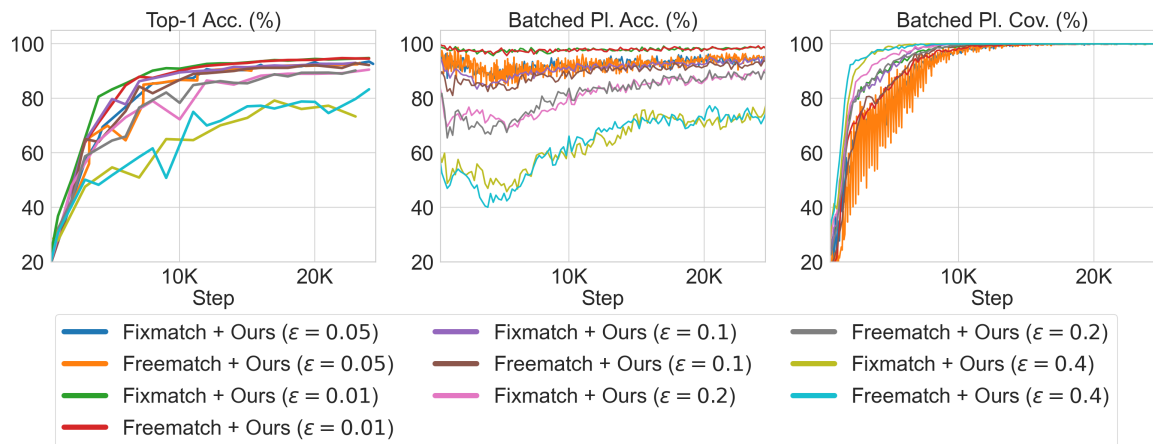


Figure 10: **(A2.)** Left to Right: Top-1 Accuracy, Batched pseudolabeling Accuracy and batched pseudolabeling coverage of our method with  $\epsilon \in \{0.01, 0.05, 0.1, 0.2, 0.4\}$  on CIFAR-10. Although having a looser constraint on the error encourages more coverage, the pseudolabeling drops as a trade-off.