
Rethinking Confidence Scores and Thresholds in Pseudolabeling-based SSL

Harit Vishwakarma^{*1} Yi Chen^{*2} Satya Sai Srinath Namburi GNVV³ Sui Jiet Tay⁴
Ramya Korlakai Vinayak² Frederic Sala¹

Abstract

Modern semi-supervised learning (SSL) methods rely on pseudolabeling and consistency regularization. Pseudolabeling is typically performed by comparing the model’s confidence scores and a predefined threshold. While several heuristics have been proposed to improve threshold selection, the underlying issues of overconfidence and miscalibration in confidence scores remain largely unaddressed, leading to inaccurate pseudolabels, degraded test accuracy, and prolonged training. We take a first-principles approach to learn confidence scores and thresholds with an explicit knob for error. This flexible framework addresses the fundamental question of optimal scores and threshold selection in pseudolabeling. Moreover, it gives practitioners a principled way to control the quality and quantity of pseudolabels. Such control is vital in SSL, where balancing pseudolabel quality and quantity directly affects model performance and training efficiency. Our experiments show that, by integrating this framework with modern SSL methods, we achieve significant improvements in accuracy and training efficiency. In addition, we provide novel insights on the trade-offs between the choices of the error parameter and the end model’s performance.

1. Introduction

The lack of high-quality labeled data is a major bottleneck in training high-accuracy models. The semi-supervised learning (SSL) paradigm tackles this problem by leveraging abundant unlabeled data alongside a limited set of labeled examples (Chapelle et al., 2006; Zhu, 2005; van Engelen

& Hoos, 2019). While SSL dates back decades and includes a wide variety of approaches, modern SSL methods frequently rely on a pair of ideas: *self-training or pseudolabeling* – where model generated labels are assigned to unlabeled data for further training (McLachlan, 1975; Amini et al., 2023; Rosenberg et al., 2005; Lee, 2013; Rizve et al., 2021) – and *consistency regularization* to enforce stability in predictions across perturbed inputs (Laine & Aila, 2017; Bachman et al., 2014; Sajjadi et al., 2016; Fan et al., 2021; Kukačka et al., 2017). SSL techniques with these ideas have strongly performed on several benchmark datasets.

While pseudolabeling is a powerful technique, its effectiveness hinges on a fundamental question: *which points should be labeled using the model’s predictions?* Since pseudolabels are derived from a model being trained, they can be highly unreliable. A naive approach that assigns pseudolabels too liberally risks injecting noisy labels into training, amplifying the model’s existing errors – a phenomenon known as confirmation bias. Conversely, an overly conservative approach that selects only the correct predictions severely limits the amount of useful training data. Both extremes can degrade SSL performance, leading to either poor generalization or slow convergence. To fully harness the potential of pseudolabeling, we need a principled approach for pseudolabeling with explicit control on these trade-offs.

A widely used strategy pseudolabels points on which the model’s confidence score exceeds a threshold. This approach provides a simple mechanism for selecting points while controlling the quality and quantity of pseudolabels via thresholds. However, the prior works based on this approach suffer from two key limitations that restrict their effectiveness. First, thresholding techniques are often heuristic-driven, lacking precise control for a target error level (Sohn et al., 2020; Zhang et al., 2021; Wang et al., 2023; Xu et al., 2021; Chen et al., 2023; Li et al., 2023). Second, commonly used scores, such as the model’s softmax outputs, tend to be unreliable. Recent studies in this vein (Loh et al., 2022; Mishra et al., 2024; Rizve et al., 2021) have highlighted issues of overconfidence and miscalibration in these scores, leading to inefficiencies in pseudolabeling.

In addition to these problems with the choices of scores and thresholds, an equally important question remains: how

^{*}Equal contribution ¹Department of Computer Sciences, University of Wisconsin-Madison, WI, USA ²Department of Electrical and Computer Engineering, University of Wisconsin-Madison, WI, USA ³GE HealthCare ⁴NYU Courant Institute. Correspondence to: Harit Vishwakarma <hvishwakarma@cs.wisc.edu>, Yi Chen <yi.chen@wisc.edu>.

much error should be allowed in pseudolabeling? As discussed earlier, a very low tolerance may pseudolabel too few points, and conversely, a high tolerance may allow for large errors, hurting the efficiency of pseudolabeling in both cases. To solve these issues, we seek *a principled solution to get confidence scores and thresholds for pseudolabeling at a given target level of error tolerance.*

We take a first principles approach to formalize the pseudolabeling objective: maximize the number of pseudolabeled points while adhering to the given error tolerance. We formulate this objective as an optimization problem over a flexible space of confidence functions and thresholds. By solving this optimization problem, we obtain confidence scores tailored to our objective of pseudolabeling. To ensure the pseudolabeling error constraint is strictly followed, we use a separate procedure to estimate thresholds on these scores using part of the validation data. Pseudolabeling with these scores and thresholds ensures we are pseudolabeling a maximal set of points that can be pseudolabeled at the given error tolerance level.

We integrate this approach into popular pseudolabeling-based SSL methods, providing two benefits. First, it provides a principled way to derive confidence scores and thresholds for any given error tolerance. Second, by enabling more precise pseudolabeling, it improves the utilization of the unlabeled data and is expected to yield an end model with higher test accuracy compared to the ad-hoc choices of scores and thresholds.

Our main contributions are summarized as follows,

1. Our work settles the question of the right choices of confidence scores and thresholds in pseudolabeling-based SSL methods by introducing a framework for learning confidence scores and thresholds. Departing from heuristic-driven or ad-hoc and unreliable choices of scores and thresholds, this framework provides principled choices of scores and thresholds for pseudolabeling with any target error tolerance.
2. We show how this framework for learning confidence scores and thresholds can work in concert with popular SSL methods such as Fixmatch (Sohn et al., 2020), FreeMatch (Wang et al., 2023), etc., and conduct an extensive empirical evaluation demonstrating that by pseudolabeling using confidence scores and thresholds learned from our method can yield significant improvements in the test accuracy.
3. Leveraging our framework’s ability to pseudolabel at any target error level, we study the impact of varying pseudolabeling error levels—from fixed to dynamic tolerance throughout training. Our results confirm the intuition that lower pseudolabeling errors lead to better end models compared to higher errors. Moreover,

among dynamic schedules, it is better to use a decreasing schedule of error tolerances.

2. Related Work

Semi-supervised learning (SSL). There is a rich literature on SSL (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 & Mitchell, 1998; Rosenberg et al., 2005; Lee, 2013; Oymak & Gulcu, 2020; Amini et al., 2023), generative models (Nigam et al., 2000; Adams & Ghahramani, 2009; Kingma et al., 2014), graph-based strategies (Blum & Chawla, 2001; Niyogi, 2013; Subramanya & 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 a high quality of pseudolabels, confidence-based thresholding is applied. Here, only the unlabeled data with confidence 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 & 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; Zheng et al., 2022; Xu et al., 2021; Roelofs et al., 2022).

Confidence functions and calibration. Miscalibration (overconfidence) in neural networks plagues various applications (Nguyen et al., 2015; Hendrycks & 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 & Ramdas, 2022; Kull et al., 2019; Zadrozny & 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) utilize negative labels and an uncertainty-aware pseudolabel selection technique. Mishra et al. (2024) incorporate a regularizer to encourage calibration.

While calibration is a reasonable goal in general, it may not be sufficient to address the overconfidence problem in SSL and other applications. In pseudolabeling, we seek the use of scores that can easily segregate the model’s correct and incorrect predictions, which is closely related to the ordinal ranking criterion (Hendrycks & Gimpel, 2017; Moon et al., 2020; Foret et al., 2021; Corbière et al., 2019). Rather than experimenting with several such choices, ideally, we would have a flexible framework that can learn confidence functions explicitly optimizing pseudolabeling objectives.

3. 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 threshold for class y .

3.1. 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}_{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 with respect 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 positive constants λ_u, λ_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, $\hat{\mathcal{L}}_s(h | D_l^b) = \frac{1}{|D_l^b|} \sum_{(\mathbf{x}, y) \in D_l^b} H(y, h, \mathbf{x})$. 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 computed. It is 1 for points having confidence score $g(\mathbf{x})[\tilde{y}]$, bigger than pre-determined threshold $\mathbf{t}[\tilde{y}]$, corresponding to the predicted class \tilde{y} . Recent methods, couple this loss and consistency regularization together by doing pseudolabeling on weakly augmented data using weak transform $\omega : \mathcal{X} \mapsto \mathcal{X}$ and then defining the cross-entropy loss on the strongly augmented data using strong transformation $\Omega : \mathcal{X} \mapsto \mathcal{X}$. The loss is,

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

Regularization. A regularization term (or a sum of multiple regularizers) is often included along with the above two losses to encourage desired behavior(s). 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 the model’s confidence scores. Including such regularizers has been fruitful in pseudolabeling-based SSL.

3.2. Quality and Quantity of Pseudolabels

Given a classifier h , the quality and quantity of the pseudolabels w.r.t. to score function g and thresholds \mathbf{t} , are:

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

$$\hat{\mathcal{P}}(g, \mathbf{t} | h, X) := \frac{1}{|X|} \sum_{(\mathbf{x}) \in X} S(\mathbf{x}, g, \mathbf{t} | h) \quad (1)$$

$$\mathcal{P}(g, \mathbf{t} | h) := \mathbb{E}_{\mathbf{x}}[S(\mathbf{x}, g, \mathbf{t} | h)]. \quad (2)$$

Pseudolabeling error (quality). This is the fraction of pseudolabeled points that got incorrect labels. This metric captures the quality of pseudolabels:

$$\hat{\mathcal{E}}(g, \mathbf{t} | h, D) := \frac{\sum_{(\mathbf{x}, y, \tilde{y}) \in D} S(\mathbf{x}, g, \mathbf{t} | h) \cdot \mathbb{1}(h(\mathbf{x}) \neq y)}{\sum_{(\mathbf{x}, y, \tilde{y}) \in D} S(\mathbf{x}, g, \mathbf{t} | h)}, \quad (3)$$

$$\mathcal{E}(g, \mathbf{t} | h) = \frac{\mathbb{E}_{\mathbf{x}}[S(\mathbf{x}, g, \mathbf{t} | h) \cdot \mathbb{1}(h(\mathbf{x}) \neq y)]}{\mathcal{P}(g, \mathbf{t} | h)}. \quad (4)$$

3.3. Our Goals

Pseudolabeling-based SSL aims to learn a classifier \hat{h}_{ssl} that generalizes well on the unseen data, i.e., has high test accuracy. This is typically achieved by pseudolabeling points using confidence scores and thresholds and incorporating them into training. However, existing choices of scores and thresholding strategies are often ad hoc and unreliable, limiting their effectiveness. Departing from these unreliable approaches, our goal is to:

- (i) Design principled solutions for confidence scores and thresholding to maximize the number of pseudolabeled points while ensuring the pseudolabeling error is at most ϵ .
- (ii) Incorporate these in the existing pseudolabeling-based SSL methods and assess whether this gives a better end model \hat{h}_{ssl} .
- (iii) Study the sensitivity of the SSL pipeline to pseudolabeling errors by leveraging the ability of our approach to explicitly ensure the pseudolabeling error remains below ϵ .

4. Methodology

In this section, we discuss our principled framework to learn scores and thresholds with explicit control of the pseudolabeling errors and use them in pseudolabeling-based SSL. This framework is an adaptation of a recent solution designed to improve confidence functions for auto-labeling (Vishwakarma et al., 2024). The detailed steps are outlined in Algorithm 3 in the Appendix.

4.1. Pseudolabeling Optimization Framework

Given the current model \hat{h}_i in the i th iteration, can we obtain confidence scores and thresholds using which we can identify a maximal set of points that can be pseudolabeled using \hat{h}_i with at most ϵ error? We begin with a theoretical formulation to learn such scores and thresholds, and then introduce its practical version.

Theoretical framework. Instead of improving 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. (2)) while keeping the pseudolabeling error (eq. (4)) below a tolerance level $\epsilon \in (0, 1)$. More specifically, given the classifier \hat{h}_i in any iteration i of SSL,

$$g_i^*, \mathbf{t}_i^* \in \arg \max_{g \in \mathcal{G}, \mathbf{t} \in T^k} \mathcal{P}(g, \mathbf{t} \mid \hat{h}_i) \text{ s.t. } \mathcal{E}(g, \mathbf{t} \mid \hat{h}_i) \leq \epsilon, \quad (5)$$

are the optimal confidence functions and thresholds for pseudolabeling using \hat{h}_i 's predictions such that the pseudolabeling error is bounded by ϵ . This frees us from arbitrary choices of confidence scores, calibration techniques, and

thresholding heuristics. Instead, solving the optimization problem over a flexible enough space will subsume specific strategies. Next, we discuss how to make this framework tractable to obtain scores and thresholds in practice.

Practical version. The optimization problem discussed earlier involves population-level quantities which are usually not accessible in practice. Thus we have to use their finite sample estimates and smooth variations to make the optimization problem tractable. Specifically, the coverage and error are estimated using a small amount of held-out labeled data (called calibration data D_{cal}) curated from the validation data. Then differentiable surrogates for the 0-1 variables are introduced. 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,

$$\tilde{\mathcal{P}}(g, \mathbf{t} \mid h, D_{\text{cal}}) := \frac{1}{|D_{\text{cal}}|} \sum_{(\mathbf{x}, y, \tilde{y}) \in D_{\text{cal}}} \sigma(\alpha, g(\mathbf{x})[\tilde{y}] - \mathbf{t}[\tilde{y}]), \quad (6)$$

$$\tilde{\mathcal{E}}(g, \mathbf{t} \mid h, D_{\text{cal}}) := \frac{\sum_{(\mathbf{x}, y, \tilde{y}) \in D_{\text{cal}}} \mathbb{1}(y \neq \tilde{y}) \sigma(\alpha, g(\mathbf{x})[\tilde{y}] - \mathbf{t}[\tilde{y}])}{\sum_{(\mathbf{x}, y, \tilde{y}) \in D_{\text{cal}}} \sigma(\alpha, g(\mathbf{x})[\tilde{y}] - \mathbf{t}[\tilde{y}])} \quad (7)$$

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 error and coverage.

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

We use 2-layer neural nets as a choice of \mathcal{G} and $T^k = [0, 1]^k$. The optimization problem (P1) is non-convex 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.

4.2. Threshold Estimation

While we can obtain both the confidence scores and thresholds by solving (P1), we propose to estimate thresholds separately on a held-out part of the validation data to avoid potential generalization issues due to learning them simultaneously from the same data D_{cal} and ensure that the pseudolabeling error constraint is strictly adhered to.

When dealing with datasets containing many classes there may not be enough samples per class to estimate reliable thresholds. Thus, to accommodate these possibilities we consider two variations of the threshold estimation procedure, (i) estimate a common (joint) threshold for all classes and (ii) estimate separate (classwise) thresholds for each class. The procedures are outlined in Algorithm 2 and Algorithm 1 in the Appendix. We discuss them briefly here.

The procedure takes in a confidence function \hat{g}_i and part of the held-out validation data referred to as D_{th} . The idea is

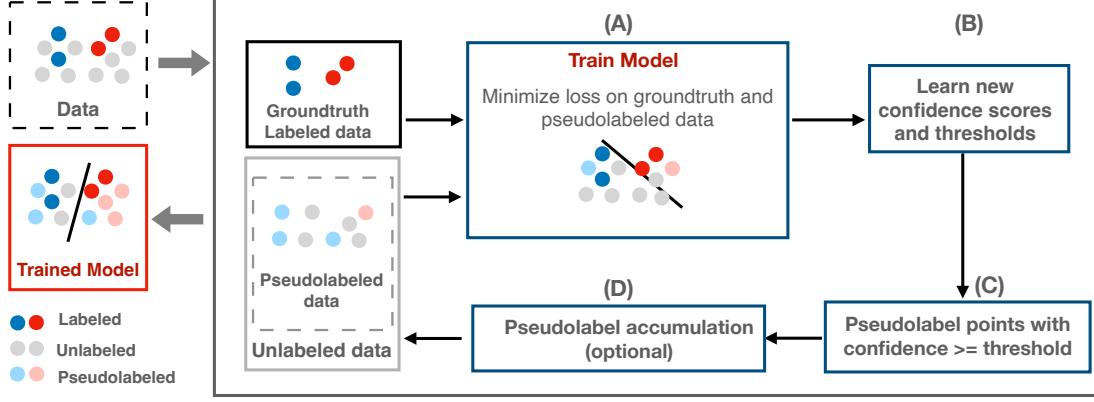


Figure 1. Workflow of pseudolabeling-based SSL with PabLO (A): Train model with standard supervised loss, consistency regularization, and other regularizers (B): Learn new confidence scores and thresholds (C): Pseudolabel points with scores greater than the estimated thresholds (D): An optional pseudolabeling accumulation to use previous pseudolabels for points that are not pseudolabeled in current round. Note, that the pseudolabels can be noisy (incorrect). The training and pseudolabeling loop continues until a pre-specified number of iterations. In the end, it outputs a model \hat{h}_{ssl} that is expected to have higher test accuracy than the model trained only on the given groundtruth training data. Note that in the end, there might be points left unlabeled, and the pseudolabels might be noisy.

to estimate errors at several thresholds on this data and then pick the smallest threshold. This can be done separately for each class to obtain classwise thresholds or a common threshold for all classes. We discuss classwise thresholding here. First the data D_{th} is partitioned into k subsets $D_{th}^{(y)}$ corresponding to each class $y \in \mathcal{Y}$. Here, we slightly abuse notation: instead of $\mathbf{t} \in T^k$, we use $t \in T$ in the estimate of pseudolabeling error at threshold t for class y . To obtain threshold $\hat{\mathbf{t}}[y]$ for class y , the procedure finds the smallest $t \in T$ such that $\hat{\mathcal{E}}(\hat{g}_i, t \mid h, D_{th}^{(y)}) + C_1 \hat{\sigma}(\hat{\mathcal{E}}) \leq \epsilon$. Here C_1 is a constant (we use $C_1 = 0.25$) and $\hat{\sigma}(z) = \sqrt{z \cdot (1 - z)}$ and $\hat{\mathcal{E}}$ is used for brevity in place of $\hat{\mathcal{E}}(\hat{g}_i, t \mid h, D_{th}^{(y)})$. The same process is used for joint threshold estimation, where a single threshold t is estimated using the entire D_{th} and the same t is used for all classes. Using the thresholds found using these procedures ensures pseudolabeling error remains below (or close to) the tolerance level ϵ .

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. Similar procedures have been used in the context of creating reliable datasets and are backed by theoretical guarantees for the quality of pseudolabels produced (Vishwakarma et al., 2023).

4.3. Pseudolabeling and Accumulation

In the usual pseudolabeling-based SSL setups, the pseudolabels inferred by the model for a mini-batch are 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 as “pseudolabel accumulation”.

Mathematically, let $\hat{y}_j^{(i-1)} = \tilde{Y}_u^{(i-1)}[j]$ and $\hat{y}_j^{(i)} = \tilde{Y}_u^{(i)}[j]$ be the previous and current (fresh) pseudolabels for j th unlabeled point. Let the corresponding masks (indicating whether the score is above the threshold) for these pseudolabels be $S_u^{(i-1)}[j]$ and $S_u^{(i)}[j] = \mathbb{1}(\hat{g}_i(\mathbf{x}_j)[\hat{y}_j^{(i)}] \geq \hat{\mathbf{t}}_i[\hat{y}_j^{(i)}])$. Then with accumulation,

$$\begin{aligned} \tilde{Y}_u^{(i)}[j] &\leftarrow S_u^{(i)}[j] \tilde{Y}_u^{(i)}[j] + (1 - S_u^{(i)}[j]) \tilde{Y}_u^{(i-1)}[j], \\ S_u^{(i)}[j] &\leftarrow S_u^{(i)}[j] \vee S_u^{(i-1)}[j]. \end{aligned}$$

Here \vee is the boolean OR operation and the steps are executed in the order. In words, if the point is pseudolabeled in the current iteration (i.e., its current mask is 1), then it will use the current pseudolabel otherwise, if the point was pseudolabeled in earlier iteration(s) it will use the pseudolabel from that iteration and mark the point as pseudolabeled. In case the point is not pseudolabeled in this iteration or any other iteration in the past, it will remain unlabeled. While it is appealing to use this trick, its use is only warranted in settings ensuring high-quality pseudolabels. We try to understand the consequences of the inclusion and exclusion of this trick in pseudolabeling-based SSL via experiments discussed in the next section.

We put together the steps for learning scores, thresholds, and performing (optional) accumulation in a common template

Dataset	Backbone Model h	k	n_u	N_{val}	N_{test}	N_l	N_{cal}	N_{th}	Augmentation
CIFAR-10	WRN-28-2	10	50K	6K	4K	40 or 250	1K	1K	Weak, Strong
CIFAR-100	WRN-28-2	100	50K	6K	4K	400 or 2500	3K	3K	Weak, Strong
SVHN	WRN-28-2	10	604,388	15,620	10,412	40 or 250	3K	3K	Weak, Strong

Table 1. Details of the dataset we use in our experiments. k is the number of classes. N_l is the number of labeled data points used for training the backbone model h . n_u is the number of unlabelled data points used for consistency regularization and pseudolabeling for all the methods. N_{val} is the number of points used for model selection in all methods. N_{test} is the number of test data points. N_{cal} is the number of points used for learning the g function. N_{th} is the number of data points used for threshold estimation.

of pseudolabeling-based SSL algorithms. We refer to this adapted method (Algorithm 3 in Appendix B) as PabLO. The high-level steps are also illustrated in Figure 1. Next, we discuss the empirical evaluation of PabLO and baselines.

5. Experiments

We conduct empirical evaluation over several settings to,

- C1.** Verify that the adaptations of popular pseudolabeling-based SSL methods with PabLO output models with better test accuracy.
- C2.** Study the effects of choice of error tolerance ϵ on test accuracy of the final model.
- C3.** Understand the role of pseudolabel accumulation in our method and baselines.

5.1. Experimental Setup

First, we briefly describe the experimental setup, with details deferred to Appendix C. The code is available on GitHub ¹.

Methods. We use two simple base methods that capture the core ideas of pseudolabeling (PL) and consistency regularization (CR). The first is *Fixmatch* (Sohn et al., 2020), which uses fixed thresholds on (maximum softmax probability) MSP scores for PL and 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)* ² (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). We provide implementations of these in the code submitted along with the paper.

Datasets. We experiment with three datasets: *CIFAR-10* (Krizhevsky et al., 2009) is an image dataset with 10 classes.

CIFAR-100 (Krizhevsky et al., 2009) is an extended version of CIFAR-10 with 100 classes. *SVHN* (Netzer et al., 2011) is a 10-class image dataset of digits from Google Street View. More details are summarized in Table 1. We use a portion of the validation data (D_{val}) for our method, split into D_{cal} , used to learn the function g , and D_{th} , used to estimate the threshold. Unless otherwise mentioned, we use N_l as 250 for CIFAR-10 and SVHN and 2500 for CIFAR-100 in our experiments.

Adjusted iterations for baselines. Empirically, our method requires more time to run compared to base SSL techniques. Therefore, we adopt the following strategy to ensure a fair comparison between the baselines and our method: First, we train our method for 25K iterations and obtain the average per iteration time, denoted as α_o . Then, we train each baseline method b for 5K iterations and obtain the average per iteration time, denoted as α_b . Using these two values, we obtain the adjusted number of iterations, $\frac{\alpha_o}{\alpha_b} \times 25000$, for the baseline method b . Coincidentally, baselines under the same dataset have similar runtimes. We, therefore, set the adjusted number of iterations on a dataset level. For CIFAR-10, the adjusted number of iterations for baselines is 37,000. For CIFAR-100, the adjusted number of iterations for baselines is 70,000. For SVHN, it is 145,000.

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). For confidence functions class \mathcal{G} , we use a class of 2-layer neural nets and provide the last two layers’ representations from h as input. We train it using SGD. The hyperparameters are deferred to Appendix C. Unless otherwise specified, our method uses pseudolabeling error tolerance $\epsilon = 5\%$.

5.2. 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 (in Table 2) the mean and standard deviation of accuracy across three runs.

¹<https://github.com/harit7/PabLO-SSL>

²We assign this name for convenience.

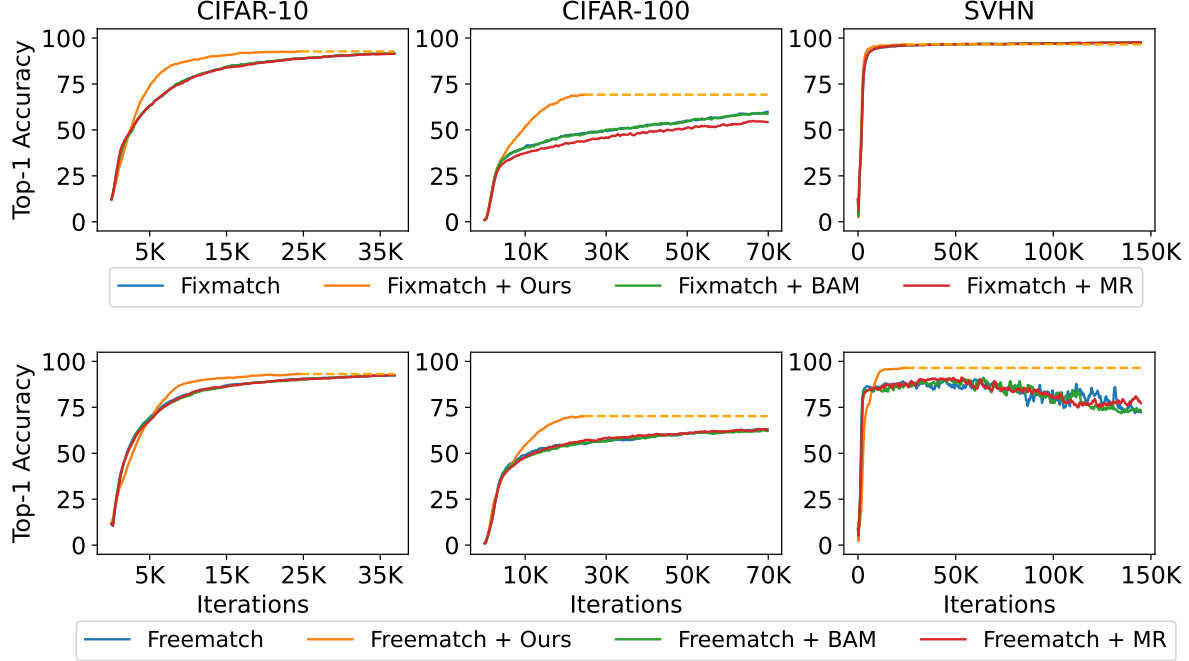


Figure 2. Top-1 accuracy of our method and baselines on CIFAR-10, CIFAR-100, and SVHN. We plot the values for every 200 steps.

C1. Test accuracy improvements. Since our method maximizes the pseudolabeling coverage and accuracy, it provides more accurate pseudolabels for model training. Therefore, we expect it to yield a model with better test accuracy than the baselines. We report the test accuracies at the end of 25K iterations in Table 2 for our methods. For the baselines, we report the test accuracies at the end of the corresponding adjusted number of iterations (well above 25K). Figure 2 illustrates how the top-1 accuracy evolves during the SSL. Similarly, Figure 4 and 5 show how batch pseudolabeling accuracy and batch pseudolabeling coverage change.

First, as expected, integrating our method into the base methods improves test accuracy across all settings. For CIFAR-10, using it with Fixmatch provides almost 2% improvement over Fixmatch alone, and using it with Freematch yields 1% improvement over Freematch. Much more significant improvements are observed in the much harder setting of CIFAR-100: a nearly 10% improvement in top-1 accuracy over Fixmatch and around 5% improvement over Freematch. SVHN is an easier setting; here, the improvements are marginal. With Fixmatch, our performance is similar to that of the baselines. But, using PabLO with Freematch improves the performance by 3%.

C2. Error tolerance affects performance. In our method, the error tolerance parameter ϵ is a knob to control the amount of noise in pseudolabels. A common wisdom in pseudolabeling is that higher noise will lead to worse performance, which is our expectation. To see this, we run our method with $\epsilon \in \{0.01, 0.05, 0.1, 0.2, 0.4\}$ in CIFAR-10

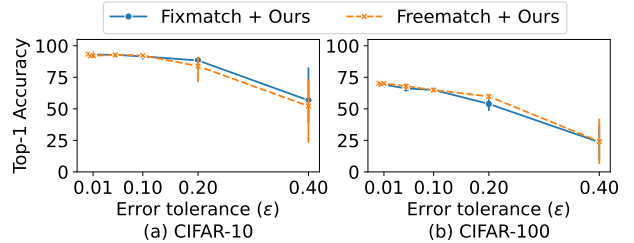


Figure 3. Top-1 accuracy of our method with different error tolerance ϵ on (a) CIFAR-10 and (b) CIFAR-100 dataset.

and CIFAR-100 settings, each with three random seeds, and report the results in Figure 3. The results are as expected — higher values of ϵ lead to degraded test accuracy due to high noise in the pseudolabels and with decreasing ϵ 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 8 and Table 9 in the Appendix.

To investigate the error tolerance further, we designed error tolerance scheduling using different error tolerances during various stages of SSL training. Table 3 summarizes the error tolerance we set at different iterations of SSL training and the corresponding top-1 accuracy for the CIFAR-10 and CIFAR-100 datasets. As we see, starting SSL with a small error tolerance and ending with a large tolerance severely impacts the performance on both CIFAR-10 and CIFAR-100 datasets. While our findings suggest a lower error tolerance is preferable, this may not hold in general. Nevertheless, our

Dataset	CIFAR-10		CIFAR-100		SVHN	
# Labels	40	250	400	2500	40	250
Fixmatch	57.16 \pm 12.12	90.8 \pm 0.78	30.38 \pm 0.68	59.09 \pm 1.10	62.27 \pm 5.57	97.57 \pm 0.08
Fixmatch + MR	54.76 \pm 5.92	90.41 \pm 0.83	16.31 \pm 1.44	54.16 \pm 0.18	74.41 \pm 16.61	97.55 \pm 0.08
Fixmatch + BaM	57.72 \pm 1.44	90.67 \pm 0.90	31.45 \pm 1.29	56.60 \pm 2.45	76.18 \pm 18.61	97.51 \pm 0.13
Fixmatch + Ours	89.72 \pm 5.31	92.69 \pm 0.74	67.44 \pm 1.14	69.10 \pm 0.45	96.58 \pm 0.14	96.54 \pm 0.13
Freematch	79.61 \pm 7.24	92.26 \pm 0.18	43.88 \pm 3.60	63.13 \pm 0.46	86.87 \pm 4.15	92.90 \pm 2.76
Freematch + MR	84.19 \pm 6.19	92.17 \pm 0.36	44.20 \pm 2.93	62.03 \pm 0.82	84.89 \pm 7.82	93.26 \pm 2.36
Freematch + BaM	88.33 \pm 0.33	92.32 \pm 0.25	44.49 \pm 3.05	62.13 \pm 2.93	81.43 \pm 14.80	91.08 \pm 3.72
Freematch + Ours	90.19 \pm 5.37	93.10 \pm 0.28	65.37 \pm 4.05	68.76 \pm 1.38	96.55 \pm 0.15	96.65 \pm 0.26
Softmatch	83.60 \pm 7.09	91.74 \pm 0.78	40.73 \pm 1.46	61.43 \pm 0.34	59.52 \pm 12.47	96.93 \pm 0.23
Softmatch + Ours	89.96 \pm 4.74	93.14 \pm 0.33	67.84 \pm 0.33	68.74 \pm 0.72	96.55 \pm 0.10	96.51 \pm 0.37
Adamatch	75.00 \pm 1.10	91.35 \pm 0.66	31.61 \pm 1.92	58.08 \pm 0.44	80.83 \pm 7.09	96.99 \pm 0.08
Adamatch + Ours	86.62 \pm 10.54	93.06 \pm 0.19	67.00 \pm 1.02	68.12 \pm 0.48	96.55 \pm 0.09	96.44 \pm 0.20

Table 2. Top-1 Accuracy for CIFAR-10, CIFAR-100 and SVHN averaged across 3 random seeds. We report the average accuracy \pm std. deviation. Here, $N_l = \{40, 250\}$ for CIFAR-10, SVHN and $N_l = \{400, 2500\}$ for CIFAR-100. The best accuracy is **bolded**.

	CIFAR-10 Top-1 Accuracy	CIFAR-100 Top-1 Accuracy
Schedule 1 Fixmatch + Ours	91.11 \pm 1.31	65.78 \pm 1.36
Schedule 1 Freematch + Ours	91.09 \pm 1.01	66.13 \pm 0.48
Schedule 2 Fixmatch + Ours	35.38 \pm 29.27	19.56 \pm 3.07
Schedule 2 Freematch + Ours	30.48 \pm 10.52	24.80 \pm 3.49

Table 3. Top-1 accuracy for the two error tolerance (ϵ) scheduling. The table reports the ϵ we use between each iteration interval and the top-1 accuracy yielded by the corresponding schedule. For schedule 1, we set $\epsilon = 0.4, 0.2, 0.1, 0.05, 0.01$ when the number of iterations is in the following interval, respectively: $[0, 5K)$, $[5K, 10K)$, $[10K, 15K)$, $[15K, 20K)$, $[20K, 25K)$. For schedule 2, we set $\epsilon = 0.01, 0.05, 0.1, 0.2, 0.4$, for the same intervals, respectively.

framework provides the flexibility to control this explicitly and thus can be tuned by practitioners for the setting at hand.

C3. Is pseudolabel accumulation helpful? Accumulation allows the methods to use old pseudolabels for points that couldn’t get pseudolabeled in the current iteration. Thus, we expect accumulation to help improve the utilization of unlabeled data and 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 and without accumulation and report the results in Table 4. We observe that our method has similar test accuracy irrespective of accumulation. However, accumulation achieves better coverage in early iterations, as observed

Method	Acc—True	Acc—False
Fixmatch	67.62 \pm 2.10	90.08 \pm 0.78
Fixmatch + MR	64.78 \pm 4.64	90.41 \pm 0.83
Fixmatch + BaM	68.10 \pm 2.02	90.67 \pm 0.90
Freematch	85.40 \pm 1.36	92.26 \pm 0.18
Freematch + MR	83.59 \pm 2.59	92.17 \pm 0.36
Freematch + BaM	85.48 \pm 3.02	92.32 \pm 0.25
Fixmatch + Ours	92.69 \pm 0.74	92.80 \pm 0.56
Freematch + Ours	93.10 \pm 0.28	91.80 \pm 1.08

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

in Figure 6 in Appendix C. These results are unsurprising since our method ensures high quality of pseudolabels while maximizing coverage; it can 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 pseudolabels generated by the baseline models are inaccurate, especially in the earlier iterations, thus degrading the overall performance. Overall, we believe accumulation will be helpful when we have pseudolabels with high accuracy. The plots for pseudolabeling coverage and accuracy over the entire run are in Figures 6, 7 in Appendix C.

Experiment with baselines using calibration data for training. We run the baselines where the amount of training data for the baseline is increased by the amount of calibration data used in our method. For the CIFAR-10 setting with 250 labels, we run the baseline now with $250 + 1000 = 1250$ labeled points for training. Similarly, for CIFAR-100 with 2500 labels setting, we run with $2500 + 3000 = 5500$

Method	CIFAR-10	CIFAR-100
FixMatch (train + cal)	92.68 \pm 0.31	64.77 \pm 0.10
FixMatch + Ours	92.69 \pm 0.74	69.10 \pm 0.45
FreeMatch (train + cal)	93.03 \pm 0.03	67.69 \pm 0.12
FreeMatch + Ours	93.10 \pm 0.28	68.76 \pm 1.38

Table 5. Comparison of FixMatch and FreeMatch with and without our method on CIFAR-10 and CIFAR-100. Here $N_l = 250$ for CIFAR-10 and 2500 for CIFAR-100. The overall training data from methods annotated as (train + cal) includes these N_l points and the calibration data N_{cal} (from Table 1) used in our method.

labels for training. The results are reported in the Table 5. We can see that even with more labeled data in training, the baselines still fall short significantly in the CIFAR-100 setting, while the performance gap in the CIFAR-10 (easier) setting narrows down.

Validation data usage, limitations, and future work. In the experiments, our method used parts of the validation data differently in comparison to the baselines. We comment on the role of validation data in the baselines, in our method, and potential limitations due to it. Recall, our goal is to address the problem of ad hoc choices of confidence scores and thresholds in pseudolabeling-based SSL. To this end, we introduced a principled solution to learn scores and thresholds that can directly achieve any specified pseudolabeling error while maximizing the number of pseudolabeled points. Given the focus of our paper, our experiments are designed to study whether using our learnable scores and thresholds can benefit the baselines. To keep our solution statistically sound, we used part of the validation data to learn the scores and thresholds. The lack (or cost) of a sufficient amount of validation data could be a limitation in practice. Future work can explore ways to reduce the amount of validation data needed for learning confidence scores and thresholds while preserving the soundness of our method. For instance, generative models or data augmentation techniques could be employed to this end.

More broadly, in semi-supervised learning (SSL) research, the cost of validation data is often not explicitly accounted for, partly due to the widespread use of standard benchmark datasets where hyperparameters have been tuned extensively over time. However, in real-world applications, new datasets typically require a substantial amount of validation data for model selection and hyperparameter tuning—both of which are essential parts of the modern training process. Developing novel benchmarks that explicitly incorporate validation data into the overall labeled data budget could help better reflect practical deployment settings and would be a valuable direction for future SSL research.

6. Conclusion

Common semi-supervised learning (SSL) methods rely on pseudolabeling, but their effectiveness is limited by unreliable confidence scores and heuristic thresholding strategies. We address these issues by introducing a principled framework for learning confidence scores and thresholds with explicit control over pseudolabeling error. We adapt existing SSL methods with this framework and empirically show that the adapted methods achieve a higher test accuracy compared to their standard versions. Additionally, we introduce pseudolabel accumulation and analyze its impact, showing that it benefits methods with reliable pseudolabels, such as those using our framework. In sum, by providing a principled, data-driven approach to obtaining scores and thresholds for pseudolabeling, our work enhances SSL methods and opens the door to more reliable and efficient pseudolabeling-based SSL.

Impact Statement

This research improves semi-supervised learning, enabling more accurate and efficient machine learning in settings where labeled data is hard to obtain by following first principles in designing thresholds and confidence functions. Our work has various potential societal implications, with no specific concerns that require special attention in this context.

Acknowledgments

This work was partly supported by funding from the American Family Data Science Institute and the Institute for Foundations of Data Science (IFDS). We thank the anonymous reviewers for their valuable feedback.

References

- Adams, R. P. and Ghahramani, Z. Archipelago: nonparametric bayesian semi-supervised learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pp. 1–8, 2009.
- Amini, M.-R., Feofanov, V., Pauletto, L., Hadjadj, L., Devijver, E., and Maximov, Y. Self-training: A survey, 2023.
- Arazo, E., Ortego, D., Albert, P., O’Connor, N. E., and McGuinness, K. Pseudo-labeling and confirmation bias in deep semi-supervised learning. In *2020 International joint conference on neural networks (IJCNN)*, pp. 1–8. IEEE, 2020.
- Bachman, P., Alsharif, O., and Precup, D. Learning with pseudo-ensembles. In *Advances in Neural Information Processing Systems*, volume 27, 2014.

- Blum, A. and Chawla, S. Learning from labeled and unlabeled data using graph mincuts. 2001.
- Blum, A. and Mitchell, T. Combining labeled and unlabeled data with co-training. In *Proc. of the eleventh annual conference on Computational learning theory*, pp. 92–100. ACM, 1998.
- Chapelle, O., Schölkopf, B., and Zien, A. (eds.). *Semi-Supervised Learning*. The MIT Press, 2006. ISBN 9780262033589.
- Chen, H., Tao, R., Fan, Y., Wang, Y., Wang, J., Schiele, B., Xie, X., Raj, B., and Savvides, M. 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>.
- Corbière, C., THOME, N., Bar-Hen, A., Cord, M., and Pérez, P. Addressing failure prediction by learning model confidence. In *Advances in Neural Information Processing Systems* 32, pp. 2902–2913. 2019.
- Djurisic, A., Bozanic, N., Ashok, A., and Liu, R. Extremely simple activation shaping for out-of-distribution detection. In *The Eleventh International Conference on Learning Representations*, 2023.
- El-Manzalawy, Y., Munoz, E. E., Lindner, S. E., and Honavar, V. Plasmosep: Predicting surface-exposed proteins on the malaria parasite using semisupervised self-training and expert-annotated data. *Proteomics*, 16(23): 2967–2976, 2016.
- Fan, Y., Kukleva, A., and Schiele, B. Revisiting consistency regularization for semi-supervised learning, 2021.
- Foret, P., Kleiner, A., Mobahi, H., and Neyshabur, B. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021.
- Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. On calibration of modern neural networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- Gupta, C. and Ramdas, A. Top-label calibration and multiclass-to-binary reductions. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=WqoBaaPHS->.
- Hendrycks, D. and Gimpel, K. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *International Conference on Learning Representations*, 2017.
- Hui, L., Belkin, M., and Wright, S. Cut your losses with squentropy. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 14114–14131, 2023.
- Joachims, T. Transductive inference for text classification using support vector machines. In Bratko, I. and Dzeroski, S. (eds.), *Proceedings of ICML-99, 16th International Conference on Machine Learning*, pp. 200–209, 1999.
- Kahn, J., Lee, A., and Hannun, A. Self-training for end-to-end speech recognition. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 7084–7088. IEEE, 2020.
- Karamanolakis, G., Mukherjee, S., Zheng, G., and Awadallah, A. H. Self-training with weak supervision. *arXiv preprint arXiv:2104.05514*, 2021.
- Kingma, D. P., Mohamed, S., Jimenez Rezende, D., and Welling, M. Semi-supervised learning with deep generative models. In *Advances in Neural Information Processing Systems*, volume 27, 2014.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Kukačka, J., Golkov, V., and Cremers, D. Regularization for deep learning: A taxonomy, 2017.
- Kull, M., Perello Nieto, M., Kängsepp, M., Silva Filho, T., Song, H., and Flach, P. Beyond temperature scaling: Obtaining well-calibrated multi-class probabilities with dirichlet calibration. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Kumar, A., Sarawagi, S., and Jain, U. 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*, pp. 2805–2814. PMLR, 10–15 Jul 2018.
- Kumar, A., Liang, P. S., and Ma, T. Verified uncertainty calibration. *Advances in Neural Information Processing Systems*, 32, 2019.
- Laine, S. and Aila, T. Temporal ensembling for semi-supervised learning. *Fifth International Conference on Learning Representations*, 2017.
- Lee, D.-H. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *ICML Workshop on Challenges in Representation Learning*, 2013.
- Li, M., Wu, R., Liu, H., Yu, J., Yang, X., Han, B., and Liu, T. Instant: Semi-supervised learning with instance-dependent thresholds. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

- Loh, C., Dangovski, R., Sudalairaj, S., Han, S., Han, L., Karlinsky, L., Soljagic, M., and Srivastava, A. On the importance of calibration in semi-supervised learning, 2022.
- Loh, C., Dangovski, R., Sudalairaj, S., Han, S., Han, L., Karlinsky, L., Soljagic, M., and Srivastava, A. Mitigating confirmation bias in semi-supervised learning via efficient bayesian model averaging. *Transactions on Machine Learning Research*, 2023.
- McLachlan, G. J. 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.
- Mishra, S., Murugesan, B., Ayed, I. B., Pedersoli, M., and Dolz, J. Do not trust what you trust: Miscalibration in semi-supervised learning, 2024.
- Moon, J., Kim, J., Shin, Y., and Hwang, S. Confidence-aware learning for deep neural networks. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pp. 7034–7044, 2020.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., Ng, A. Y., et al. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, pp. 7. Granada, Spain, 2011.
- Nguyen, A., Yosinski, J., and Clune, J. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 427–436, 2015.
- Nigam, K., McCallum, A. K., Thrun, S., and Mitchell, T. Text classification from labeled and unlabeled documents using em. *Machine learning*, 39:103–134, 2000.
- Niyogi, P. Manifold regularization and semi-supervised learning: Some theoretical analyses. *Journal of Machine Learning Research*, 14(5), 2013.
- Oliver, A., Odena, A., Raffel, C. A., Cubuk, E. D., and Goodfellow, I. Realistic evaluation of deep semi-supervised learning algorithms. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Oymak, S. and Gulcu, T. C. Statistical and algorithmic insights for semi-supervised learning with self-training. *arXiv preprint arXiv:2006.11006*, 2020.
- Rizve, M. N., Duarte, K., Rawat, Y. S., and Shah, M. In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning. In *International Conference on Learning Representations*, 2021.
- Roelofs, B., Berthelot, D., Sohn, K., Carlini, N., and Kurakin, A. Adamatch: A unified approach to semi-supervised learning and domain adaptation. In *The Tenth International Conference on Learning Representations*, 2022.
- Rosenberg, C., Hebert, M., and Schneidman, H. Semi-supervised self-training of object detection models. In *Seventh IEEE Workshops on Applications of Computer Vision (WACV/MOTION’05) - Volume 1*, volume 1, pp. 29–36, 2005. doi: 10.1109/ACVMOT.2005.107.
- Sajjadi, M., Javanmardi, M., and Tasdizen, T. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pp. 1171–1179, 2016.
- Scudder, H. Probability of error of some adaptive pattern-recognition machines. *IEEE Transactions on Information Theory*, 11(3):363–371, 1965.
- Singh, A., Nowak, R., and Zhu, J. Unlabeled data: Now it helps, now it doesn't. In *Advances in Neural Information Processing Systems*, volume 21. Curran Associates, Inc., 2008.
- Sohn, K., Berthelot, D., Carlini, N., Zhang, Z., Zhang, H., Raffel, C. A., Cubuk, E. D., Kurakin, A., and Li, C.-L. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608, 2020.
- Subramanya, A. and Talukdar, P. P. *Graph-based semi-supervised learning*. Springer Nature, 2022.
- van Engelen, J. E. and Hoos, H. H. A survey on semi-supervised learning. *Machine Learning*, 109:373 – 440, 2019.
- Vapnik, V. N., Vapnik, V., et al. Statistical learning theory. 1998.
- Vishwakarma, H., Lin, H., Sala, F., and Vinayak, R. K. Promises and pitfalls of threshold-based auto-labeling. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Vishwakarma, H., Tay, S. J., Namburi, S. S. S., Sala, F., Vinayak, R. K., et al. Pearls from pebbles: Improved confidence functions for auto-labeling. *arXiv preprint arXiv:2404.16188*, 2024.
- 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., Shinozaki, T., Schiele, B., Wang, J., Xie, X., and Zhang, Y. Usb: A unified semi-supervised learning benchmark

- for classification. In *Thirty-sixth Conference on Neural Information Processing Systems, Datasets and Benchmarks Track*, 2022.
- Wang, Y., Chen, H., Heng, Q., Hou, W., Fan, Y., Wu, Z., Wang, J., Savvides, M., Shinozaki, T., Raj, B., Schiele, B., and Xie, X. 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.
- Xie, Q., Dai, Z., Hovy, E., Luong, T., and Le, Q. Un-supervised data augmentation for consistency training. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- Xu, Y., Shang, L., Ye, J., Qian, Q., Li, Y.-F., Sun, B., Li, H., and Jin, R. Dash: Semi-supervised learning with dynamic thresholding. In *International Conference on Machine Learning*, pp. 11525–11536. PMLR, 2021.
- Zadrozny, B. and Elkan, C. Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 694–699, 2002.
- Zhang, B., Wang, Y., Hou, W., Wu, H., Wang, J., Okumura, M., and Shinozaki, T. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information Processing Systems*, 34: 18408–18419, 2021.
- Zheng, M., You, S., Huang, L., Wang, F., Qian, C., and Xu, C. Simmatch: Semi-supervised learning with similarity matching. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 14471–14481, 2022.
- Zhu, X. Semi-supervised learning literature survey. In *University of Wisconsin-Madison, Department of Computer Sciences*, 2005.

Supplementary Material

The supplementary material is organized as follows. First, we summarize the notations in Table 6 in Appendix A, then we provide formal algorithms in Appendix B and additional experimental results and details are provided in Appendix C.

A. Glossary

The notations are summarized in Table 6 below.

B. Detailed Algorithms

Algorithm 1 Estimate Pseudolabeling Thresholds Classwise

Input: Confidence function \hat{g}_i , classifier \hat{h}_i , Part of validation data $D_{\text{th}}^{(i)}$ for threshold estimation, pseudolabeling error tolerance ϵ , space of thresholds T , label space \mathcal{Y} .

Output: Pseudolabeling thresholds $\hat{\mathbf{t}}_i$, where $\hat{\mathbf{t}}_i[y]$ is the threshold for class y .

for $y \in \mathcal{Y}$ **do**

 # Extract the set of points $D_{\text{th}}^{(y)}$ for which the groundtruth class is y .

$D_{\text{th}}^{(y)} \leftarrow \{(\mathbf{x}', y') \in D_{\text{th}} : y' = y\}$

$T'_y \leftarrow T \cup \{\infty\}$.

 # Estimate pseudolabeling error at each threshold on class specific data $D_{\text{th}}^{(y)}$. Pick the smallest threshold with the sum of the estimated error and C_1 times the std. deviation is below ϵ . Here C_1 is set to 0.25 and $\hat{\sigma}(z) = \sqrt{z(1-z)}$.

$\hat{\mathbf{t}}_i[y] \leftarrow \min\{t \in T'_y : \hat{\mathcal{E}}(\hat{g}_i, t \mid \hat{h}_i, D_{\text{th}}^{(y)}) + C_1 \hat{\sigma}(\hat{\mathcal{E}}(\hat{g}_i, t \mid \hat{h}_i, D_{\text{th}}^{(y)})) \leq \epsilon\}$,

end for

return $\hat{\mathbf{t}}_i$

Algorithm 2 Estimate Pseudolabeling Threshold Jointly for All Classes

Input: Confidence function \hat{g}_i , classifier \hat{h}_i , Part of validation data $D_{\text{th}}^{(i)}$ for threshold estimation, pseudolabeling error tolerance ϵ , space of thresholds T , label space \mathcal{Y} .

Output: Pseudolabeling thresholds $\hat{\mathbf{t}}_i$, where $\hat{\mathbf{t}}_i[y]$ is the threshold for class y .

$T' \leftarrow T \cup \{\infty\}$

Estimate pseudolabeling error at each threshold on the entire set D_{th} . Pick the smallest threshold with the sum of the estimated error and C_1 times $\hat{\sigma}$ is below ϵ . Here C_1 is set to 0.25 and $\hat{\sigma}(z) = \sqrt{z(1-z)}$.

$t \leftarrow \min\{t \in T' : \hat{\mathcal{E}}(\hat{g}_i, t \mid \hat{h}_i, D_{\text{th}}) + C_1 \hat{\sigma}(\hat{\mathcal{E}}(\hat{g}_i, t \mid \hat{h}_i, D_{\text{th}})) \leq \epsilon\}$.

for $y \in \mathcal{Y}$ **do**

$\hat{\mathbf{t}}_i[y] \leftarrow t$

end for

return $\hat{\mathbf{t}}_i$

Algorithm 3 Pseudolabeling Based SSL with PabLO

Input: Labeled data for training D_l , validation data D_{val} , unlabeled pool X_u , error tolerance ϵ , use-accumulation flag, num_iters, batch size B , replication factor μ , weak ω and strong Ω augmentations, number of calibration points N_{cal} , num. of threshold estimation points N_{th} , frequency of invoking PabLO F , space of thresholds T , label space \mathcal{Y} .

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

Set initial pseudolabels and masks to 0.

$\tilde{Y}_u^{(0)} \leftarrow [0, 0, \dots, 0]$, $S_u^{(0)} \leftarrow [0, 0, \dots, 0]$, $i \leftarrow 1$

Draw calibration and threshold estimation sets from D_{val} .

$D_{\text{cal}}, D_{\text{th}} \leftarrow \text{DrawRandomly}(D_{\text{val}}, N_{\text{cal}}, N_{\text{th}})$.

Training loop with pseudolabeling.

while $i \leq \text{num_iters}$ **do**

Draw batches D_l^b, X_u^b of labeled and unlabeled points, I_u^b denotes the indices corresponding to points in X_u^b .

$D_l^b, X_u^b, I_u^b \leftarrow \text{DrawRandomBatch}(\mu D_l, \mu X_u, B)$

Create weak and strong augmentations of X_u^b .

$X_{u,w}^b, X_{u,s}^b \leftarrow \omega(X_u^b), \Omega(X_u^b)$

/ Begin Pseudolabeling Block **/**

Perform pseudolabeling using PabLO.

if $i \% F = 0$ **then**

Get \hat{g}_i by solving optimization (P1).

$\hat{g}_i, \hat{\mathbf{t}}_i \leftarrow \text{SolveOptProblemP1}(\hat{h}_i, D_{\text{cal}})$

Estimate pseudolabeling thresholds.

if estimate threshold classwise **then**

Use Algorithm 1.

$\hat{\mathbf{t}}_i \leftarrow \text{ClasswiseThreshold}(\hat{g}_i, \hat{h}_i, D_{\text{th}}, \epsilon, T, \mathcal{Y})$

else

Use Algorithm 2.

$\hat{\mathbf{t}}_i \leftarrow \text{JointThreshold}(\hat{g}_i, \hat{h}_i, D_{\text{th}}, \epsilon, T, \mathcal{Y})$

end if

Compute fresh psuedolabels $\tilde{Y}_u^{(i)}$ and pseudolabeling masks $S_u^{(i)}$ for all points in X_u .

$\tilde{Y}_u^{(i)} \leftarrow \hat{h}_i(\omega(X_u))$, $S_u^{(i)} \leftarrow \mathbb{1}(\hat{g}_i(\omega(X_u)) \geq \hat{\mathbf{t}})$

if use-accumulation **then**

Apply pseudolabel accumulation if enabled.

$\tilde{Y}_u^{(i)} \leftarrow S_u^{(i)} \tilde{Y}_u^{(i)} + (1 - S_u^{(i)}) \tilde{Y}_u^{(i-1)}$

$S_u^{(i)} \leftarrow S_u^{(i)} \vee S_u^{(i-1)}$

end if

else

$\hat{g}_i, \hat{\mathbf{t}}_i = \hat{g}_{i-1}, \hat{\mathbf{t}}_{i-1}$

end if

/ End Pseudolabeling Block **/**

Extract pseudolabels and masks for the current unlabeled batch. Then compute supervised and unsupervised losses.

$\tilde{Y}_u^b, S_u^b \leftarrow \tilde{Y}_u[I_u^b], S_u[I_u^b]$

$\hat{\mathcal{L}}_s(\hat{h}_i) \leftarrow \text{supervised_loss}(h, D_l^b)$

$\hat{\mathcal{L}}_u(\hat{h}_i) \leftarrow \text{unsupervised_loss}(h, X_{u,w}^b, X_{u,s}^b, \tilde{Y}_u^b, S_u^b)$

$\hat{\mathcal{L}}(\hat{h}_i) \leftarrow \hat{\mathcal{L}}_s(\hat{h}_i) + \lambda_u \hat{\mathcal{L}}_u(\hat{h}_i)$

Perform a gradient descent step to get new model \hat{h}_{i+1} .

$\hat{h}_{i+1} \leftarrow \text{SGD_update}(\hat{\mathcal{L}}(\hat{h}_i)); \quad i \leftarrow i + 1$

Evaluate model on D_{val} to keep track of the best model.

if $i \% \text{eval_freq} = 0$ **then**

$\text{eval_acc} \leftarrow \text{evaluate_model}(\hat{h}_i, D_{\text{val}})$

If eval_acc is best so far then $\hat{h}_{\text{ssl}} = \hat{h}_i$

end if

end while

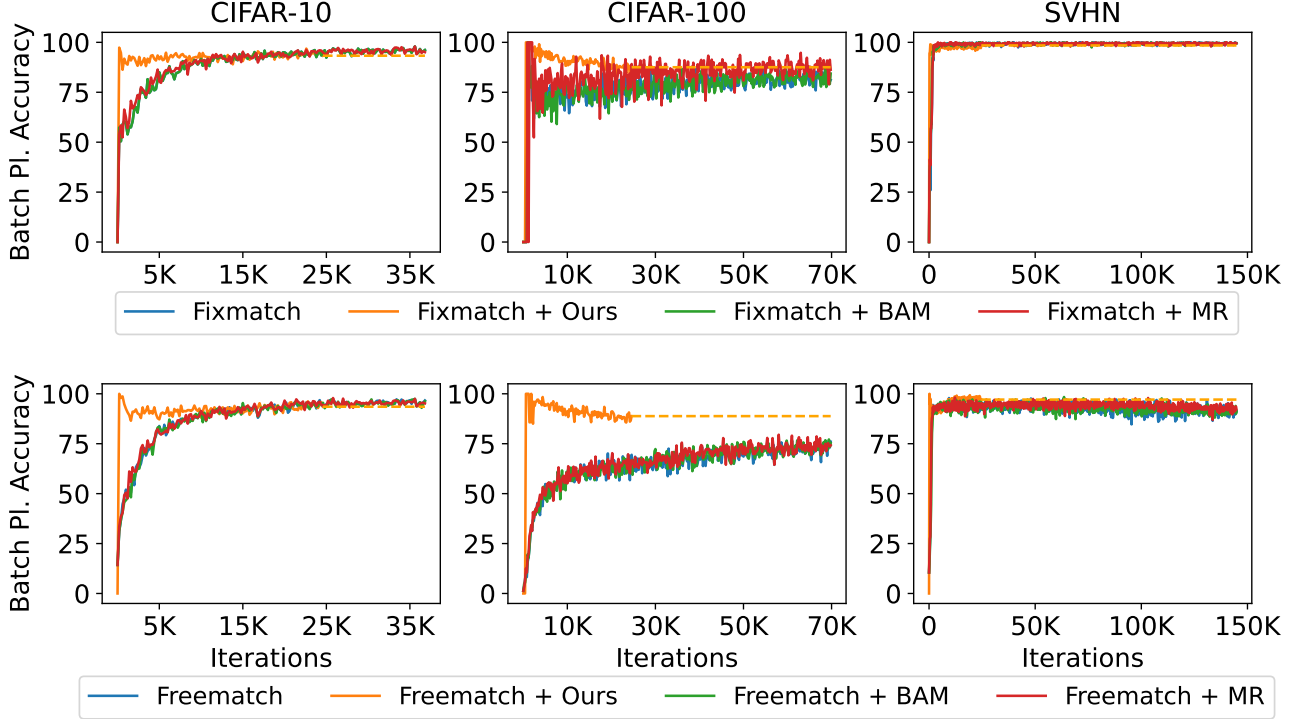


Figure 4. Batch pseudolabel accuracy of our method and baselines on CIFAR-10, CIFAR-100, and SVHN. We plot the values for every 200 steps.

C. Additional Experiments and Details

Compute. We ran all of our experiments on a high-throughput system with various GPUs. Therefore, each individual experiment task may be scheduled among NVIDIA A100 SXM4-40GB, NVIDIA A100 SXM4-80GB, NVIDIA L40, and NVIDIA H100 80GB HBM3. We measured the runtime of our algorithm on a desktop with a single NVIDIA RTX 4090. On CIFAR-10, it took about 0.203 seconds for each iteration for our method and around 0.140 seconds for the baselines. On CIFAR-100, it took about 0.396 seconds for each iteration for our method and around 0.143 seconds for the baselines. On SVHN, it took about 1.275 seconds for each iteration for our method and around 0.225 seconds for the baselines.

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 10.

Note that the number of epochs we used to train the function g and to estimate \mathbf{t} is dynamic. That is, its actual value depends on and is proportional to the current number of iterations of the SSL training. More concretely, at iteration i of SSL training, we use $\min(\lfloor i/25 \rfloor, \text{max epoch})$ number of epochs to find g and \mathbf{t} .

We additionally conduct the following ablation study to study our technique’s dependence on the amount of data used in learning g and thresholds.

A2. 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 N_{th} in $\{250, 500, 750, 1000\}$ on CIFAR-10 setting for 3 random seeds and report the result in Fig 8. 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 to Table 7 for more details.

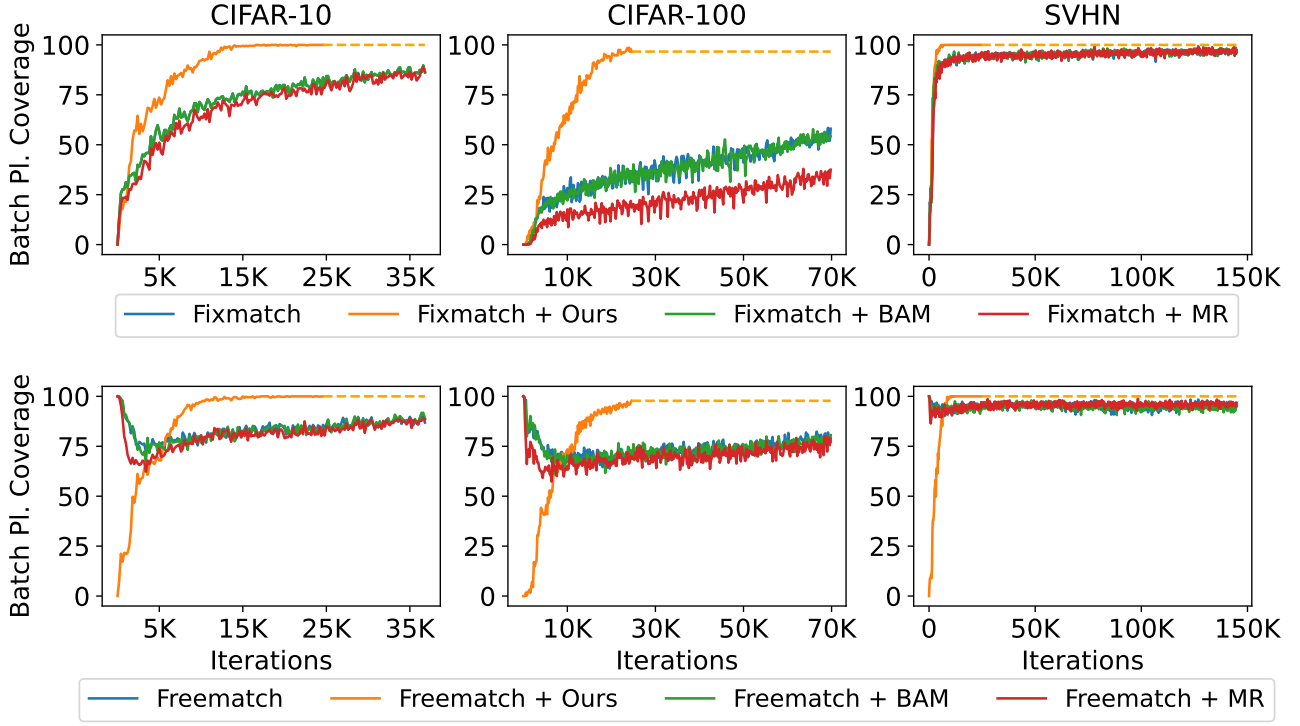


Figure 5. Batch pseudolabel coverage of our method and baselines on CIFAR-10, CIFAR-100, and SVHN. We plot the values for every 200 steps.

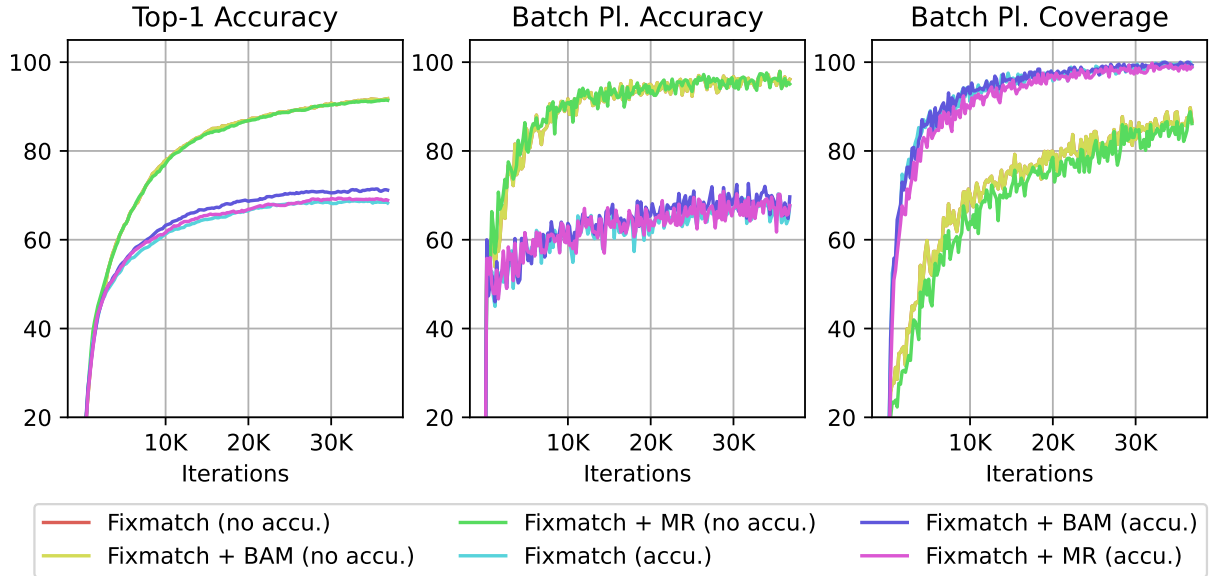


Figure 6. (A1.) Left to Right: Top-1 Accuracy, Batched pseudolabeling Accuracy, and batched pseudolabeling coverage of Fixmatch with and without pseudolabeling accumulation enabled on CIFAR-10. It can be seen that enabling pseudolabeling accumulation worsen the performance of baseline methods in terms of accuracy and coverage.

Symbol	Definition
\mathcal{X}	feature space.
\mathcal{Y}	label space i.e. $1, 2, \dots k$.
\mathcal{H}	hypothesis space (model class for the classifiers).
\mathcal{G}	space of confidence functions.
k	number of classes.
\mathbf{x}, y	\mathbf{x} is a datapoint in \mathcal{X} and y is its true label (if available).
h	a model $h : \mathcal{X} \rightarrow \mathcal{Y}$.
g	confidence function $g : \mathcal{X} \rightarrow T^k \subseteq \mathbb{R}^k$
\hat{y}	hard label prediction.
\tilde{y}	\tilde{y} is used as pseudolabel.
\hat{h}_{ssl}	a best learned model using SSL.
ϵ	pseudolabeling error tolerance.
g_i^*	optimal confidence function at i iteration.
\mathbf{t}_i^*	optimal threshold at i iteration.
X_u	available unlabeled data drawn from the distribution P_x over \mathcal{X} .
X_u^b	batch of unlabeled data.
D_l	set of labeled data points drawn from the distribution P_{xy} .
D_l^b	batch of labeled data.
D_{val}	validation data.
D_{cal}	calibration data; part of validation data used to optimize surrogate functions.
D_{th}	part of validation data to estimate threshold \mathbf{t} .
\mathbf{t}	k dimensional vector of thresholds representing for k classes.
$\mathbf{t}[y]$	y th entry of \mathbf{t} i.e. the threshold for class y .
n_u	number of unlabeled points, i.e. size of X_u used for consistency regularization and pseudolabeling.
N_l	number of labeled points, i.e. size of D_l . Usual SSL setting has, $N_l \ll n_u$.
N_{val}	number of points used for model selection.
N_{test}	number of test data points.
N_{cal}	number of points used for learning the g function.
N_{th}	number of data points used for threshold estimation.
$\hat{\mathcal{L}}_s$	supervised loss.
$\hat{\mathcal{L}}_u$	unsupervised loss with weighted importance λ_u .
$\hat{\mathcal{L}}_r$	sum of regularization terms for supervised and unsupervised loss with weighted importance λ_r .
$H(y, h, \mathbf{x})$	standard cross-entropy loss.
$S(\mathbf{x}, g, \mathbf{t} \mid h)$	pseudolabeling mask.
ω	weak transformation, $\omega : \mathcal{X} \mapsto \mathcal{X}$.
Ω	strong transformation, $\Omega : \mathcal{X} \mapsto \mathcal{X}$.
α_o, α_b	average time taken by our method and baseline methods.
$\hat{\mathcal{P}}(g, \mathbf{t} \mid h, X)$	estimated pseudolabeling coverage, see eq. (1).
$\mathcal{P}(g, \mathbf{t} \mid h)$	population level pseudolabeling coverage, see eq. (2).
$\hat{\mathcal{E}}(g, \mathbf{t} \mid h, D)$	estimated pseudolabeling error, see eq. (3).
$\mathcal{E}(g, \mathbf{t} \mid h)$	population level pseudolabeling error, see eq. (4).
$\tilde{\mathcal{P}}(g, \mathbf{t} \mid h, D)$	surrogate estimated pseudolabeling coverage, see eq. (6).
$\tilde{\mathcal{E}}(g, \mathbf{t} \mid h, D)$	surrogate estimated pseudolabeling error, see eq. (7).
λ	hyperparamter controlling the importance of pseudolabeling coverage and error in (P1).

Table 6. Glossary of variables and symbols used in this paper.

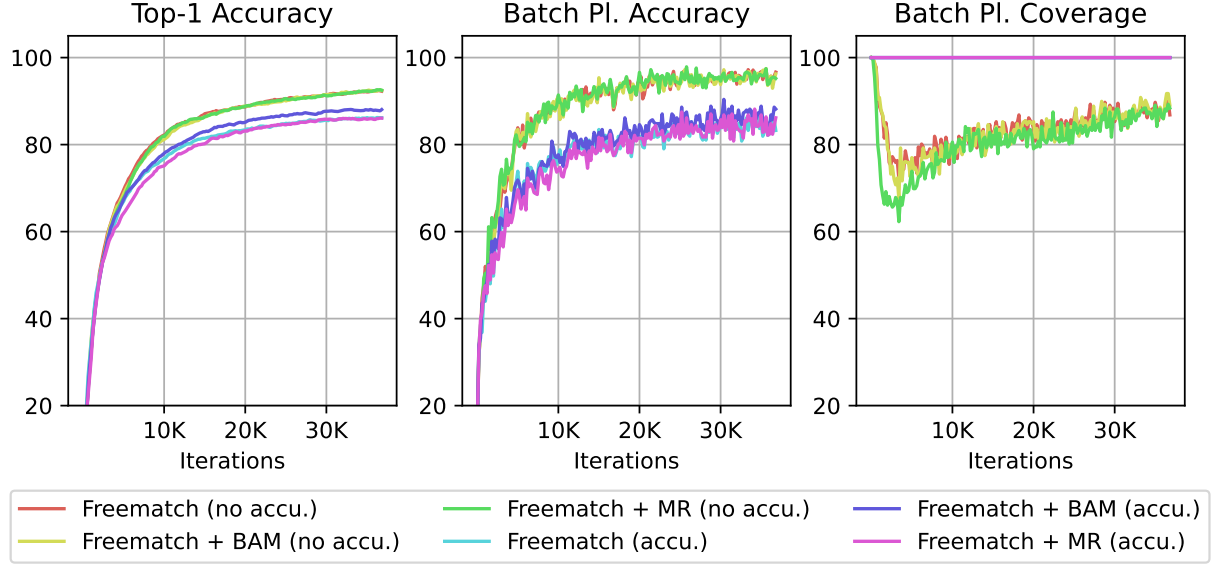


Figure 7. (A1.) Left to Right: Top-1 Accuracy, Batched pseudolabeling Accuracy, and batched pseudolabeling coverage of Freematch with and without pseudolabeling accumulation enabled on CIFAR-10. It can be seen that enabling pseudolabeling accumulation worsen the performance of baseline methods in terms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

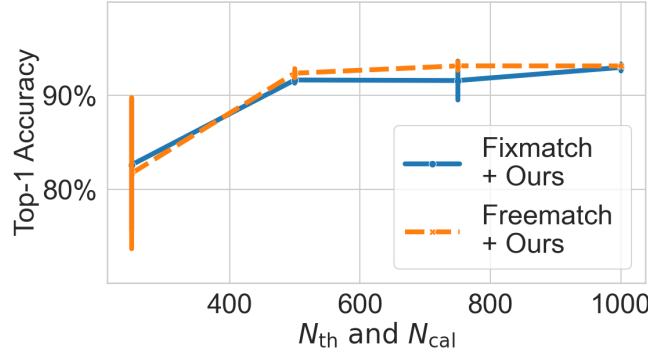


Figure 8. Top-1 accuracy of our method with different N_{th} and N_{cal} .

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

Table 7. Results on CIFAR-10 with varying N_{cal} and N_{th} .

Method	$\epsilon = 0.01$	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.4$
Fixmatch + Ours	93.05 ± 0.54	91.54 ± 0.95	88.35 ± 2.90	56.72 ± 22.25
Freematch + Ours	92.11 ± 1.18	92.31 ± 0.16	83.89 ± 10.36	52.17 ± 25.36

Table 8. Results on CIFAR-10 with varying ϵ .

Method	$\epsilon = 0.01$	$\epsilon = 0.1$	$\epsilon = 0.2$	$\epsilon = 0.4$
Fixmatch + Ours	69.19 ± 1.13	65.01 ± 0.34	53.88 ± 8.15	23.58 ± 18.21
Freematch + Ours	70.13 ± 0.67	64.95 ± 1.41	59.83 ± 1.32	24.09 ± 17.22

Table 9. Results on CIFAR-100 with varying ϵ .

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

Table 10. Hyperparameters used for our method.