

# DOTMATCH: SIMPLIFIED SEMI-SUPERVISED LEARNING WITH THE LOG DOT PRODUCT LOSS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Semi-supervised learning (SSL) algorithms typically work by generating supervisory signals for unsupervised data using the model being trained, but such supervisory signals are generally imperfect, thus various techniques have been proposed to balance the signal-to-noise ratio, such as confidence-based pseudo-labeling, consistency regularization and entropy regularization. However, these methods often require careful tuning of hyperparameters, such as the confidence threshold in pseudo-labeling and the regularization strength in regularization methods, which is often a challenging task, particularly with limited labeled data available for validation. In this paper, we introduce DotMatch, an SSL algorithm that is capable of balancing the signal-to-noise ratio without any algorithm specific hyperparameters. Specifically, we introduce a novel consistency loss on unsupervised data to replace the cross-entropy loss, called the log dot product (LDP) loss, which is simply the negative log of the dot product between the predicted label distributions of weak and strong augmented views of an input. We empirically and theoretically demonstrate that the LDP loss enjoys several benefits for SSL, compared to the cross-entropy loss with soft target: nonconfident examples have low impacts on model updates, as in confidence-based pseudo-labeling methods such as SoftMatch; predictions are encouraged to have a low entropy, as in entropy-regularized methods; and interestingly, its gradient is appropriately scaled relative to the gradient of the supervised loss, thus requiring no regularization constant. We additionally combine the LDP loss with distribution alignment to ensure the distribution of predictions on unlabeled data match that of the labeled data. Extensive experiments show that DotMatch is competitive with state-of-the-art baselines without needing to tune any algorithm-specific hyperparameters for different datasets. Code is available at: [open source upon acceptance].

## 1 INTRODUCTION

Deep learning has seen great success in many supervised learning problems (He et al., 2016; Vaswani et al., 2017; Dong et al., 2018), but requires lots of labeled training data which is often difficult to obtain. Semi-supervised learning (SSL) (Zhu, 2005; Zhu & Goldberg, 2009) eases this burden by enabling the use of supplemental unlabeled data which is much easier to obtain, for example, collect lots of images without manually annotating them. Many SSL algorithms use a form of self-training (McLachlan, 1975; Rosenberg et al., 2005; Lee et al., 2013; Xie et al., 2020b) where, during training, supervision for the unlabeled data is provided by the model itself. Since the model is only capable of providing inaccurate supervision, the goal is to maintain a high signal-to-noise ratio whereby we get abundant highly accurate supervision with minimal noise (similar to the quantity-quality tradeoff (Chen et al., 2023)).

Several techniques and their combinations are often utilized in SSL algorithms, such as (a) pseudo-labeling (Lee et al., 2013; Sohn et al., 2020), which generates labels for unlabeled data using the current model, then uses such pseudo-labels to further train the model; (b) consistency regularization (Bachman et al., 2014; Sajjadi et al., 2016b; Laine & Aila, 2017) which uses the model’s prediction after randomly perturbing the input or the model itself as a supervision signal; (c) entropy minimization (Grandvalet & Bengio, 2004a), which ensures that the model learns to make confident predictions on unlabeled data; (d) distribution alignment (Bridle et al., 1991; Berthelot et al., 2020), which encourages the model to match the average predicted label distribution on unlabeled examples with the distribution of observed labels. However, these methods often require

sophisticated balancing of signal-to-noise ratio. For example, many state-of-the-art SSL algorithms perform pseudo-labeling and propose carefully crafted pseudo-label confidence threshold schemes to improve performance (Zhang et al., 2021; Wang et al., 2023; Chen et al., 2023; Li et al., 2023) rather than using a constant threshold (Sohn et al., 2020). In addition, existing methods often require careful tuning of hyperparameters, such as the regularization strength in regularization methods.

We propose a SSL algorithm, DotMatch, that provides a simple approach to control the signal-to-noise ratio, and avoids the challenging task of tuning algorithm-specific hyperparameters as it does not have such hyperparameters. DotMatch uses a new loss function, called the log dot product (LDP) loss, which serves as an alternative generalization of the CE loss with a single label to soft weighted labels. Instead of considering the dot product between the soft target distribution and the logarithm of the predicted label distribution as is done with the CE loss, the LDP loss instead computes the logarithm of the dot product between the two distributions. As we shall demonstrate empirically and theoretically in Section 4, this approach has several advantages in the context of SSL: low-confidence examples naturally contribute less to the gradient update without requiring additional confidence-based loss weights as is typically used with pseudo-labeling (Sohn et al., 2020); the model will automatically learn to produce high confidence predictions (Grandvalet & Bengio, 2004a; Sajjadi et al., 2016a) without needing to rely on inaccurate one-hot pseudo-labels, additional entropy regularization which requires tuning of the regularization strength, or target sharpening which requires tuning the sharpening hyperparameter; and the gradients are appropriately scaled relative to the gradients of the supervised CE loss, thus no additional regularization constant is necessary. We combine the LDP loss with advanced data augmentation techniques (Xie et al., 2020a) in DotMatch to encourage the model to make the same predictions for different augmented versions of the same input. Finally, distribution alignment (Berthelot et al., 2020) is applied to the targets in the LDP consistency regularizer to ensure that the average predicted label distribution on unlabeled examples matches that of the labeled data, and does not collapse into a single class in order to achieve high confidence and consistent predictions. DotMatch demonstrates strong performance in extensive experiments. Notably, this is achieved without the need to tune algorithm specific hyperparameters such as confidence threshold or regularization strength, allowing it to perform well on new problems without requiring extra labeled data for validation.

Our contributions are as follows:

- We provide a motivating example to show explicitly that down-weighting the loss on low-confidence unlabeled examples leads to improved gradient quality in SSL, and in turn improves the predictive performance of learned models.
- We introduce the novel log dot product loss for consistency regularization which naturally reduces the influence of low-confidence examples on the model updates without requiring additional confidence-based loss weights, encourages the model to produce confident predictions on unlabeled data without introducing any hyperparameters that require tuning, and does not require any regularization constant due to the scale of the gradient matching that of the supervised CE loss.
- We further propose an SSL algorithm, DotMatch, which uses the LDP loss for consistency regularization combined with distribution alignment to encourage the distribution of predictions on unlabeled examples to match that of the labeled examples.
- A novel theoretical analysis of SSL consistency loss functions is performed to explain their efficacy, and explore the relationship between loss gradients and confidence.
- Extensive experiments on SSL benchmark problems demonstrate that DotMatch is competitive with state-of-the-art baselines despite its simplicity. The ablation study shows that the components of DotMatch positively contribute to performance.

## 2 RELATED WORK

Our work is closely related to many ideas and algorithms from the SSL literature. We present some of these in this section, with broader introductions to SSL available in, for example, (Zhu, 2005; van Engelen & Hoos, 2020).

Consistency regularization (Bachman et al., 2014; Sajjadi et al., 2016b; Laine & Aila, 2017) is commonly used in SSL algorithms as a natural way to incorporate unlabeled data into the training

process of a supervised learning model. At a high level, the goal is to encourage the model to produce similar outputs for similar inputs. One way to achieve this is by training the model to make the same prediction on different augmented versions of the same input. This approach is particularly effective when using high quality data augmentation strategies (Xie et al., 2020a). Specifically, the model is encouraged to match the prediction on a strongly augmented input with that of a weakly augmented version of the same input. Some SSL methods have proposed alternative measures of consistency, such as encouraging consistency between different versions of the model (Chen et al., 2022; Huang et al., 2023), unlabeled and labeled datasets (Huang et al., 2024; Heidari et al., 2024), a subset of the class predictions (Wu & Cui, 2024; Yang et al., 2023), examples from the same cluster (Liu et al., 2025), adding inconsistency regularization between different examples (Deng et al., 2024), or using ideas from optimal transport to measure consistency (Tan et al., 2024). DotMatch continues the work along this direction by proposing a new consistency loss.

Entropy regularization (Grandvalet & Bengio, 2004a) is another important idea in SSL that encourages the model to produce high-confidence predictions on unlabeled examples. Combining entropy regularization with consistency regularization ensures that the model does not converge to undesirable local minima where perfect consistency may be achieved by predicting labels for all examples uniformly at random. Target sharpening (Berthelot et al., 2019; Xie et al., 2020a; Berthelot et al., 2020) is one method that has been proposed to incorporate entropy minimization into SSL. Another technique is pseudo-labeling (Lee et al., 2013) with one-hot pseudo-labels. Since pseudo-labels generated by the model are generally inaccurate, especially early in training, it is common to only utilize pseudo-labels that the model is confident about. Some confidence threshold schemes include using a constant threshold (Sohn et al., 2020), dynamic thresholds which either depend on the pseudo-label class (Zhang et al., 2021; Wang et al., 2023) or not (Chen et al., 2023), thresholding an alternative measure of confidence (Min et al., 2024), or using auxiliary models to measure confidence or provide thresholds (Li et al., 2023; 2024; Fang et al., 2024). Since the LDP loss naturally minimizes the entropy of predictions and focuses more on highly confident examples, we do not use any target sharpening, one hot pseudo-labeling, confidence thresholding or auxiliary models in DotMatch.

Distribution Alignment (DA) (Berthelot et al., 2020; 2022) modifies the pseudo-labels such that the average pseudo-label distribution aligns better with the labeled data label distribution. This effect can also be achieved with loss regularization (Bridle et al., 1991; Wang et al., 2023). DA is particularly beneficial when the number of classes is large (Sohn et al., 2020), as the model predictions are susceptible to collapse into a single class in order to achieve consistent and low entropy predictions. We use the former type of DA in DotMatch as it does not introduce any additional regularization hyperparameters that require tuning.

Objective functions similar to the LDP-based consistency loss were proposed in (Hsu et al., 2019; Cao et al., 2022), except different problem settings were considered in both works so the purpose of the objectives is different, they did not incorporate weak-strong augmentations for consistency or DA, and we provide additional insights into the properties of the loss through an extensive analysis of the gradients.

### 3 PROBLEM SETTING

Let  $\mathbf{x} \in \mathbb{R}^d$  denote an input and  $\mathbf{y} \in \{0, 1\}^K \cap \Delta$  be a one-hot class label, where  $\Delta = \{\mathbf{q} \in [0, 1]^K : \mathbf{q}^\top \mathbf{1} = 1\}$  is the standard  $K - 1$  simplex. Following standard practice (Zhang et al., 2021; Wang et al., 2022), we consider the setting where data comes in mini-batches of labeled and unlabeled examples having data augmentation applied. In each training iteration we receive a batch of labeled examples  $\mathcal{L} = \{(\mathbf{w}(\mathbf{x}_i), \mathbf{y}_i)\}_{i=1}^{n_L}$  with each  $(\mathbf{x}_i, \mathbf{y}_i)$  distributed according to some unknown joint distribution  $p(\mathbf{x}, \mathbf{y})$ , and a batch of unlabeled examples  $\mathcal{U} = \{(\mathbf{w}(\mathbf{x}_i), \mathbf{s}(\mathbf{x}_i))\}_{i=1}^{n_U}$  with each  $\mathbf{x}_i$  being drawn from the marginal distribution over inputs  $p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y})$ . The weak ( $\mathbf{w}$ ) and strong ( $\mathbf{s}$ ) augmentation functions  $\mathbf{w}, \mathbf{s} : \mathbb{R}^d \rightarrow \mathbb{R}^d$  are stochastic. Some example weak augmentations are random crop, flip, shift, and strong augmentations include CutOut (DeVries & Taylor, 2017), CTAugmnet (Berthelot et al., 2020), RandAugment (Cubuk et al., 2020), or data modality agnostic augmentations such as MixUp (Zhang et al., 2018) or adversarial perturbations (Miyato et al., 2019). For the model we use a neural network  $p(\cdot; \theta) : \mathbb{R}^d \rightarrow \Delta$  with softmax output layer, and learnable parameters collected in the vector  $\theta$ . We consider SSL objective functions of

the form

$$L(\theta) = \frac{1}{|\mathcal{L}|} \sum_{(\mathbf{w}, \mathbf{y}) \in \mathcal{L}} \ell_l(\mathbf{p}(\mathbf{w}; \theta), \mathbf{y}) + \frac{\lambda}{|\mathcal{U}|} \sum_{(\mathbf{w}, \mathbf{s}) \in \mathcal{U}} w(\max(\mathbf{p}_w)) \ell_u(\mathbf{p}(\mathbf{s}; \theta), \mathbf{p}_w) \quad (1)$$

for some classification loss functions  $\ell_l, \ell_u : \Delta \times \Delta \rightarrow \mathbb{R}_{\geq 0}$ , and non-decreasing weight function  $w : [1/K, 1] \rightarrow [0, 1]$ . We abuse the notations  $\mathbf{w}$  and  $\mathbf{s}$  to denote a weak and a strong augmented example respectively, and we use  $\mathbf{p}_x = \mathbf{p}(\mathbf{x}; \theta)$  to denote the prediction on example  $\mathbf{x}$ , which is not used in gradient computation. The unlabeled loss weight  $\lambda$  is typically set to 1. Many existing SSL algorithms have objectives of this form, for example, (Sohn et al., 2020; Zhang et al., 2021; Chen et al., 2023; Wang et al., 2023). At test time, the predicted label for new example  $\mathbf{x} \in \mathbb{R}^d$  is given by  $\operatorname{argmax}(\mathbf{p}_x)$ . Throughout the paper we use  $\mathbf{e}_i \in \{0, 1\}^K \cap \Delta$  to denote the one-hot label for class  $i$ .

## 4 DOTMATCH

In this section we first use a motivating example to highlight an important SSL algorithm design consideration: down-weighting low-confidence unlabeled examples in the loss, which leads to improved gradient quality and, in turn, improved predictive performance. We then introduce our new LDP loss function and SSL algorithm, DotMatch, and give theoretical justification for their efficacy.

### 4.1 A MOTIVATING EXAMPLE

A desirable quality of an SSL algorithm is to have unlabeled examples with low-confidence pseudo-labels down-weighted in the computation of the loss to reduce the impact of confirmation bias during training, since higher confidence pseudo-labels tend to be correct more often than low-confidence ones. This property can equivalently be viewed through the lens of gradient-based learning: we want the contribution to the gradient from low-confidence examples to be smaller than that of high-confidence examples, so that the model update will be less influenced by low-confidence examples. Many recent SSL algorithms that weight the loss for each unlabeled example based on the prediction confidence have this property.

We compare one representative algorithm, SoftMatch (Chen et al., 2023), against a simple SSL algorithm that replaces the SoftMatch weighting function with the constant function  $w(\alpha) = 1$  for each unlabeled example, which we refer to as CE(hard), since  $\ell_u$  in this case is the CE loss with hard one-hot pseudo-label. We train each algorithm separately on EMNIST (Cohen et al., 2017) with 4 labeled examples per class. For a detailed description of the experiment settings, we kindly refer you to [Section 5](#). After 5000 training iterations, before training is finished, we record the target confidence and gradient norm

$$(\max(\mathbf{p}_w), \|\nabla_{\theta} w(\max(\mathbf{p}_w)) \ell_u(\mathbf{p}(\mathbf{s}; \theta), \mathbf{p}_w)\|_2)$$

for each unlabeled example in the batch. We also look at the quality of the gradients associated with unlabeled data as measured by the proportion of total gradient norm coming from correctly pseudo-labeled examples

$$\text{Quality}(w, \ell_u) = \frac{\frac{1}{|\mathcal{U}|} \sum_{(\mathbf{w}, \mathbf{s}, \mathbf{y}) \in \mathcal{U}} \mathbb{1}(\mathbf{y} = \mathbf{e}_{\operatorname{argmax}(\mathbf{p}_w)}) \|\nabla_{\theta} w(\max(\mathbf{p}_w)) \ell_u(\mathbf{p}(\mathbf{s}; \theta), \mathbf{p}_w)\|_2}{\frac{1}{|\mathcal{U}|} \sum_{(\mathbf{w}, \mathbf{s}, \mathbf{y}) \in \mathcal{U}} \|\nabla_{\theta} w(\max(\mathbf{p}_w)) \ell_u(\mathbf{p}(\mathbf{s}; \theta), \mathbf{p}_w)\|_2}. \quad (2)$$

We further report the test accuracy after training is completed at 50 000 iterations. The quality measure in [Equation \(2\)](#) is similar in spirit to the quality of pseudo-labels from (Chen et al., 2023), but applied to gradient norms. Note that the gradient quality is not computable in general SSL problems due to the use of ground truth labels for the unlabeled examples, and is only used here for illustrative purposes. Results are reported in [Figure 1](#).

We learn a few things from the results. Firstly, incorrect pseudo-labels tend to have lower confidence than correct ones for both methods. SoftMatch reduces the influence of incorrectly pseudo-labeled examples on the gradient by explicitly down-weighting the contribution of low-confidence pseudo-labels, whereas CE(hard) has no mechanism for doing this and simply treats all unlabeled examples equally. This results in SoftMatch having higher quality gradients compared to CE(hard). Overall, down-weighting the gradient norm of low-confidence unlabeled examples reduces confirmation bias from incorrect pseudo-labels during training, in turn improving performance of the final model.

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

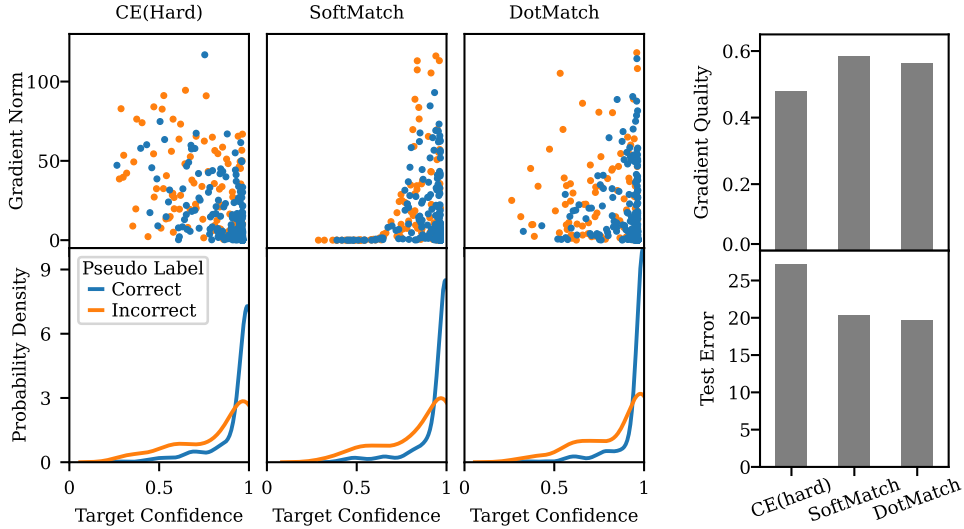


Figure 1: Comparison of the unlabeled example gradients and performance of CE(hard), SoftMatch, and our algorithm DotMatch, which is motivated by the differences between CE(hard) and SoftMatch. Left: Gradient norm (top) and density (bottom) vs confidence. CE(hard) has most large-gradient unconfident examples, SoftMatch least, and DotMatch in-between. Right: Gradient quality (top) and test error (bottom). SoftMatch and DotMatch both higher gradient quality and lower test error compared to CE(hard). Results recorded after 5000 training iterations on EMNIST (Cohen et al., 2017) with 4 labels per class, except test error, which was recorded after completing 50k training iterations.

The LDP loss used in our DotMatch algorithm is motivated by the above observation on SoftMatch’s ability to down-weight low-confidence examples in learning, and Figure 1 demonstrates DotMatch’s ability to do so too. In the next section we introduce the LDP loss and provides some theoretical insights on its ability to down-weight the gradient norm of low-confidence examples, which in turn produces high quality gradients and better performance of the final learned model.

#### 4.2 LOG DOT PRODUCT LOSS

The de facto unlabeled loss function  $l_u$  is the CE loss with a single pseudo-label generated from the model. However, this ignores the full information about the predicted label distribution. One way to generalize the CE loss with a single label  $\ell(\mathbf{q}, \mathbf{e}_y) = -\log q_y$  for some fixed  $y \in \{1, \dots, K\}$ ,  $\mathbf{q} \in \Delta$ , to multiple weighted labels  $\mathbf{p} \in \Delta$  is  $\ell(\mathbf{q}, \mathbf{p}) = -\mathbf{p}^\top \log \mathbf{q}$ . We refer to this loss as CE(soft). However, it is not the only generalization. Another option is to look at the negative log of the dot product between distributions  $\ell(\mathbf{q}, \mathbf{p}) = -\log(\mathbf{p}^\top \mathbf{q})$ . We refer to this loss function as the log dot product loss. This loss has many interesting properties that we highlight below.

**Consistency and Entropy Minimization** A prediction  $\mathbf{q}$  that minimizes the LDP loss with target  $\mathbf{p}$  is one that puts all of its probability mass on the set of maximizers of the target

$$\operatorname{argmin}_{\mathbf{q} \in \Delta} -\log(\mathbf{p}^\top \mathbf{q}) = \{\mathbf{q} \in \Delta : q_k = 0 \forall k \notin \operatorname{argmax}(\mathbf{p})\}.$$

This occurs when the dot product inside the logarithm achieves its maximum value. In the common situation where there is a unique maximizer in  $\mathbf{p}$ , the argmin is the singleton containing  $\mathbf{e}_{\operatorname{argmax}(\mathbf{p})}$ . In this common situation, the LDP loss is minimized when the target and prediction are both equal and have zero entropy. Compare with CE(soft) loss:

$$\operatorname{argmin}_{\mathbf{q} \in \Delta} -\mathbf{p}^\top \log(\mathbf{q}) = \mathbf{p},$$

which only encourages low entropy predictions if the target  $\mathbf{p}$  has low entropy (for example, one-hot pseudo-label). LDP loss has an entropy minimization effect without the need for one-hot targets, temperature scaled softmax targets, or additional entropy regularization loss term.

Table 1: Minimizers and gradients of different loss functions, where  $\mathbf{p}$  denotes the target distribution,  $\hat{y}$  the index of  $\mathbf{p}$  with largest value,  $\mathbf{z}$  the logits,  $\boldsymbol{\theta}$  the model parameters,  $\mathbf{q}$  the predicted distribution (either as a function of  $\mathbf{z}$  or  $\boldsymbol{\theta}$ ), and  $\odot$  the element-wise product. Note that  $\mathbf{q} = \text{softmax}(\mathbf{z})$ .

Name	$\ell(\mathbf{q}, \mathbf{p})$	$\text{argmin}_{\mathbf{q}} \ell(\mathbf{q}, \mathbf{p})$	$\nabla_{\mathbf{z}} \ell(\mathbf{q}(\mathbf{z}), \mathbf{p})$	$\nabla_{\boldsymbol{\theta}} \ell(\mathbf{q}(\boldsymbol{\theta}), \mathbf{p})$
CE(soft)	$-\mathbf{p}^\top \log \mathbf{q}$	$\mathbf{p}$	$\mathbf{q} - \mathbf{p}$	$-\sum_{k=1}^K \frac{p_k}{q_k(\boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} q_k(\boldsymbol{\theta})$
CE(hard)	$-e_{\hat{y}}^\top \log \mathbf{q}$	$e_{\hat{y}}$	$\mathbf{q} - e_{\hat{y}}$	$-\frac{1}{q_{\hat{y}}(\boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} q_{\hat{y}}(\boldsymbol{\theta})$
LDP	$-\log(\mathbf{p}^\top \mathbf{q})$	$e_{\hat{y}}$	$\mathbf{q} - \frac{\mathbf{p} \odot \mathbf{q}}{\mathbf{p}^\top \mathbf{q}}$	$-\sum_{k=1}^K \frac{p_k}{\mathbf{p}^\top \mathbf{q}(\boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} q_k(\boldsymbol{\theta})$

**Relationship between confidence and gradient norm** The loss minimizers and gradients with respect to both logits  $\mathbf{z} \in \mathbb{R}^K$  and  $\boldsymbol{\theta}$  are provided in Table 1, with derivations in Appendix A. Looking at the form of the gradients with respect to  $\boldsymbol{\theta}$  in Table 1, we see that the LDP loss implicitly scales the learning rate in an instance-dependent way similarly as the CE loss variants. The scaling factor is smallest when the prediction and target are equal to the same one-hot vector  $\mathbf{p} = \mathbf{q}(\boldsymbol{\theta}) = e_k$  for any  $k \in \{1, \dots, K\}$ , and grows larger as  $\mathbf{p}$  and  $\mathbf{q}$  become more orthogonal (have non-overlapping support). For the special case  $\mathbf{p} = e_k$  for some  $k \in \{1, \dots, K\}$ , all three loss functions coincide with loss gradient  $\mathbf{q} - e_k$ . For CE(soft) loss it is relatively well known that we get zero gradient only when  $\mathbf{q} = \mathbf{p}$ . The minimizer and maximizer of the LDP gradient norm are given in Theorem 1, which tells us that the LDP loss implicitly focuses more on examples with high-confidence and ignores examples with low-confidence. All proofs are in the Appendix.

**Theorem 1.** Let  $\mathbf{z} \in \mathbb{R}^K$  be a fixed vector of logits,  $\mathbf{q}(\mathbf{z}) = \text{softmax}(\mathbf{z}) \in \Delta$  a prediction for the target  $\mathbf{p} \in \Delta$ ,  $q_k(\mathbf{z})$  the  $k$ -th entry of vector  $\mathbf{q}(\mathbf{z})$ , and  $m \in \text{argmin}_k q_k(\mathbf{z})$ . Then, the Euclidean norm of the gradient of the LDP loss with respect to logits satisfies

$$0 \leq \|\nabla_{\mathbf{z}}[-\log(\mathbf{p}^\top \mathbf{q}(\mathbf{z}))]\|_2 \leq \|\mathbf{q}(\mathbf{z}) - e_m\|_2 = \|\nabla_{\mathbf{z}}[-e_m^\top \log(\mathbf{q}(\mathbf{z}))]\|_2. \quad (3)$$

In addition:

- (a) The minimum is achieved when the target is the uniform distribution  $\mathbf{p} = \mathbf{1}/K$ .
- (b) The maximum is achieved when the target  $\mathbf{p}$  is a one-hot vector for label  $m$ .

In Theorem 2 we show that the gradient of the LDP loss has a scaling factor that naturally decreases to 0 as the uncertainty of the target distribution increases.

**Theorem 2.** Let  $\mathbf{q}(\boldsymbol{\theta}) \in \Delta$  denote the prediction for target  $\mathbf{p} \in \Delta$  when viewed as a function of the learnable model weights  $\boldsymbol{\theta}$ . The Euclidean norm of the gradient of the LDP loss with respect to model weights satisfies

$$\|\nabla_{\boldsymbol{\theta}}[-\log(\mathbf{p}^\top \mathbf{q}(\boldsymbol{\theta}))]\|_2 \leq \frac{\max(\mathbf{p}) - \min(\mathbf{p})}{\min(\mathbf{p})} \sum_{k=1}^K \|\nabla_{\boldsymbol{\theta}} q_k(\boldsymbol{\theta})\|_2. \quad (4)$$

The scaling factor  $\frac{\max(\mathbf{p}) - \min(\mathbf{p})}{\min(\mathbf{p})}$  depends on the target distribution  $\mathbf{p}$  only, and tends to 0 when  $\mathbf{p}$  tends to the uniform distribution. From the theorem’s proof, we can tighten the bound by replacing the denominator  $\min(\mathbf{p})$  in the scaling factor by  $\mathbf{p}^\top \mathbf{q}(\boldsymbol{\theta})$ . We use the form in the theorem as it gives a simpler interpretation of the effect of target distribution uncertainty on the gradient scale.

Proposition 1 is an analogue of Theorem 1 for several unlabeled data loss functions used in many popular SSL algorithms such as (Sohn et al., 2020; Chen et al., 2023; Wang et al., 2023).

**Proposition 1.** Let  $\mathbf{z}, \mathbf{q}(\mathbf{z})$  and  $\mathbf{p}$  be defined as in Theorem 1,  $p_k$  be the  $k$ -th element of vector  $\mathbf{p}$ ,  $M \in \text{argmax}_k p_k$  be any maximizer of  $\mathbf{p}$ , and  $w : [1/K, 1] \rightarrow [0, 1]$  be non-decreasing. Then, the Euclidean norm of the gradient with respect to logits of loss functions of the form  $\ell(\mathbf{q}(\mathbf{z}), \mathbf{p}) = w(p_M)(-\log q_M(\mathbf{z}))$  satisfy the following bound:

$$w(1/K) \|\mathbf{q}(\mathbf{z}) - e_M\|_2 \leq \|\nabla_{\mathbf{z}} w(p_M)(-\log q_M(\mathbf{z}))\|_2 \leq w(1) \|\mathbf{q}(\mathbf{z}) - e_M\|_2. \quad (5)$$

In addition:

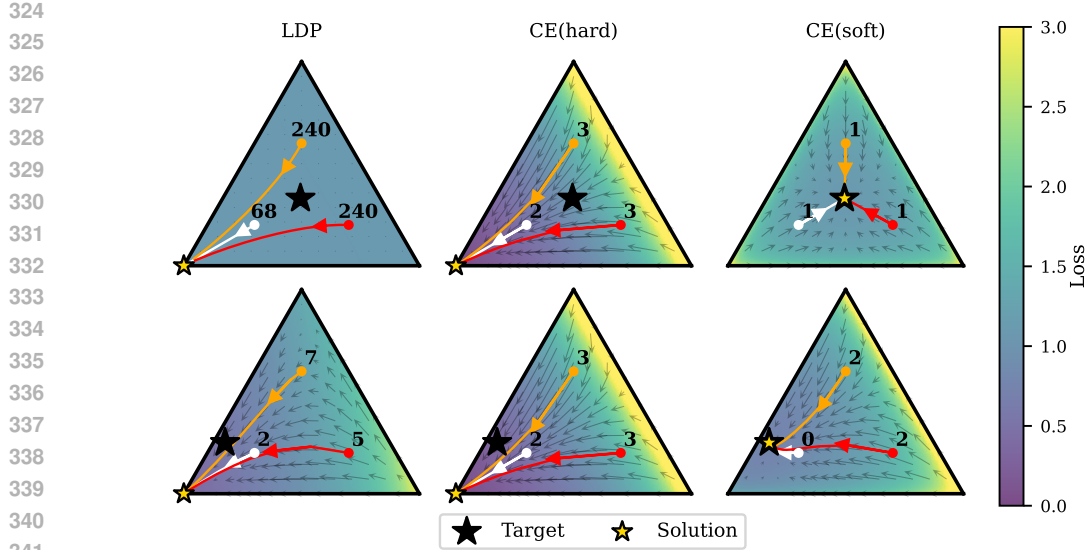


Figure 2: Gradient descent trajectories for LDP, CE(hard), and CE(soft) losses. Explicitly, each trajectory tracks the softmax of the updates  $z \leftarrow z - \nabla_z \ell(\text{softmax}(z), \mathbf{p})$ , with  $\mathbf{p}$  being the target,  $\mathbf{q}^* = \text{argmin}_q \ell(\mathbf{q}, \mathbf{p})$  the solution, and trajectory starting points  $z \in \{\log(0.4e_i + 0.2)\}_{i=1}^3$ . Number of iterations until convergence  $\|\text{softmax}(z) - \mathbf{q}^*\|_2^2 < 0.1$  also indicated. Target 1  $(0.34, 0.33, 0.33)^\top$ , target 2  $(0.7, 0.25, 0.05)^\top$  in row 1, 2 respectively.

(a) The minimum is achieved when the target is the uniform distribution  $\mathbf{p} = \mathbf{1}/K$ .

(b) The maximum is achieved when the target is any one-hot vector  $\mathbf{p} = \mathbf{e}_k$ ,  $k \in \{1, \dots, K\}$ .

A special, but common, case of Proposition 1 is when the weighting function is a step function  $w(\alpha) = \mathbf{1}(\alpha \geq \tau)$  for some  $\tau \in [1/K, 1]$ . In this case, like the LDP loss, the lower bound in Equation (5) is zero, but unlike LDP loss, the upper bound is  $\|\mathbf{q}(z) - \mathbf{e}_M\|$ , which is bounded from above by  $\|\mathbf{q}(z) - \mathbf{e}_m\|$ . Consequently, compared to commonly used loss functions of the form considered in Proposition 1, the LDP loss is able to put extra priority on examples in the extreme case where targets have maximum confidence associated with the class that has the least confidence according to the prediction.

To emphasize the relationship between confidence and gradient norm for different loss functions, we provide visualizations of gradient flow in Figure 2. LDP overall gives a flatter loss landscape compared to both variants of the CE loss, particularly when the target is close to uniform. As a result, gradient updates are generally more conservative for LDP loss, depending on the confidence of the target. A specific example from Figure 2: consider the trajectory beginning at  $(0.2, 0.6, 0.2)^\top$  (trajectory starting at top center) with target  $(0.34, 0.33, 0.33)^\top$  (top row), it takes 240 gradient descent steps to converge to the solution using LDP loss, but only 3 and 1 iterations for CE(hard) and CE(soft), respectively, despite the high uncertainty in the target. As the entropy of the target decreases, each loss function eventually becomes equivalent. However, the LDP loss tends to produce slightly different trajectories depending on the full target distribution, particularly when compared to CE(hard), which shares the same solution as LDP but ignores the complete target distribution.

We demonstrate how the gradients of the LDP loss are scaled relative to the two variants of the CE loss, CE(Soft) and CE(Hard), in Figure 3. Gradients of the LDP loss are appropriately scaled in the sense that the gradient norm  $\|\nabla_z \log(\mathbf{p}^\top \mathbf{q}(z))\|_2$  is vanishing when the target distribution is close to uniform, close to the gradient norm of CE(Hard) when the target distribution is far from uniform, and is interpolated approximately monotonically in between. The vanishing behavior is expected from Theorem 2 (when applied to the logits), and also note that when  $\mathbf{p} = \mathbf{e}_{\text{argmin}(\mathbf{q})}$ , the CE(Hard) loss gradient  $\|\nabla_z \mathbf{e}_{\text{argmax}(\mathbf{p})}^\top \log(\mathbf{q}(z))\|_2$  upper-bounds the LDP gradient norm  $\|\nabla_z \log(\mathbf{p}^\top \mathbf{q}(z))\|_2$ , by Theorem 1 and illustrated in Figure 3. As a result, when using the LDP loss for consistency regularization in SSL, we do not require any additional regularization constant to ensure the gradient

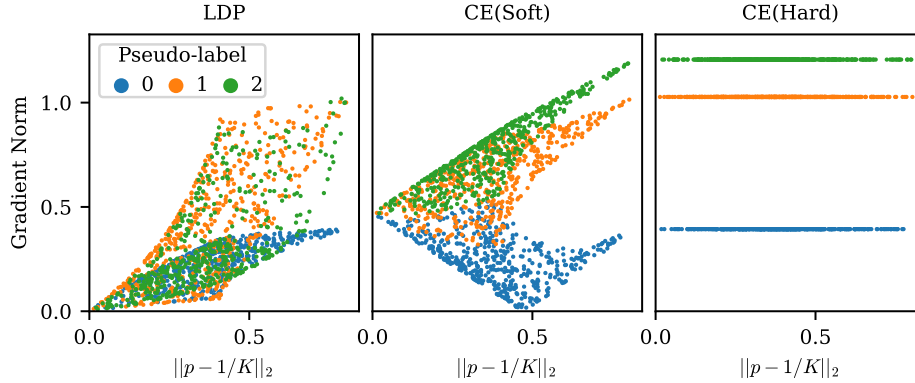


Figure 3: Comparing the distance  $\|\mathbf{p} - \mathbf{1}/K\|_2$  between target vectors  $\mathbf{p}$  and the uniform distribution with the norm  $\|\nabla_z \ell(\mathbf{q}(z), \mathbf{p})\|_2$  of the gradient with respect to logits for LDP, CE(Soft) and CE(Hard). We fix the prediction vector as  $\mathbf{q} = (0.7, 0.25, 0.05)^\top$ . Target vectors are randomly sampled from the probability simplex  $\Delta$ . Pseudo-labels of each target (index with largest element in  $\mathbf{p}$ ) are indicated with different colors. Gradient norms of LDP loss are naturally scaled to be approximately monotonic in the distance between the target distribution and the uniform distribution, with the gradient being exactly equal  $\mathbf{0}$  when the target is the uniform distribution. Neither of the CE loss variants have this same monotonic scaling property.

of the consistency regularizer is scaled appropriately relative to the supervised loss. While the gradient norm is generally smallest when both the target and prediction agree on the most likely label for all three loss functions, the CE(Soft) and CE(Hard) loss functions do not share the desirable scaling property where gradient norms are smaller for highly uncertain targets; in particular, even when the target is completely random and thus should be ignored, CE(Soft) and CE(Hard) still have non-zero gradients. This is why when using CE(Soft) or CE(Hard), additional weights are required to manually shrink the gradient norms when the target is close to uniform.

#### 4.3 DOTMATCH OBJECTIVE FUNCTION

The DotMatch objective function for a single batch of semi-supervised data is

$$L(\theta) = -\frac{1}{|\mathcal{L}|} \sum_{(\mathbf{w}, \mathbf{y}) \in \mathcal{L}} \log(\mathbf{y}^\top \mathbf{p}(\mathbf{w}; \theta)) - \frac{1}{|\mathcal{U}|} \sum_{(\mathbf{w}, \mathbf{s}) \in \mathcal{U}} \log(\text{DA}(\mathbf{p}_w)^\top \mathbf{p}(\mathbf{s}; \theta)), \quad (6)$$

where  $\text{DA}(\mathbf{p}) = \text{sumnorm}\left(\left(\mathbf{p} + \varepsilon \mathbf{1}\right) \odot \frac{\boldsymbol{\pi} + \varepsilon \mathbf{1}}{\boldsymbol{\pi}_t + \varepsilon \mathbf{1}}\right)$ ,  $\text{sumnorm}(\mathbf{x}) = \mathbf{x} / \sum_j x_j$ , division is element-wise,  $\boldsymbol{\pi} \in \Delta$  is the class prior estimated using all available labels,  $\hat{\boldsymbol{\pi}}_t = m \hat{\boldsymbol{\pi}}_{t-1} + (1 - m) \frac{1}{|\mathcal{U}|} \sum_{(\mathbf{w}, \mathbf{s}) \in \mathcal{U}} \mathbf{p}_w$  is the EMA of the prior estimated using batches of unlabeled data during training,  $m \in [0, 1]$  is the EMA momentum, and  $\varepsilon$  is a small constant to avoid numerical instabilities which we set to  $1e-6$  in our experiments. DotMatch inherits the nice properties of the LDP loss, namely consistency, entropy minimization and small gradient norms for low-confidence unlabeled examples. DotMatch also addresses the quantity-quality tradeoff (Chen et al., 2023) by always utilizing all unlabeled examples (high quantity), while implicitly down-weighting the contribution from low-confidence examples (high quality). We compare the performance of DotMatch against state-of-the-art baselines in the next section, and perform an ablation study.

## 5 EXPERIMENTS

**Baselines** We compare DotMatch against two supervised learning baselines: Supervised, which only uses available labeled examples for training to provide a lower bound on performance, and FullySupervised, which uses all labels including true labels for unlabeled examples to give an upper bound on performance. Strong semi-supervised baselines include single model SSL algorithms that have implementations provided with the original papers: FixMatch (Sohn et al., 2020), SoftMatch (Chen et al., 2023), FreeMatch (Wang et al., 2023), FlatMatch (Huang et al., 2023), and InterLUDE (Huang et al., 2024).

Table 2: Test error rates (mean±std% over 3 seeds) of various SSL methods. Best in **bold**, overlapping 95% CI underlined. FullySupervised results for STL10 omitted since not all labels available.

Dataset # Label	MNIST			EMNIST			CIFAR10			SVHN			STL10		#1	
	10	40	250	188	1175	10	40	250	400	2500	40	250	40	1000		
Supervised	67.82±3.26	<u>28.96±2.06</u>	7.37±0.84	85.69±1.00	51.94±0.96	23.01±0.77	82.49±1.24	77.01±1.64	56.62±0.41	89.45±0.86	60.05±0.17	83.46±2.29	25.37±2.79	75.49±0.75	32.57±1.79	0
FixMatch	44.27±12.83	<u>4.78±3.75</u>	1.29±0.04	45.58±1.11	23.88±1.48	<u>16.07±0.74</u>	78.75±3.96	21.03±7.98	<u>5.41±0.24</u>	51.49±1.40	29.33±0.09	<b>5.85±0.07</b>	<u>2.31±0.04</u>	43.65±5.96	<u>6.81±0.26</u>	1
SoftMatch	7.31±2.62	<u>3.65±2.36</u>	1.33±0.06	46.42±4.12	<u>20.23±0.08</u>	15.49±0.41	37.75±0.06	6.75±0.20	<u>5.34±0.09</u>	43.12±1.24	<u>27.60±0.34</u>	<u>21.58±12.24</u>	2.54±0.04	26.01±6.40	<b>6.53±0.06</b>	1
FreeMatch	<b>4.70±2.39</b>	<u>2.52±0.85</u>	1.34±0.08	51.88±0.06	<u>20.36±1.89</u>	15.83±0.64	<u>30.44±7.19</u>	<b>5.41±0.04</b>	<b>5.28±0.24</b>	<u>43.46±2.31</u>	<b>27.12±0.29</b>	<u>17.94±11.16</u>	5.01±1.19	<u>28.65±3.07</u>	<u>6.54±0.04</u>	4
FlatMatch	<u>24.06±18.17</u>	<b>2.30±0.99</b>	1.66±0.07	52.70±7.27	<u>21.03±1.71</u>	15.83±0.29	<u>29.26±9.44</u>	6.35±0.46	6.06±0.37	46.09±1.94	29.13±0.26	17.32±4.65	9.63±0.76	<b>21.00±1.24</b>	7.76±0.09	2
InterLUDE	67.49±4.62	<u>13.16±0.22</u>	<b>1.10±0.19</b>	71.45±7.21	32.44±3.10	16.98±0.31	55.78±8.95	27.70±3.63	<u>5.75±0.41</u>	58.99±2.91	46.28±1.76	30.59±12.65	<b>2.25±0.02</b>	42.78±3.35	7.87±0.20	2
DotMatch (Ours)	<u>5.74±1.56</u>	<u>2.96±2.47</u>	<u>1.23±0.17</u>	<b>36.20±3.39</b>	<b>19.66±1.36</b>	<b>15.46±0.89</b>	<b>25.32±16.42</b>	10.43±2.85	<u>5.67±0.23</u>	<b>42.96±0.86</b>	28.19±0.52	<u>15.15±15.42</u>	2.49±0.07	28.65±4.36	7.07±0.08	5
FullySupervised		0.90±0.04			11.42±0.11			4.92±0.09		26.93±0.00		2.46±0.19				-

**Datasets** We choose datasets with a wide variety of input size, number of classes and number of training examples: MNIST (Lecun et al., 1998), EMNIST(balanced) (Cohen et al., 2017), CIFAR10, CIFAR100 (Krizhevsky, 2009), SVHN (Netzer et al., 2011), and STL10 (Coates et al., 2011). A summary of benchmark datasets is provided in [Appendix E.1](#).

**Hyperparameters** All experiments were run using the USB (Wang et al., 2022) classic CV benchmark (equivalent to TorchSSL (Zhang et al., 2021)), apart from FlatMatch and InterLUDE for which we adapted their provided code bases. The architectures used for datasets with color images: WRN-28-2 (Zagoruyko & Komodakis, 2016) for CIFAR10 and SVHN, WRN-28-8 for CIFAR100, WRN-37-2 for STL10 (Zhou et al., 2020). On black and white image datasets we used WRN-10-1 for MNIST and WRN-22-1 for EMNIST, modified to have only one input channel. At test time we make predictions using EMA of the model  $\text{argmax}(\mathbf{p}(\mathbf{x}; \theta_{\text{EMA}}))$ . EMA momentum  $m$  is set to 0.999. Number of training iterations  $T$  for CIFAR10, CIFAR100, SVHN, STL10 is  $2 \times 10^5$ , similar to USB, except here we train all models from scratch instead of starting with pre-trained models. Number of training iterations for MNIST and EMNIST is  $5 \times 10^4$ . Learning rate follows the cosine schedule  $0.03 \cos(\frac{7\pi t}{16T})$ , where  $t$  is the training iteration number. Optimizer is SGD with nesterov momentum 0.9. Labeled data batch size 64, unlabeled data batch size 448. See [Appendix E.2](#) for a detailed breakdown of all hyperparameter settings, including algorithm specific hyperparameters used for the baseline methods.

**Benchmark** We compare all algorithms on the datasets with different numbers of labeled examples per class. For example, CIFAR10 with 10 labels means there is 1 labeled example per class. Each algorithm is trained using 3 different random seeds. Following standard practice (Zhang et al., 2021; Chen et al., 2023; Wang et al., 2023; Huang et al., 2023), we report average of best test error rate across all checkpoints in [Table 2](#). Checkpoints occur at every 1000 training iterations. Test error measured after training is completed can be found in [Table 7](#) in [Appendix F](#).

**Results** As shown in [Table 2](#), DotMatch achieves the best average test error rate more often than any of the baseline algorithms (5 for DotMatch versus 4 for second best FreeMatch). It performs particularly well with extremely limited labeled data, either 1 or 4 labels per class. In these settings it achieved the best or equivalent to the best test accuracy 5 times which ties with SoftMatch and FreeMatch. Surprisingly, FixMatch demonstrated the strongest performance on SVHN with 40 labels, with both small average test error and small variance across the random seeds. The small variance is likely due to the simplicity of using a constant confidence threshold rather than a dynamic one. DotMatch, SoftMatch and FreeMatch all had equivalently strong performance on this problem, however with wider confidence intervals.

Notably, DotMatch performs the best across all label settings for the EMNIST dataset, which, to the best of our knowledge, has not been used in such an SSL benchmark previously. In particular, DotMatch is significantly better than all baseline methods on EMNIST with 1 labeled example per class. This is likely because DotMatch has a significant advantage of not requiring tuning for any algorithm specific hyperparameters, whereas we copied the recommended algorithm specific hyperparameters for baseline methods from CIFAR10 since there was not enough labeled data available to do proper tuning.

**Runtime** We compared the runtime per iteration of DotMatch against baseline methods and report the results in [Table 8](#) in the Appendix. All algorithms were implemented in a unified code base and run on identical hardware. DotMatch is consistently the fastest across all datasets, with FixMatch generally being the second fastest, followed by FreeMatch then SoftMatch.

Table 3: Test error rates (mean $\pm$ std over 3 seeds) for different variations of DotMatch on EMNIST with varying amounts of labeled data. Best in **bold**, overlapping 95% CI underlined.

Loss	DA	EMNIST		
		47	188	1175
CE(soft)	✗	65.87 $\pm$ 0.48	28.84 $\pm$ 0.76	17.73 $\pm$ 0.95
CE(hard)	✗	50.77 $\pm$ 2.80	27.20 $\pm$ 2.32	16.95 $\pm$ 0.58
LDP	✗	56.73 $\pm$ 1.77	24.99 $\pm$ 1.92	<u>15.61<math>\pm</math>0.60</u>
CE(soft)	✓	79.21 $\pm$ 1.23	43.28 $\pm$ 2.45	21.09 $\pm$ 0.72
CE(hard)	✓	46.21 $\pm$ 2.44	26.03 $\pm$ 2.06	16.90 $\pm$ 0.34
LDP	✓	<b>36.20<math>\pm</math>3.30</b>	<b>19.66<math>\pm</math>1.36</b>	<b>15.46<math>\pm</math>0.59</b>

Overall, DotMatch generally performs at least as well as the state-of-art baselines on most problems, while being simpler, faster, and not requiring any tuning of algorithm specific hyperparameters like regularization strength or confidence threshold schedule.

### 5.1 ABLATION

We conduct experiments on the utility of the components of DotMatch, namely the choice of loss function and use of DA. We compare three different choices for the unlabeled loss function  $\ell_u$ : CE(soft), CE(hard), and LDP, each with and without DA applied to the targets. We train on EMNIST using the same experiment settings from the main benchmark.

As shown in Table 3, DotMatch (LDP+DA) achieves the best result in every case, with only LDP without DA matching the performance when there are 1175 labeled examples, or 25 per class. Other than LDP loss, CE(hard) with DA generally gives the second best results across all settings, which validates that entropy minimization (in this case via one-hot pseudo-labeling) is an important design consideration for SSL algorithms, which is lacking when using CE(soft). Adding DA consistently gives performance at least as good, if not better, when combined with CE(hard) or LDP loss, but hurts performance when using CE(soft). We suspect this is because it is most sensitive to changes in the target since the minimizer also changes, whereas for CE(hard) and LDP the minimizer only changes if the pseudo-label changes, which is not always the case when using DA.

## 6 CONCLUSION

We proposed the log dot product loss function for consistency regularization in SSL, and used it in the novel algorithm DotMatch. Theoretical analysis shows that DotMatch exhibits desirable qualities of an SSL algorithm, namely consistency, entropy minimization and small gradient norms for low-confidence unlabeled examples. Notably, no algorithm specific hyperparameters need to be tuned for DotMatch, making it much simpler than existing algorithms. Minimizing the number of hyperparameters in SSL algorithms is particularly important since there is often insufficient labeled examples available for validation purposes. Our DotMatch demonstrated strong performance in an extensive benchmark against state-of-the-art baselines, and we showed that the proposed configuration is the best among natural alternatives in the ablation study.

Since the unlabeled data loss used in DotMatch based on LDP loss does not require any labels, and the properties of consistency, entropy minimization and DA are desirable for any classifier, it is also suitable as a plug-in regularizer in other weakly supervised learning (WSL) (Zhou, 2017) settings to help improve performance without introducing additional hyperparameters. For example, in learning from complementary labels (Ishida et al., 2017; Feng et al., 2020), noisy labels (Frenay & Verleysen, 2014; Song et al., 2023), partial labels (Grandvalet & Bengio, 2004b; Lv et al., 2024), positive-unlabeled data (Elkan & Noto, 2008), similarity labels (Bao et al., 2018; Hsu et al., 2019) or unlabeled data with label proportions (Tang et al., 2023). A limitation of this study is that we did not consider SSL problems with large class imbalances, mainly because it is difficult to estimate the true imbalanced label distribution in interesting and challenging settings with extremely scarce labeled data. In principle, DotMatch automatically adapts to imbalanced scenarios, provided the labeled data label distribution is an accurate reflection of the true label distribution for unlabeled data. We leave the thorough investigation into imbalanced SSL problems and use of LDP in other WSL settings to future work.

## REFERENCES

- 540  
541  
542 Philip Bachman, Quais Alsharif, and Doina Precup. Learning with pseudo-ensembles. In Z. Ghahra-  
543 mani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger (eds.), *Advances in Neural*  
544 *Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- 545 Han Bao, Gang Niu, and Masashi Sugiyama. Classification from pairwise similarity and unlabeled  
546 data. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference*  
547 *on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 452–461.  
548 PMLR, 2018.
- 549 David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A  
550 Raffel. Mixmatch: A holistic approach to semi-supervised learning. In H. Wallach, H. Larochelle,  
551 A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information*  
552 *Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 553 David Berthelot, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and  
554 Colin Raffel. Remixmatch: Semi-supervised learning with distribution matching and augmenta-  
555 tion anchoring. In *International Conference on Learning Representations*, 2020.
- 556 David Berthelot, Rebecca Roelofs, Kihyuk Sohn, Nicholas Carlini, and Alexey Kurakin. Adamatch:  
557 A unified approach to semi-supervised learning and domain adaptation. In *International Confer-*  
558 *ence on Learning Representations*, 2022.
- 559 John Bridle, Anthony Heading, and David MacKay. Unsupervised classifiers, mutual information  
560 and 'phantom targets'. In J. Moody, S. Hanson, and R.P. Lippmann (eds.), *Advances in Neural*  
561 *Information Processing Systems*, volume 4. Morgan-Kaufmann, 1991.
- 562 Kaidi Cao, Maria Brbic, and Jure Leskovec. Open-world semi-supervised learning. In *International*  
563 *Conference on Learning Representations*, 2022.
- 564 Baixu Chen, Junguang Jiang, Ximei Wang, Pengfei Wan, Jianmin Wang, and Mingsheng Long.  
565 Debaised self-training for semi-supervised learning. In S. Koyejo, S. Mohamed, A. Agarwal,  
566 D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*,  
567 volume 35, pp. 32424–32437. Curran Associates, Inc., 2022.
- 568 Hao Chen, Ran Tao, Yue Fan, Yidong Wang, Jindong Wang, Bernt Schiele, Xing Xie, Bhiksha Raj,  
569 and Marios Savvides. Softmatch: Addressing the quantity-quality tradeoff in semi-supervised  
570 learning. In *The Eleventh International Conference on Learning Representations*, 2023.
- 571 Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised  
572 feature learning. In Geoffrey Gordon, David Dunson, and Miroslav Dudík (eds.), *Proceedings*  
573 *of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of  
574 *Proceedings of Machine Learning Research*, pp. 215–223, Fort Lauderdale, FL, USA, 11–13 Apr  
575 2011. PMLR.
- 576 Gregory Cohen, Saeed Afshar, Jonathan Tapson, and André van Schaik. EMNIST: Extending mnist  
577 to handwritten letters. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pp.  
578 2921–2926, 2017. doi: 10.1109/IJCNN.2017.7966217.
- 579 Ekin Dogus Cubuk, Barret Zoph, Jon Shlens, and Quoc Le. Randaugment: Practical automated  
580 data augmentation with a reduced search space. In H. Larochelle, M. Ranzato, R. Hadsell, M.F.  
581 Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp.  
582 18613–18624. Curran Associates, Inc., 2020.
- 583 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-  
584 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,  
585 pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 586 Qinyi Deng, Yong Guo, Zhibang Yang, Haolin Pan, and Jian Chen. Boosting semi-supervised  
587 learning with contrastive complementary labeling. *Neural Networks*, 170:417–426, 2024. ISSN  
588 0893-6080. doi: <https://doi.org/10.1016/j.neunet.2023.11.052>.

- 594 Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks  
595 with cutout. *arXiv preprint arXiv:1708.04552*, 2017.  
596
- 597 Linhao Dong, Shuang Xu, and Bo Xu. Speech-transformer: A no-recurrence sequence-to-sequence  
598 model for speech recognition. In *2018 IEEE International Conference on Acoustics, Speech and*  
599 *Signal Processing (ICASSP)*, pp. 5884–5888, 2018. doi: 10.1109/ICASSP.2018.8462506.
- 600 Charles Elkan and Keith Noto. Learning classifiers from only positive and unlabeled data. In  
601 *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and*  
602 *Data Mining (KDD2008)*, pp. 213–220, 2008.  
603
- 604 Shijie Fang, Qianhan Feng, and Tong Lin. VCC-INFUSE: Towards accurate and efficient selection  
605 of unlabeled examples in semi-supervised learning. In Kate Larson (ed.), *Proceedings of the*  
606 *Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, pp. 3953–3961.  
607 International Joint Conferences on Artificial Intelligence Organization, 8 2024. doi: 10.24963/  
608 ijcai.2024/437.
- 609 Lei Feng, Takuo Kaneko, Bo Han, Gang Niu, Bo An, and Masashi Sugiyama. Learning with mul-  
610 tiple complementary labels. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th*  
611 *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learn-*  
612 *ing Research*, pp. 3072–3081. PMLR, 2020.
- 613 Benoit Frenay and Michel Verleysen. Classification in the presence of label noise: A survey. *IEEE*  
614 *Transactions on Neural Networks and Learning Systems*, 25(5):845–869, 2014. doi: 10.1109/  
615 TNNLS.2013.2292894.
- 616 Yves Grandvalet and Yoshua Bengio. Semi-supervised learning by entropy minimization. In L. Saul,  
617 Y. Weiss, and L. Bottou (eds.), *Advances in Neural Information Processing Systems*, volume 17.  
618 MIT Press, 2004a.
- 619 Yves Grandvalet and Yoshua Bengio. Learning from partial labels with minimum entropy. CIRANO  
620 Working Papers 2004s-28, CIRANO, 2004b.  
621
- 622 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
623 nition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*  
624 *(CVPR)*, June 2016.
- 625 Marzi Heidari, Hanping Zhang, and Yuhong Guo. Reinforcement learning guided semi-supervised  
626 learning. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang  
627 (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 136990–137009.  
628 Curran Associates, Inc., 2024.  
629
- 630 Yen-Chang Hsu, Zhaoyang Lv, Joel Schlosser, Phillip Odom, and Zsolt Kira. Multi-class classi-  
631 fication without multi-class labels. In *International Conference on Learning Representations*,  
632 2019.
- 633 Zhe Huang, Xiaowei Yu, Dajiang Zhu, and Michael C Hughes. InterLUDE: Interactions between  
634 labeled and unlabeled data to enhance semi-supervised learning. In Ruslan Salakhutdinov, Zico  
635 Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp  
636 (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of  
637 *Proceedings of Machine Learning Research*, pp. 20452–20473. PMLR, 21–27 Jul 2024.
- 638 Zhuo Huang, Li Shen, Jun Yu, Bo Han, and Tongliang Liu. Flatmatch: Bridging labeled data  
639 and unlabeled data with cross-sharpness for semi-supervised learning. In A. Oh, T. Naumann,  
640 A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Pro-*  
641 *cessing Systems*, volume 36, pp. 18474–18494. Curran Associates, Inc., 2023.
- 642 Takashi Ishida, Gang Niu, Weihua Hu, and Masashi Sugiyama. Learning from complementary la-  
643 bels. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and  
644 R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Asso-  
645 ciates, Inc., 2017.
- 646 Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University  
647 of Toronto, 2009.

- 648 Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. In *International*  
649 *Conference on Learning Representations*, 2017.
- 650 Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recog-  
651 nition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. doi: 10.1109/5.726791.
- 652  
653 Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for  
654 deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3,  
655 pp. 896. Atlanta, 2013.
- 656 Muyang Li, Runze Wu, Haoyu Liu, Jun Yu, Xun Yang, Bo Han, and Tongliang Liu. InstanT: Semi-  
657 supervised learning with instance-dependent thresholds. In A. Oh, T. Naumann, A. Globerson,  
658 K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*,  
659 volume 36, pp. 2922–2938. Curran Associates, Inc., 2023.
- 660 Siyuan Li, Weiyang Jin, Zedong Wang, Fang Wu, Zicheng Liu, Cheng Tan, and Stan Z. Li. Semire-  
661 ward: A general reward model for semi-supervised learning. In *The Twelfth International Con-*  
662 *ference on Learning Representations*, 2024.
- 663 Shuyang Liu, Ruiqiu Zheng, Yunhang Shen, Zhou Yu, Ke Li, Xing Sun, and Shaohui Lin. Prob-  
664 ability-density-aware semi-supervised learning. *Proceedings of the AAAI Conference on Ar-*  
665 *tificial Intelligence*, 39(18):18943–18950, Apr. 2025. doi: 10.1609/aaai.v39i18.34085.
- 666  
667 Jiaqi Lv, Yangfan Liu, Shiyu Xia, Ning Xu, Miao Xu, Gang Niu, Min-Ling Zhang, Masashi  
668 Sugiyama, and Xin Geng. What makes partial-label learning algorithms effective? In A. Globerson,  
669 L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in*  
670 *Neural Information Processing Systems*, volume 37, pp. 89513–89534. Curran Associates, Inc.,  
671 2024.
- 672  
673 G. J. McLachlan. Iterative reclassification procedure for constructing an asymptotically optimal rule  
674 of allocation in discriminant analysis. *Journal of the American Statistical Association*, 70(350):  
675 365–369, 1975. doi: 10.1080/01621459.1975.10479874.
- 676 Zeping Min, Jinfeng Bai, and Chengfei Li. Leveraging local variance for pseudo-label selection in  
677 semi-supervised learning. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(13):  
678 14370–14378, Mar. 2024. doi: 10.1609/aaai.v38i13.29350.
- 679 Takeru Miyato, Shin-Ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: A  
680 regularization method for supervised and semi-supervised learning. *IEEE Transactions on Pattern*  
681 *Analysis and Machine Intelligence*, 41(8):1979–1993, 2019. doi: 10.1109/TPAMI.2018.2858821.
- 682  
683 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading  
684 digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning*  
685 *and Unsupervised Feature Learning 2011*, 2011.
- 686 Chuck Rosenberg, Henry Schneiderman, and Martial Hebert. Semi-Supervised Self-Training of  
687 Object Detection Models. In *Proceedings of the Seventh IEEE Workshops on Application of*  
688 *Computer Vision*, volume 1, pp. 29–36. IEEE Computer Society, 2005. doi: 10.1109/ACVMOT.  
689 2005.107.
- 690 Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Mutual exclusivity loss for semi-supervised  
691 deep learning. In *2016 IEEE International Conference on Image Processing (ICIP)*, pp. 1908–  
692 1912, 2016a. doi: 10.1109/ICIP.2016.7532690.
- 693  
694 Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transforma-  
695 tions and perturbations for deep semi-supervised learning. In D. Lee, M. Sugiyama, U. Luxburg,  
696 I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 29.  
697 Curran Associates, Inc., 2016b.
- 698 Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel,  
699 Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised  
700 learning with consistency and confidence. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan,  
701 and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 596–608.  
Curran Associates, Inc., 2020.

- 702 Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy  
703 labels with deep neural networks: A survey. *IEEE Transactions on Neural Networks and Learning*  
704 *Systems*, 34(11):8135–8153, 2023. doi: 10.1109/TNNLS.2022.3152527.
- 705  
706 Zhiquan Tan, Kaipeng Zheng, and Weiran Huang. OTMatch: Improving semi-supervised learning  
707 with optimal transport. In *Forty-first International Conference on Machine Learning*, 2024.
- 708  
709 Yuting Tang, Nan Lu, Tianyi Zhang, and Masashi Sugiyama. Multi-class classification from multiple  
710 unlabeled datasets with partial risk regularization. In Emtiyaz Khan and Mehmet Gonen (eds.),  
711 *Proceedings of The 14th Asian Conference on Machine Learning*, volume 189 of *Proceedings of*  
712 *Machine Learning Research*, pp. 990–1005. PMLR, 2023.
- 713  
714 Jesper E van Engelen and Holger H Hoos. A survey on semi-supervised learning. *Machine Learning*,  
109(2):373–440, February 2020.
- 715  
716 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
717 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg,  
718 S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural*  
719 *Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 720  
721 Yidong Wang, Hao Chen, and Yue Fan et al. USB: A unified semi-supervised learning benchmark  
722 for classification. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh  
723 (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 3938–3961. Curran  
724 Associates, Inc., 2022.
- 725  
726 Yidong Wang, Hao Chen, Qiang Heng, Wenxin Hou, Yue Fan, Zhen Wu, Jindong Wang, Marios Sav-  
727 vides, Takahiro Shinozaki, Bhiksha Raj, Bernt Schiele, and Xing Xie. Freematch: Self-adaptive  
728 thresholding for semi-supervised learning. In *The Eleventh International Conference on Learning*  
*Representations*, 2023.
- 729  
730 Zhiyu Wu and Jinshi Cui. Allmatch: Exploiting all unlabeled data for semi-supervised learning. In  
731 Kate Larson (ed.), *Proceedings of the Thirty-Third International Joint Conference on Artificial*  
732 *Intelligence, IJCAI-24*, pp. 5245–5253. International Joint Conferences on Artificial Intelligence  
733 Organization, 8 2024. doi: 10.24963/ijcai.2024/580.
- 734  
735 Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and Quoc Le. Unsupervised data augmentation  
736 for consistency training. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin  
737 (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 6256–6268. Curran  
738 Associates, Inc., 2020a.
- 739  
740 Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V. Le. Self-training with noisy student  
741 improves imagenet classification. In *Proceedings of the IEEE/CVF Conference on Computer*  
742 *Vision and Pattern Recognition (CVPR)*, June 2020b.
- 743  
744 Lihe Yang, Zhen Zhao, Lei Qi, Yu Qiao, Yinghuan Shi, and Hengshuang Zhao. Shrinking class  
745 space for enhanced certainty in semi-supervised learning. In *Proceedings of the IEEE/CVF Inter-*  
746 *national Conference on Computer Vision (ICCV)*, pp. 16187–16196, October 2023.
- 747  
748 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In Edwin R. Hancock Richard  
749 C. Wilson and William A. P. Smith (eds.), *Proceedings of the British Machine Vision Conference*  
750 *(BMVC)*, pp. 87.1–87.12. BMVA Press, September 2016. ISBN 1-901725-59-6. doi: 10.5244/C.  
751 30.87.
- 752  
753 Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jindong Wang, Manabu Okumura, and  
754 Takahiro Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo la-  
755 beling. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.),  
*Advances in Neural Information Processing Systems*, volume 34, pp. 18408–18419. Curran As-  
sociates, Inc., 2021.
- Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. MixUp: Beyond empir-  
ical risk minimization. In *International Conference on Learning Representations*, 2018.

756 Tianyi Zhou, Shengjie Wang, and Jeff Bilmes. Time-consistent self-supervision for semi-supervised  
757 learning. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Con-*  
758 *ference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp.  
759 11523–11533. PMLR, 13–18 Jul 2020.

760  
761  
762 Zhi-Hua Zhou. A brief introduction to weakly supervised learning. *National Science Review*, 5(1):  
763 44–53, 08 2017. ISSN 2095-5138. doi: 10.1093/nsr/nwx106.

764  
765 X. Zhu and A. Goldberg. *Introduction to Semi-Supervised Learning*. Synthesis Lectures on  
766 Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2009. ISBN  
767 9781598295481.

768  
769  
770 Xiaojin Zhu. Semi-supervised learning literature survey. Technical Report 1530, Computer Sci-  
771 ences, University of Wisconsin-Madison, 2005.

## 772 773 A DERIVATION OF RESULTS FROM TABLE 1

774  
775 Let  $\mathbf{q}(\mathbf{z}) = \text{softmax}(\mathbf{z}) \in \Delta$  be the predicted label distribution for the target  $\mathbf{p} \in \Delta$ ,  $\mathbf{z} \in \mathbb{R}^K$  are  
776 logits, and  $\hat{y} = \text{argmax}(\mathbf{p})$ . We first derive the gradient of a single element  $\nabla_{\mathbf{z}} q_i(\mathbf{z})$ , which we will  
777 use to simplify other calculations. To do this, consider the partial derivative of a single element  $q_i$   
778 with respect to a single element of  $\mathbf{z}$ ,  $z_j$ .

$$\begin{aligned} 780 \frac{\partial q_i}{\partial z_j} &= q_i \frac{\partial \log q_i}{\partial z_j} \\ 781 &= q_i \frac{\partial}{\partial z_j} \log \left( \frac{\exp(z_i)}{\sum_{k=1}^K \exp(z_k)} \right) \\ 782 &= q_i \frac{\partial}{\partial z_j} \left( z_i - \log \left( \sum_{k=1}^K \exp(z_k) \right) \right) \\ 783 &= q_i \left( \mathbb{1}(i=j) - \frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)} \right) \\ 784 &= q_i (\mathbb{1}(i=j) - q_j). \end{aligned}$$

785  
786  
787  
788  
789  
790  
791  
792  
793 The gradient of  $q_i$  with respect to  $\mathbf{z}$  is thus

$$794 \nabla_{\mathbf{z}} q_i(\mathbf{z}) = q_i(\mathbf{z})(\mathbf{e}_i - \mathbf{q}(\mathbf{z})),$$

795  
796  
797 and the full Jacobian matrix

$$798 \nabla_{\mathbf{z}} \mathbf{q}(\mathbf{z}) = \text{diag}(\mathbf{q}(\mathbf{z})) - \mathbf{q}(\mathbf{z})\mathbf{q}(\mathbf{z})^\top,$$

799  
800  
801 where  $\text{diag}(\mathbf{q})$  is the diagonal matrix with elements of  $\mathbf{q}$  on the diagonal.

802  
803 **CE(soft)** The loss is

$$\begin{aligned} 804 \ell(\mathbf{q}, \mathbf{p}) &= -\mathbf{p}^\top \log \mathbf{q} \\ 805 &= \mathbf{p}^\top \log \left( \frac{\mathbf{p}}{\mathbf{q}} \right) - \mathbf{p}^\top \log \mathbf{p} \\ 806 &= \text{KL}(\mathbf{p} \parallel \mathbf{q}) + \text{Ent}(\mathbf{p}). \end{aligned}$$

810 Minimizing the CE loss over  $\mathbf{q}$  is equivalent to minimizing the KL-divergence, which happens when  
 811  $\mathbf{q} = \mathbf{p}$ . Now we derive the gradient of the CE(soft) loss with respect to logits:

$$\begin{aligned}
 812 \quad \nabla_{\mathbf{z}} - \mathbf{p}^\top \log \mathbf{q}(\mathbf{z}) &= -\nabla_{\mathbf{z}} \sum_{k=1}^K p_k \log q_k(\mathbf{z}) \\
 813 &= -\sum_{k=1}^K p_k \nabla_{\mathbf{z}} \log q_k(\mathbf{z}) \\
 814 &= -\sum_{k=1}^K p_k \frac{\nabla_{\mathbf{z}} q_k(\mathbf{z})}{q_k(\mathbf{z})} \\
 815 &= -\sum_{k=1}^K p_k \frac{q_k(\mathbf{z})(\mathbf{e}_k - \mathbf{q}(\mathbf{z}))}{q_k(\mathbf{z})} \\
 816 &= -\sum_{k=1}^K p_k (\mathbf{e}_k - \mathbf{q}(\mathbf{z})) \\
 817 &= \mathbf{q}(\mathbf{z}) - \mathbf{p}.
 \end{aligned}$$

818 Now consider  $\mathbf{q}$  as a general function of parameters  $\boldsymbol{\theta}$ . The gradient with respect to  $\boldsymbol{\theta}$  is:

$$\begin{aligned}
 819 \quad \nabla_{\boldsymbol{\theta}} - \mathbf{p}^\top \log \mathbf{q}(\boldsymbol{\theta}) &= -\nabla_{\boldsymbol{\theta}} \sum_{k=1}^K p_k \log q_k(\boldsymbol{\theta}) \\
 820 &= -\sum_{k=1}^K \frac{p_k}{q_k(\boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} q_k(\boldsymbol{\theta}).
 \end{aligned}$$

821 CE(HARD)

822 The loss is  $\ell(\mathbf{q}, \mathbf{p}) = -\mathbf{e}_{\hat{y}}^\top \log \mathbf{q} = -\log q_{\hat{y}}$ . This is minimized when  $q_{\hat{y}}$  is maximized, that is,  $\mathbf{q}$   
 823 assigns all probability mass to class  $\hat{y}$ ,  $\mathbf{q} = \mathbf{e}_{\hat{y}}$ .

824 Since CE(hard) is a special case of CE(soft) with target  $\mathbf{p} = \mathbf{e}_{\hat{y}}$ , we can plug this value into the  
 825 gradients for CE(soft) to get the gradients of CE(hard):

$$\begin{aligned}
 826 \quad \nabla_{\mathbf{z}} - \mathbf{e}_{\hat{y}}^\top \log \mathbf{q}(\mathbf{z}) &= \mathbf{q}(\mathbf{z}) - \mathbf{e}_{\hat{y}}, \\
 827 \quad \nabla_{\boldsymbol{\theta}} - \mathbf{e}_{\hat{y}}^\top \log \mathbf{q}(\boldsymbol{\theta}) &= -\frac{1}{q_{\hat{y}}(\boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} q_{\hat{y}}(\boldsymbol{\theta}).
 \end{aligned}$$

828 LDP

829 The LDP loss is  $\ell(\mathbf{q}, \mathbf{p}) = -\log(\mathbf{p}^\top \mathbf{q}) = -\log(\sum_{k=1}^K p_k q_k)$ . To minimize this loss over  $\mathbf{q}$ , we  
 830 need to maximize the dot product  $\sum_{k=1}^K p_k q_k$ . This is achieved when  $\mathbf{q}$  puts all of its mass on the  
 831 largest elements of  $\mathbf{p}$ ,

$$\underset{\mathbf{q} \in \Delta}{\operatorname{argmin}} -\log(\mathbf{p}^\top \mathbf{q}) = \{\mathbf{q} \in \Delta : q_k = 0 \forall k \notin \operatorname{argmax}(\mathbf{p})\}.$$

832 In the usual case where the largest element in  $\mathbf{p}$  is unique, we have  $\underset{\mathbf{q} \in \Delta}{\operatorname{argmin}} -\log(\mathbf{p}^\top \mathbf{q}) = \mathbf{e}_{\hat{y}}$ .

833 Now compute the gradient of LDP loss with respect to logits:

$$\begin{aligned}
 834 \quad \nabla_{\mathbf{z}} - \log(\mathbf{p}^\top \mathbf{q}(\mathbf{z})) &= -\frac{\nabla_{\mathbf{z}} \mathbf{q}(\mathbf{z})^\top \mathbf{p}}{\mathbf{p}^\top \mathbf{q}(\mathbf{z})} \\
 835 &= -\frac{(\operatorname{diag}(\mathbf{q}(\mathbf{z})) - \mathbf{q}(\mathbf{z})\mathbf{q}(\mathbf{z})^\top) \mathbf{p}}{\mathbf{p}^\top \mathbf{q}(\mathbf{z})} \\
 836 &= -\frac{\mathbf{p} \odot \mathbf{q}(\mathbf{z}) - \mathbf{q}(\mathbf{z})\mathbf{q}(\mathbf{z})^\top \mathbf{p}}{\mathbf{p}^\top \mathbf{q}(\mathbf{z})} \\
 837 &= \mathbf{q}(\mathbf{z}) - \frac{\mathbf{p} \odot \mathbf{q}(\mathbf{z})}{\mathbf{p}^\top \mathbf{q}(\mathbf{z})}.
 \end{aligned}$$

The gradient with respect to  $\theta$  is

$$\begin{aligned}
\nabla_{\theta} -\log(\mathbf{p}^{\top} \mathbf{q}(\theta)) &= -\nabla_{\theta} \log \left( \sum_{k=1}^K p_k q_k(\theta) \right) \\
&= -\frac{\nabla_{\theta} \sum_{k=1}^K p_k q_k(\theta)}{\sum_{k=1}^K p_k^{\top} q_k(\theta)} \\
&= -\frac{\sum_{k=1}^K p_k \nabla_{\theta} q_k(\theta)}{\sum_{k=1}^K p_k q_k(\theta)} \\
&= -\sum_{k=1}^K \frac{p_k}{\mathbf{p}^{\top} \mathbf{q}(\theta)} \nabla_{\theta} q_k(\theta).
\end{aligned}$$

## B PROOF OF THEOREM 1

**Theorem 1.** Let  $\mathbf{z} \in \mathbb{R}^K$  be a fixed vector of logits,  $\mathbf{q}(\mathbf{z}) = \text{softmax}(\mathbf{z}) \in \Delta$  a prediction for the target  $\mathbf{p} \in \Delta$ ,  $q_k(\mathbf{z})$  the  $k$ -th entry of vector  $\mathbf{q}(\mathbf{z})$ , and  $m \in \text{argmin}_k q_k(\mathbf{z})$ . Then, the Euclidean norm of the gradient of the LDP loss with respect to logits satisfies

$$0 \leq \|\nabla_{\mathbf{z}}[-\log(\mathbf{p}^{\top} \mathbf{q}(\mathbf{z}))]\|_2 \leq \|\mathbf{q}(\mathbf{z}) - \mathbf{e}_m\|_2 = \|\nabla_{\mathbf{z}}[-\mathbf{e}_m^{\top} \log(\mathbf{q}(\mathbf{z}))]\|_2. \quad (3)$$

In addition:

- (a) The minimum is achieved when the target is the uniform distribution  $\mathbf{p} = \mathbf{1}/K$ .
- (b) The maximum is achieved when the target  $\mathbf{p}$  is a one-hot vector for label  $m$ .

*Proof.* We use the short-hand notation  $\mathbf{q}$  to represent  $\mathbf{q}(\mathbf{z})$ . If there exist a probability vector  $\mathbf{p}$  that is a stationary point of the LDP loss, then this will minimize the gradient norm. To find the stationary point, set the gradient equal to the zero vector and rearrange:

$$\begin{aligned}
\mathbf{0} &= \mathbf{q} - \frac{\mathbf{p} \odot \mathbf{q}}{\mathbf{p}^{\top} \mathbf{q}}, \\
\mathbf{p}^{\top} \mathbf{q} \mathbf{q} &= \mathbf{p} \odot \mathbf{q}.
\end{aligned}$$

We have equality between two vectors, so each element should be equal, that is,  $\mathbf{p}^{\top} \mathbf{q} q_k = p_k q_k$  for all  $k$ . The equation is satisfied when  $q_k = 0$  and corresponding  $p_k$  can be arbitrary. However, we assume  $\mathbf{q}$  is a softmax vector, so no element can be exactly equal to zero. In the case  $q_k \neq 0$ , we can divide both sides by  $q_k$ , giving  $\mathbf{p}^{\top} \mathbf{q} = p_k$ . This implies  $p_k$  is constant for all  $k$ , or  $\mathbf{p} = \mathbf{1}/K$  is the uniform distribution.

To find the  $\mathbf{p}$  that maximizes the norm of the gradient  $\|\mathbf{q} - \mathbf{p} \odot \mathbf{q} / \mathbf{p}^{\top} \mathbf{q}\|_2$ , note that  $\mathbf{x} = \mathbf{p} \odot \mathbf{q} / \mathbf{p}^{\top} \mathbf{q}$  is a probability vector. This implies that (a)  $\|\mathbf{x}\|_2^2 \leq 1$ , where equality holds when  $\mathbf{x}$  is a one-hot vector; and (b)  $\mathbf{q}^{\top} \mathbf{x} \geq q_m$ , where equality holds when  $\mathbf{x}$  is the one-hot vector  $\mathbf{e}_m$ . Hence we have

$$\|\mathbf{q} - \mathbf{x}\|_2^2 = \|\mathbf{q}\|_2^2 - 2\mathbf{q}^{\top} \mathbf{x} + \|\mathbf{x}\|_2^2 \leq \|\mathbf{q}\|_2^2 - 2q_m + 1 = \|\mathbf{q} - \mathbf{e}_m\|_2^2,$$

where the maximum is achieved when  $\mathbf{x} = \mathbf{e}_m$ ; equivalently, when  $\mathbf{p} \odot \mathbf{q} / \mathbf{p}^{\top} \mathbf{q} = \mathbf{e}_m$ , or  $\mathbf{p} = \mathbf{e}_m$ .  $\square$

## C PROOF OF THEOREM 2

**Theorem 2.** Let  $\mathbf{q}(\theta) \in \Delta$  denote the prediction for target  $\mathbf{p} \in \Delta$  when viewed as a function of the learnable model weights  $\theta$ . The Euclidean norm of the gradient of the LDP loss with respect to model weights satisfies

$$\|\nabla_{\theta}[-\log(\mathbf{p}^{\top} \mathbf{q}(\theta))]\|_2 \leq \frac{\max(\mathbf{p}) - \min(\mathbf{p})}{\min(\mathbf{p})} \sum_{k=1}^K \|\nabla_{\theta} q_k(\theta)\|_2. \quad (4)$$

918 *Proof.* The gradient of the LDP loss with respect to model weights  $\theta$  can be read from Table 1 row  
919 3 column 5:

$$920 \quad \nabla_{\theta}[-\log(\mathbf{p}^{\top} \mathbf{q}(\theta))] = -\sum_{k=1}^K \frac{p_k}{\mathbf{p}^{\top} \mathbf{q}(\theta)} \nabla_{\theta} q_k(\theta).$$

921 Since  $\sum_{k=1}^K q_k(\theta) = 1$ , we have  $\nabla_{\theta} \sum_{k=1}^K q_k(\theta) = \mathbf{0}$ , and

$$922 \quad \nabla_{\theta}[-\log(\mathbf{p}^{\top} \mathbf{q}(\theta))] = -\frac{1}{\mathbf{p}^{\top} \mathbf{q}(\theta)} \sum_{k=1}^K (p_k - \min(\mathbf{p})) \nabla_{\theta} q_k(\theta).$$

923 Applying the triangle inequality to the norm of the gradient gives

$$924 \quad \|\nabla_{\theta}[-\log(\mathbf{p}^{\top} \mathbf{q}(\theta))]\|_2 \leq \frac{1}{\mathbf{p}^{\top} \mathbf{q}(\theta)} \sum_{k=1}^K (p_k - \min(\mathbf{p})) \|\nabla_{\theta} q_k(\theta)\|_2. \quad (7)$$

925 We can simplify the upper bound in Equation (7) by applying  $\mathbf{p}^{\top} \mathbf{q}(\theta) \geq \min(\mathbf{p})$  and  $p_k \leq \max(\mathbf{p})$   
926 to give the result

$$927 \quad \|\nabla_{\theta}[-\log(\mathbf{p}^{\top} \mathbf{q}(\theta))]\|_2 \leq \frac{\max(\mathbf{p}) - \min(\mathbf{p})}{\min(\mathbf{p})} \sum_{k=1}^K \|\nabla_{\theta} q_k(\theta)\|_2.$$

928 □

## 929 D PROOF OF PROPOSITION 1

930 **Proposition 1.** Let  $\mathbf{z}$ ,  $\mathbf{q}(\mathbf{z})$  and  $\mathbf{p}$  be defined as in Theorem 1,  $p_k$  be the  $k$ -th element of vector  $\mathbf{p}$ ,  
931  $M \in \operatorname{argmax}_k p_k$  be any maximizer of  $\mathbf{p}$ , and  $w : [1/K, 1] \rightarrow [0, 1]$  be non-decreasing. Then,  
932 the Euclidean norm of the gradient with respect to logits of loss functions of the form  $\ell(\mathbf{q}(\mathbf{z}), \mathbf{p}) =$   
933  $w(p_M)(-\log q_M(\mathbf{z}))$  satisfy the following bound:

$$934 \quad w(1/K) \|\mathbf{q}(\mathbf{z}) - \mathbf{e}_M\|_2 \leq \|\nabla_{\mathbf{z}} w(p_M)(-\log q_M(\mathbf{z}))\|_2 \leq w(1) \|\mathbf{q}(\mathbf{z}) - \mathbf{e}_M\|_2. \quad (5)$$

935 In addition:

- 936 (a) The minimum is achieved when the target is the uniform distribution  $\mathbf{p} = \mathbf{1}/K$ .
- 937 (b) The maximum is achieved when the target is any one-hot vector  $\mathbf{p} = \mathbf{e}_k$ ,  $k \in \{1, \dots, K\}$ .

938 *Proof.* The gradient of the loss can be read from Table 1 and scaled by the weight function,  
939  $\nabla_{\mathbf{z}} \ell(\mathbf{q}(\mathbf{z}), \mathbf{p}) = w(p_M)(\mathbf{q}(\mathbf{z}) - \mathbf{e}_M)$ . When viewed as a function of  $\mathbf{p}$ , the norm of this gradi-  
940 ent  $w(p_M) \|\mathbf{q}(\mathbf{z}) - \mathbf{e}_M\|_2$  takes on its extreme values when the weight function  $w$  is minimized or  
941 maximized. Since  $w$  is non-decreasing, it attains its minimum value at  $1/K$  when  $\mathbf{p} = \mathbf{1}/K$  is the  
942 uniform distribution, and attains its maximum value when  $\mathbf{p}$  is a one-hot vector with  $p_M = 1$ . □

## 943 E EXPERIMENT SETTINGS

### 944 E.1 DATASETS

945 We provide some summary statistics about each of the datasets included in the benchmark in Table 4.

### 946 E.2 HYPERPARAMETERS

947 To aid in reproducibility, we present all hyperparameters used in the experiments. Algorithm-  
948 independent hyperparameters are provided in Table 5.

949 For the black and white datasets MNIST and EMNIST, the models are modified to have one input  
950 channel rather than three, and only the subset of the augmentations in RandAugment that do not rely  
951 on color are used.

952 For algorithm dependent hyperparameters we generally copy values from the original papers: Fix-  
953 Match  $\tau = 0.95$ , SoftMatch uses adaptive  $\mu_t$  and  $\sigma_t$  in the weight function, FreeMatch fairness

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

Table 4: Benchmark Datasets.

Dataset	$n$ train	$n$ test	$K$	$d$
MNIST (Lecun et al., 1998)	60 000	10 000	10	$28 \times 28$
EMNIST (Cohen et al., 2017)	112 800	18 800	47	$28 \times 28$
CIFAR10 (Krizhevsky, 2009)	50 000	10 000	10	$3 \times 32 \times 32$
CIFAR100 (Krizhevsky, 2009)	50 000	10 000	100	$3 \times 32 \times 32$
SVHN (Netzer et al., 2011)	604 388	26 032	10	$3 \times 32 \times 32$
STL10 (Coates et al., 2011)	105 000	8 000	10	$3 \times 96 \times 96$
ImageNet (Deng et al., 2009)	1 281 167	50 000	1 000	$3 \times 224 \times 224$

Table 5: Algorithm independent hyperparameters.

Dataset	MNIST	EMNIST	CIFAR-10	CIFAR-100	STL-10	SVHN	ImageNet
Model	WRN-10-1	WRN-22-1	WRN-28-2	WRN-28-8	WRN-37-2	WRN-28-2	ResNet50
Train iters	50 000	50 000	200 000	200 000	200 000	200 000	1 048 576
Weight decay		5e-4		1e-3		5e-4	3e-4
Labeled batch size			64				128
Unlabeled batch size			64				128
Learning rate			0.03 $\cos(\frac{7\pi t}{16T})$				
SGD momentum			0.9				
EMA momentum			0.999				
Weak Augmentation			Random Crop and Flip				
Strong Augmentation			RandAugment				

loss weight  $\lambda_{SAF} = 0.01$  for MNIST, EMNIST, CIFAR10-10, CIFAR100-400, STL10-40, SVHN, and  $\lambda_{SAF} = 0.05$  for all other cases. FlatMatch perturbation magnitude  $\rho = 0.1$ , InterLUDE delta consistency loss weight  $\lambda_{DC} = 1.0$ , fusion strength  $\alpha = 0.1$ . All experiments were run on a cluster containing NVIDIA V100, A100 and A30 gpus.

The number of learnable parameters for each different architecture are provided in Table 6.

Table 6: Architectures.

Name	#Params
WRN-10-1	77 578
WRN-22-1	274 415
WRN-28-2	1 467 626
WRN-28-8	23 401 028
WRN-37-2	5 929 450
ResNet50	25 557 032

## F ADDITIONAL EXPERIMENT RESULTS

The test accuracy recorded after training is completed is presented in Table 7. The conclusions are similar as for Table 2. Overall DotMatch achieves the best or equivalent to the best performance 11 times out of 15 experiments, only matched by other state-of-the-art methods SoftMatch and FreeMatch. DotMatch either achieves the best or equivalent to the best performance in all of the most challenging settings with the least amount of labeled data, the only algorithm to do so. It also has consistently strong performance on MNIST and EMNIST, which are datasets not typically included in benchmarks from other works.

We compared the runtime of DotMatch against baseline methods and report the results in Table 8. Time is recorded as the average seconds per iteration for the first 1000 training iterations. All algorithms are implemented using the USB code base, and hardware is kept identical: 1 NVIDIA A100 40GB, 128GB RAM. DotMatch is consistently the fastest algorithm across all datasets, with FixMatch generally being the second fastest, followed by FreeMatch then SoftMatch. Interestingly, FixMatch was the slowest on ImageNet out of the four algorithms, but the rest of the rankings on ImageNet remain the same as for other datasets.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

Table 7: Final test error rates (mean $\pm$ std% over 3 seeds) of various SSL methods recorded after training completed. Best in **bold**, overlapping 95% CI underlined. FullySupervised results for STL10 omitted since not all labels available.

Dataset	MNIST			EMNIST			CIFAR10			CIFAR100			SVHN			STL10	
# Label	10	40	250	47	188	1175	10	40	250	400	2500	40	250	40	1000	1000	
Supervised	71.31 $\pm$ 5.54	30.36 $\pm$ 2.97	8.98 $\pm$ 1.37	85.69 $\pm$ 1.00	52.68 $\pm$ 1.00	23.30 $\pm$ 0.65	84.05 $\pm$ 1.46	79.15 $\pm$ 1.15	62.22 $\pm$ 0.36	89.88 $\pm$ 1.23	61.63 $\pm$ 0.52	87.37 $\pm$ 1.55	32.68 $\pm$ 0.72	77.95 $\pm$ 1.93	33.81 $\pm$ 1.55	-	
FixMatch	44.41 $\pm$ 12.81	<u>4.96</u> $\pm$ 3.98	1.34 $\pm$ 0.04	<u>48.24</u> $\pm$ 2.07	24.14 $\pm$ 1.37	<u>16.11</u> $\pm$ 0.76	85.84 $\pm$ 0.01	22.10 $\pm$ 0.33	<u>5.46</u> $\pm$ 0.16	52.58 $\pm$ 1.33	29.93 $\pm$ 0.12	<b>5.88</b> $\pm$ 4.10	<u>2.37</u> $\pm$ 0.08	48.00 $\pm$ 8.67	<u>7.07</u> $\pm$ 0.38	-	
SoftMatch	<u>8.51</u> $\pm$ 2.93	<u>3.69</u> $\pm$ 2.39	1.38 $\pm$ 0.04	49.32 $\pm$ 2.91	<b>20.60</b> $\pm$ 0.39	<u>15.61</u> $\pm$ 0.34	<u>42.17</u> $\pm$ 1.95	6.91 $\pm$ 0.29	<u>5.51</u> $\pm$ 0.07	<b>45.65</b> $\pm$ 1.33	<u>28.73</u> $\pm$ 0.41	<u>24.07</u> $\pm$ 12.61	2.55 $\pm$ 0.04	34.71 $\pm$ 9.71	<b>6.60</b> $\pm$ 0.11	-	
FreeMatch	<b>4.97</b> $\pm$ 2.49	<u>2.88</u> $\pm$ 1.14	1.37 $\pm$ 0.06	54.81 $\pm$ 1.22	<u>21.09</u> $\pm$ 2.12	<u>15.86</u> $\pm$ 0.29	32.56 $\pm$ 2.83	<b>5.51</b> $\pm$ 0.09	<b>5.39</b> $\pm$ 0.21	<u>46.57</u> $\pm$ 2.48	<b>28.47</b> $\pm$ 0.28	<u>28.77</u> $\pm$ 10.07	13.80 $\pm$ 1.43	34.32 $\pm$ 20.76	<u>6.64</u> $\pm$ 0.07	-	
FlatMatch	26.49 $\pm$ 23.12	<b>2.45</b> $\pm$ 1.06	1.70 $\pm$ 0.06	53.84 $\pm$ 4.27	<u>21.33</u> $\pm$ 1.66	<u>15.84</u> $\pm$ 0.27	30.92 $\pm$ 21.46	6.54 $\pm$ 0.56	6.10 $\pm$ 0.34	<u>47.46</u> $\pm$ 1.72	29.58 $\pm$ 0.11	21.45 $\pm$ 7.23	13.69 $\pm$ 0.15	<b>21.12</b> $\pm$ 1.32	7.94 $\pm$ 0.31	-	
Interlude	67.51 $\pm$ 4.62	<u>13.30</u> $\pm$ 0.32	<b>1.10</b> $\pm$ 0.10	73.32 $\pm$ 4.43	32.89 $\pm$ 2.96	17.08 $\pm$ 0.27	<u>55.87</u> $\pm$ 2.08	28.19 $\pm$ 4.07	<u>5.80</u> $\pm$ 0.38	60.41 $\pm$ 3.08	48.36 $\pm$ 1.62	33.25 $\pm$ 15.24	<b>2.26</b> $\pm$ 0.03	43.66 $\pm$ 5.86	7.92 $\pm$ 0.13	-	
DotMatch	6.40 $\pm$ 3.46	<u>2.98</u> $\pm$ 2.46	<u>1.25</u> $\pm$ 0.18	<b>44.76</b> $\pm$ 1.82	<u>20.83</u> $\pm$ 1.16	<b>15.55</b> $\pm$ 0.87	<b>29.09</b> $\pm$ 18.63	10.71 $\pm$ 2.77	<u>5.86</u> $\pm$ 0.33	<u>46.04</u> $\pm$ 0.65	29.51 $\pm$ 0.58	<u>16.48</u> $\pm$ 14.47	2.53 $\pm$ 0.08	<u>36.40</u> $\pm$ 12.41	7.23 $\pm$ 0.09	-	
FullySupervised		0.93 $\pm$ 0.06			11.42 $\pm$ 0.11			5.05 $\pm$ 0.12		26.93 $\pm$ 0.00		2.49 $\pm$ 0.20				-	

Table 8: Training time (seconds/iteration) for various SSL methods. Best in **bold**, second best in underline. Algorithms are implemented using a unified codebase and run using identical hardware.

Dataset	MNIST	EMNIST	CIFAR10	CIFAR100	SVHN	STL10	ImageNet
# Label	250	188	250	400	250	1000	100 000
FixMatch	<u>0.284</u>	<u>0.303</u>	<u>0.415</u>	<u>0.556</u>	<u>0.412</u>	<u>0.679</u>	1.758
SoftMatch	0.289	0.313	0.444	0.666	0.443	0.763	1.620
FreeMatch	0.286	0.305	0.418	0.563	0.416	0.694	<u>1.579</u>
DotMatch	<b>0.280</b>	<b>0.293</b>	<b>0.393</b>	<b>0.460</b>	<b>0.396</b>	<b>0.624</b>	<b>1.538</b>