DMT: Dynamic Mutual Training for Semi-Supervised Learning

Zhengyang Feng^{a,1}, Qianyu Zhou^{a,1}, Qiqi Gu^a, Xin Tan^a, Guangliang Cheng^b, Xuequan Lu^c, Jianping Shi^b, Lizhuang Ma^a

^a Shanghai Jiao Tong University, Shanghai, China
^b Sense Time Research, Shanghai, China
^c Deakin University, Australia

Abstract

Recent semi-supervised learning methods use pseudo supervision as core idea, especially self-training methods that generate pseudo labels. However, pseudo labels are unreliable. Self-training methods usually rely on single model prediction confidence to filter low-confidence pseudo labels, thus remaining high-confidence errors and wasting many low-confidence correct labels. In this paper, we point out it is difficult for a model to counter its own errors. Instead, leveraging inter-model disagreement between different models is a key to locate pseudo label errors. With this new viewpoint, we propose mutual training between two different models by a dynamically re-weighted loss function, called Dynamic Mutual Training (DMT). We quantify inter-model disagreement by comparing predictions from two different models to dynamically re-weight loss in training, where a larger disagreement indicates a possible error and corresponds to a lower loss value. Extensive experiments show that DMT achieves state-of-theart performance in both image classification and semantic segmentation. Our codes are released at https://github.com/voldemortX/DST-CBC.

Keywords: dynamic mutual training, inter-model disagreement, noisy pseudo label, semi-supervised learning

¹Equal contribution.

1. Introduction

In recent years, with the rise of deep learning, substantial improvements have been shown in various computer vision tasks, e.g. image classification [1, 2] and semantic segmentation [3, 4]. However, deep learning methods require a large amount of annotated data to learn generalized representations. Although a large-scale dataset is easily gathered from cameras or web pages, the labor cost for labeling such a dataset has become unbearable in many real-world applications. For example, it takes 1.5 hours for a human annotator to label a high-resolution image of urban street scenes with pixel-wise annotations [5]. In this work, we focus on semi-supervised learning to alleviate the label costs, by taking semantic segmentation and image classification as examples.

Semi-supervised learning labels only a small part of the dataset (labeled subset), and exploits the remaining part as unlabeled data (unlabeled subset). To learn without labels, a natural idea is "bootstrapping" (pulling oneself up by one's own bootstraps) [6], i.e. using self (pseudo) supervisions. Two lines of approaches have achieved good performance on both semi-supervised image classification and semantic segmentation: entropy minimization (i.e. self-training) [7, 8] and consistency regularization [9, 10]. Recently, hybrid methods that combine those two directions, MixMatch with sharpening [11], s4GAN + MLMT with separate network branches [12], show state-of-the-art performance on image classification and semantic segmentation, respectively. Our method is based on offline self-training with data augmentation as consistency regularization, a hybrid method applicable to both tasks (Section 2.2).

Nevertheless, bootstrapping methods face a common issue, that is, pseudo supervisions tend to have classification errors. To address this, previous self-training methods [7, 8, 13] select pseudo labels by confidence, e.g. the predicted probability from a model trained on the labeled subset. However, methods based on the common assumption that higher confidence corresponds to cleaner labels still have drawbacks. We conduct a pseudo labeling experiment (Fig. 1) to illustrate this issue. As shown in Fig. 1, using confidence to select pseudo

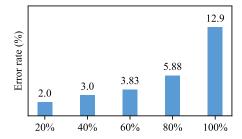


Figure 1: Pseudo label error statistics. We report pseudo label error rates on 1,000 random images from CIFAR-10, using a 28-layer WideResNet model trained with 4,000 labeled samples only. The overall error rate is 12.90%, error rates for top-20%, top-40%, top-60%, top-80% images according to prediction confidence are also plotted. It can be observed that high-confidence errors do exist and lots of data will be discarded to achieve a low error rate, e.g. < 3% means discarding 60% data.

labels suffers from two limitations. First, low-confidence correct pseudo labels are often ignored, i.e. in order to achieve a low label error rate for pseudo supervision, a large portion of low-confidence correct pseudo labels have to be discarded. Second, high confidence errors do exist. We can observe that even pseudo labels with top-20% confidence still have some errors. Moreover, pseudo supervision error from the model itself (termed as *self-error*) can be extremely harmful in semi-supervised learning (Section 3).

To address these limitations, we propose a new viewpoint from inter-model disagreement. In particular, no matter what pseudo label selection metric is employed, there is only one model to counter its own errors. Instead, two different models with disagreements on classification decisions may be able to identify each other's errors. For instance, in image classification, model A provides a pseudo label on unlabeled image x for model B to learn, and we can quantify the disagreement between A and B by their prediction statistics and assign lower sample importance to this image if their disagreement is larger. Since the possibility of different models confidently making the same mistakes is lower, most incorrect pseudo labels will have less impact on learning.

To quantify and exploit the inter-model disagreement between different mod-

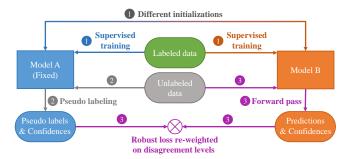


Figure 2: Overview of our training framework. There are two different models trained on the labeled subset, and one model provides pseudo supervisions for the other. A noise-robust loss (dynamic loss) is introduced in the semi-supervised training stage, leveraging the intermodel disagreement based on two models' predictions and confidences. The training process is conducted sequentially along the serial numbers 1-3.

els, we propose Dynamic Mutual Training (DMT) with a noise-robust loss (Fig. 2). First, we instantiate two different models by different weight initializations or training on two different labeled subsets. Then, one model provides unreliable pseudo labels for the other on unlabeled data. Finally, to quantify the inter-model disagreement when using those noisy pseudo labels, we define three disagreement cases and corresponding loss re-weighting strategies based on the relation of prediction confidence between the two models. In this way, loss value is dynamically weighted lower when the disagreement is larger. Thus, we can endure label noise in training. For semantic segmentation, by assigning higher weights on high-quality pixels in an image and suppressing low-quality pixels, pseudo supervision quality can be overall improved on each unlabeled image, leading to larger performance gains. Furthermore, inspired by curriculum learning, or easy to hard [14], we apply DMT iteratively to gradually exploit unlabeled data for better performance.

Note that other disagreement-based semi-supervised learning methods use different models and learn by maximizing their agreement on unlabeled data [15, 16]. In contrast, we view disagreement [17] as a principle, i.e. the inter-model disagreement provides **the possibility of learning** (Fig. 3). So long as this principle is leveraged properly, the exact method formulation can vary. Thus,

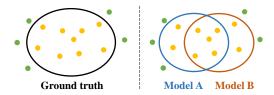


Figure 3: An illustration of the inter-model disagreement on a toy binary classification experiment. Positive and negative samples are marked as yellow and green, respectively. Model A and B have inter-model disagreements that enable the possibility to improve the performance up to the ground truth.

in this paper we employ the inter-model disagreement to specifically combat pseudo label noise by loss re-weighting, rather than penalizing disagreement as a learning objective. Since our method formulation mostly concerns the loss function, it is also compatible with other approaches, e.g. entropy minimization and consistency regularization.

Our main contributions are summarized as follows.

- We analyze the pseudo label noise problem for semi-supervised learning and propose a new method from the new viewpoint of inter-model disagreement, i.e. instead of single model confidence, disagreement between models may indicate possible pseudo label errors.
- To quantify the inter-model disagreement, we propose a general and efficient bootstrapping approach, called Dynamic Mutual Training (DMT).
 DMT exploits the relation between different model predictions by a noise-robust loss function where a larger inter-model disagreement corresponds to a lower loss weighting. The performance of DMT is further enhanced by casting it into an iterative framework.
- We demonstrate the effectiveness of our approach in different tasks and datasets, i.e. semi-supervised image classification on CIFAR-10 and semisupervised semantic segmentation on PASCAL VOC 2012 and Cityscapes.
 Through extensive comparisons and ablations, the proposed method shows

state-of-the-art performance on both tasks. In the harder semantic segmentation task, our method even surpasses manual annotation under a certain setting.

2. Preliminaries

Semi-supervised learning methods are often based on certain prior knowledge, thus models can "bootstrap" themselves with extra unlabeled data for better generalization. There are mainly three types of prior knowledge.

- 1. **Entropy minimization.** Predictive entropy is minimized for a model to make decisions on unlabeled data, which is apparently better than not making any decisions at all [18].
- 2. Consistency regularization. Prediction should remain consistent when unlabeled data is perturbed by data augmentation [9].
- 3. **Disagreement-based.** Multiple classifiers should reach an agreement on unlabeled data predictions [15].

First, we focus on entropy minimization [18] and summarize self-training methods as two types, online and offline. Then, we explain how consistency regularization can be integrated. Finally, in Section 4, we further integrate the disagreement principle in our proposed method.

2.1. Entropy Minimization and Self-training

Entropy H is defined as:

$$H = -\sum_{c=1}^{C} p_c \log p_c , \qquad (1)$$

where p_c is the predicted probability for class c, and C is the total number of classes. Since $\sum_{c=1}^{C} p_c = 1$, H approaches the minimum value of 0 when one class is 1 and other classes are 0. Thus, entropy minimization encourages the model to make a certain decision.

Self-training takes the most probable class as a pseudo label and train models on unlabeled data, which is a common approach to achieve the minimum entropy. Note that here we do not consider soft pseudo labels (labels as probability vectors instead of a hard label or one-hot vector), since they do not directly correspond to entropy minimization and we do not observe decent performance of using soft pseudo labels.

Denote $c^* \leftarrow \arg\max_c F(c|x)$, pseudo label l is defined as:

$$l = \begin{cases} c^*, & F(c^*|x) > T\\ ignored, & otherwise \end{cases}$$
 (2)

where x is an image and $F(\cdot)$ is a classification model that predicts a probability distribution. T is the threshold for selection (e.g. fixed value such as 0.5 or ranked value such as the 20-th percentile). We summarize self-training for semi-supervised learning as two types by when pseudo labels are generated.

Online self-training generates pseudo labels after each forward pass in a network. Pseudo labels are directly selected based on some selection metric online and provide supervision for the immediate backward pass [7].

Offline self-training firstly generates all pseudo labels, with a model trained only on the labeled subset (often followed by some form of selection process). Then the model is fine-tuned/re-trained on all labels (pseudo and human-labeled). This procedure can be applied iteratively by relabeling unlabeled data with the most recently trained model [13].

2.2. Self-training and Consistency Regularization

We find that the consistency regularization plays a similar role to data augmentation in general. As shown in Fig. 4, with the same set of augmentation transforms, general data augmentation in offline self-training is similar to an "anchored" version of explicit consistency-based methods such as Mean Teacher [9]. Anchoring may bring better performance overall if the anchor is positive, and degradation if the anchor is negative. In summary, consistency regularization and data augmentation have similar principles, and we observe good performance with data augmentation in offline self-training, which is the same as fully-supervised training after pseudo labeling.

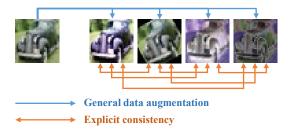


Figure 4: Comparison between general data augmentation and explicit consistency regularization. It can be observed that they have similar principles and similar outcomes. General data augmentation can be seen as a simple form of "anchoring". Since the pseudo label on the unperturbed image has higher accuracy, it should be used as a semantic anchor for other perturbations. A similar concept of augmentation anchoring is also mentioned in [19].

However, for online self-training, it can be non-trivial to impose consistency regularization. Multiple different perturbations on input require multiple forward passes and induce higher computational cost, as demonstrated in many consistency-based methods [20, 9]. Online self-training is also more sensitive to pseudo label noise, which will be elaborated in Section 3. Thus, we investigate the offline self-training case in our method.

3. Noisy Pseudo Label

Pseudo labels are noisy, since they are not generated by human annotators. Pseudo label noise is unique, coming from the model itself, different from random noise and noise from crowd-sourcing or search engines which are widely explored and modeled according to the types and levels of noise [21, 22, 23, 24]. Note that in other semi-supervised learning methods without explicit pseudo labels, self-supervision noise also exists. For instance, in consistency-based Mean Teacher [9] where the Mean Squared Error (MSE) loss is used to enforce consistency of student model output to teacher model output, the pseudo-supervision (probability distribution) provided by the teacher is noisy. In this work, we discuss noisy pseudo labels (one-hot vectors) in self-training, which is more straightforward.

To address pseudo label noise, self-training methods adopt prediction confidence as noise indicator, i.e. selecting high confidence pseudo labels. The simplest policy is confidence thresholding, where only pseudo labels with confidence higher than a fixed threshold are selected [7]. Zou et al. further use different thresholds for each class [13]. Hung et al. consider the confidence of a specialized two-class discriminator to discriminate between real and fake (wrong) labels [8]. Although noise rate is roughly lower when confidence is higher, it is difficult to attain sufficiently clean pseudo labels without discarding a large portion of unlabeled data (Fig. 1), proved by observations on semantic segmentation tasks [8] where only $27\% \sim 36\%$ pixels can be pseudo-labeled without performance degradation on PASCAL VOC 2012 and [25] where they only use $30\% \sim 50\%$ pseudo labels on Cityscapes. This goes against the purpose of semi-supervised learning, which is using rather than discarding unlabeled data. By contrast, we define different weights to samples instead of discarding them.

Other than wasting unlabeled data, selected pseudo labels still have some noise. This is particularly troublesome for online self-training. Because modern networks are trained with random and heavy data augmentation and the model in training changes after each parameter update, its predictions and errors also change. The change of errors can be extremely confusing when a model keeps minimizing entropy (fitting) on new errors, resulting in self-error accumulation throughout training. For offline self-training, although self-errors remain an obstacle, it is not as severe as that in online self-training, since the pseudo label errors are fixed throughout training. However, there is always the same dilemma: a model cannot possibly correct itself, except for ignoring some uncertain predictions. Thus, we introduce a new viewpoint from inter-model disagreement as illustrated in Fig. 3. Since it is possible to use one model to correct another, and obviously different models with sufficient disagreement between them are less likely to make the same mistakes, especially the same high-confidence mistakes.

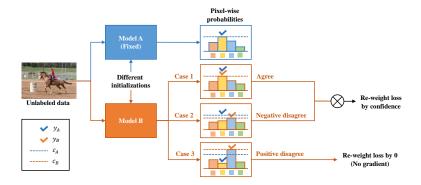


Figure 5: Dynamic Mutual Training (DMT). There are two sufficiently different models (F_A, F_B) . F_A is used to generate pseudo labels (y_A) and record corresponding confidences (c_A) before mutual training (procedure marked by blue lines), and F_B (with prediction y_B and confidence c_B) is then trained by them. There are three possible cases in mutual training and three corresponding loss re-weighting strategies based on the two models' disagreement degree, defined in Eq. 3. All histograms above are for illustration purposes, rather than real outputs.

4. Method

In this section, we first present Dynamic Mutual Training (DMT) in Section 4.1 with the dynamic loss (Section 4.1.1) and techniques for how to initialize models with disagreement in different tasks (Section 4.1.2). Then we cast DMT to an iterative learning framework for better performance (Section 4.2).

4.1. Dynamic Mutual Training

We propose Dynamic Mutual Training (DMT), to quantify the inter-model disagreement and enable noise-robust training, illustrated in Fig. 5. First, we train two different models F_A and F_B on the labeled subset from two different initializations/sub-samplings. Then, one model, e.g. F_A , is fixed and generates pseudo labels and confidences on the unlabeled subset. And the other, F_B , fine-tunes on all data (labeled and pseudo labeled) with our dynamically weighted cross-entropy loss. In the same way, F_B can train F_A .

4.1.1. The Dynamic Loss

We propose the dynamic loss, where the quantified inter-model disagreement serves as the dynamic loss weight. Taking image classification for example, we assume F_A trains F_B , let \mathcal{X} , \mathcal{U} denote labeled and unlabeled (pseudo labeled) samples in a batch of size N and let u be an unlabeled image in \mathcal{U} , we define its pseudo label as $y_A \leftarrow \arg\max_y F_A(y|u)$, with confidence $c_A \leftarrow F_A(y_A|u)$. And the prediction in training is $y_B \leftarrow \arg\max_y F_B(y|u)$, with confidence $c_B \leftarrow F_B(y_B|u)$. Let $p_B \leftarrow F_B(y_A|u)$ be the predicted probability of class y_A by F_B . The dynamic loss weight ω_u is defined as:

$$\omega_{u} = \begin{cases}
p_{B}^{\gamma_{1}}, & y_{A} = y_{B} \\
p_{B}^{\gamma_{2}}, & y_{A} \neq y_{B}, c_{A} \geq c_{B} \\
0, & y_{A} \neq y_{B}, c_{A} < c_{B}.
\end{cases}$$
(3)

The dynamic loss on unlabeled samples $\mathcal{L}_{\mathcal{U}}$ is then defined as:

$$\mathcal{L}_{\mathcal{U}} = \frac{1}{N} \sum_{u, y_A \in \mathcal{U}} \omega_u CE(y_A, F_B(u)), \tag{4}$$

where $CE(\cdot)$ is the cross-entropy loss.

Intuitively, for the pseudo labeled data, there are three different cases in training:

- 1. **Agreement.** F_B agrees with the pseudo label.
- 2. Negative disagreement. F_B disagrees with the pseudo label but the confidence on F_B 's decision is lower than the pseudo label's.
- 3. Positive disagreement. F_B disagrees with the pseudo label and has higher confidence.

In cases 1 and 2, we use the current model's predicted probability p_B on the pseudo labeled class as weight, perceived as the quantified disagreement, i.e. a higher p_B means F_B has a higher agreement with F_A . In case 3, we set the dynamic weight to 0 because the pseudo label is probably incorrect.

Dynamic weights are further re-scaled by hyper-parameters γ_1, γ_2 , a higher γ magnifies confidence differences (take smaller steps on low-confidence examples) and suppresses gradients overall. It can be interpreted that a relatively higher γ_1 represents a more emphasized entropy minimization, a higher γ_2 represents a more emphasized mutual learning. High γ values are often better for high-noise scenarios, or to maintain larger inter-model disagreement.

Note that training uses the labeled subset along with the pseudo-labeled data, and the loss for labeled data $\mathcal{L}_{\mathcal{X}}$ remains unchanged, i.e. the typical cross-entropy loss:

$$\mathcal{L}_{\mathcal{X}} = \frac{1}{N} \sum_{x,gt \in \mathcal{X}} CE(gt, F_B(x)), \tag{5}$$

where x and gt denote image and ground truth pairs. The combined loss \mathcal{L} is defined as:

$$\mathcal{L} = \mathcal{L}_{\mathcal{X}} + \mathcal{L}_{\mathcal{U}}.\tag{6}$$

The above example is given on classification. With regard to semantic segmentation, $\omega_u^{H\times W}$ is a pixel-wise map (H for height and W for width), the re-weighting strategy remains the same and applies on each pixel.

Remark. Naively, one can directly use p_B as the dynamic weight without distinguishing the three cases or scaling. But we observe that the naive approach does not work well under severe pseudo label noise. Also, other design of the dynamic loss brings no noticeable improvement (Section 5.4).

4.1.2. Initialize Disagreement

A key problem for leveraging the inter-model disagreement is how to initialize sufficiently different models. For simpler tasks such as CIFAR-10 image classification, it is simple to just randomly initialize different models for sufficient disagreement. However, for tasks that require pre-trained weights to work well, e.g. semantic segmentation, sufficiently different off-the-shelf pre-trained weights are hard to obtain, and the extra amount of time needed for a new pre-training is too costly compared to the task at hand. Thus, we mainly use

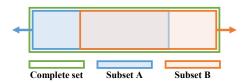


Figure 6: Difference maximized sampling. The complete set is randomly shuffled at first, and subset A and B are drawn with an equal size but with fewest overlapped samples.

pre-trained weights from different datasets (ImageNet and COCO) in semantic segmentation. For some extreme cases where labeled data is scarce and one set of pre-trained weights is clearly superior (about 100-200 labeled images), we use the better pre-trained weights, and train two models from different sub-subsets of the labeled subset by difference maximized sampling (Fig. 6).

4.2. Iterative Framework

In this section we cast DMT into an iterative framework for better performance. Curriculum learning [14], or easy to hard, has been explored in semi-supervised image classification [26] and unsupervised domain adaptive semantic segmentation (a similar setup to semi-supervised learning) [13] for better bootstrapping performance. Specifically, the same bootstrapping algorithm is repeated for multiple iterations, each iteration explores a harder setting, e.g. more pseudo labels with lower confidence. Inspired by their successes, we also perform DMT iteratively to achieve better performance.

```
Algorithm 1: Pseudo code for iterative DMT process in image classification.

Input: Unlabeled subset S_u, labeled subset S_l.

Output: Final best model F.

1 Randomly initialize F^0

2 Train F^0 on S_l

3 \alpha = \{20\%, 40\%, 60\%, 80\%, 100\%\}

4 foreach iteration i \in \{1, 2, 3, 4, 5\} do

5 Predict and save top \alpha_i images on S_u with F^{i-1} \rightarrow pseudo labeled set S_p

Randomly initialize F^i with a previously unused random seed

7 Train F^i on both S_l and latest S_p with the dynamic loss

8 F = F^5
```

4.2.1. Image Classification

For this task, we first train on the labeled subset, then conduct DMT iteratively for multiple times; each time we select more top-confident pseudo labels from the unlabeled subset and re-train a randomly initialized model for sufficient disagreement, same as concurrent work Curriculum Labeling [26]. Pseudo code is shown in Alg. 1. However, the model re-trained from scratch provides little meaningful information at early training stage, thus we use a sigmoid-like function for γ values inspired by [9]. Concretely, with the total training steps t_{max} , at step t, $\gamma = \gamma_{max} e^{5(1-\frac{t}{t_{max}})^2}$.

```
Algorithm 2: Pseudo code for iterative DMT process in semantic segmentation.
   Input: Unlabeled subset S_u, labeled sub S_l.
   Output: Final best model F.
 1 if start from different pre-trained weights then
       Initialize F_A^0 and F_B^0 with different pre-trained weights
       Train F_A^0 on S_l
       Train F_B^0 on S_l
 4
 5 else
       Initialize F_A^0 and F_B^0 with the same pre-trained weights
       S_A, S_B = DifferenceMaximizedSampling(S_l)
       Train F_A^0 on S_A
       Train F_B^0 on S_B
10 \alpha = \{20\%, 40\%, 60\%, 80\%, 100\%\}
11 foreach iteration i \in \{1, 2, 3, 4, 5\} do
       Predict and save top \alpha_i pixels from each classes on S_u with F_A^{i-1} \to \text{pseudo labeled set } S_p
       Fine-tune F_A^i from F_A^{i-1} on both S_l and latest S_p with the dynamic loss
13
       Predict and save top \alpha_i pixels from each classes on \mathcal{S}_u with F_R^{i-1} \to \text{pseudo labeled set } \mathcal{S}_p
14
       Fine-tune F_B^i from F_B^{i-1} on both S_l and latest S_p with the dynamic loss
15
16 F = \text{best}(F_A^5, F_B^5)
```

4.2.2. Semantic Segmentation

Motivated by the fact that some classes are much easier to learn than others in semantic segmentation, CBST [13] proposed an iterative self-training scheme by using more top-confident pseudo labels from each class at each iteration. Furthermore, unlike image classification, fine-tuning performs reasonably well in semantic segmentation, and converges faster than re-training. Inspired by CBST, we conduct two separate fine-tunings between two differently initialized

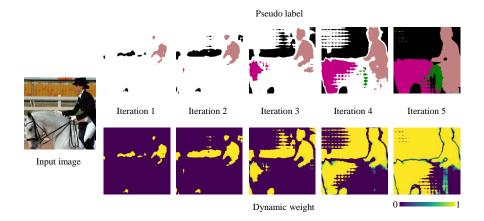


Figure 7: Illustration of iterative DMT training on PASCAL VOC 2012. We show pseudo labels from one model and dynamic weights generated when training another. At each iteration, more pixels are pseudo labeled. For some apparently incorrect regions we can observe rather low dynamic weights. Models are trained with 1/20 manually labeled data in the training set. White regions are ignored in labels and corresponding dynamic weights are shown as 0.

models, as shown in Alg. 2. In this setup, two models train each other equally and the inter-model disagreement is more straightforwardly utilized. To better illustrate how DMT is used iteratively for this harder task, we provide some examples of pseudo labels and dynamic weights during training in Fig. 7.

However, there is some difference between our iterative framework and CBST, since CBST does not select top-confident pseudo labels by direct ranking like us. Instead, it first uses class-wise thresholds (defined by ranking) to re-normalize softmax predictions by dividing the thresholds class-wise, then pixels with confidence over 1 are selected. In most cases, this is the same as direct ranking, while in extreme cases, the predicted class would change, e.g. for a binary classification task, softmax result [0.6, 0.4] that originally predicts class 0, re-normalized by ranked thresholds [0.61, 0.39], will be changed to prediction class 1. This kind of mistake happen only when using nearly all pseudo labels, which does not affect the original CBST experiments since they use less than half pseudo labels. While we use up to all data, re-normalization brings a slight degradation to final

performance. Therefore, we use direct ranking instead of re-normalization.

We set 5 iterations to our method for training. Since fine-tuning converges very fast, the total number of training steps for each one of the two models remains similar to a fully-supervised training on the entire dataset.

5. Experiments

In this section, we first specify dataset configurations (Section 5.1) and implementation details (Section 5.2). Then we compare the proposed DMT with state-of-the-art methods on image classification and semantic segmentation (Section 5.3). Finally, we analyze DMT by conducting a set of ablation studies (Section 5.4).

5.1. Datasets

For image classification, we employ the commonly used CIFAR-10 [27] dataset. For semantic segmentation, we evaluate our method on the popular PASCAL VOC 2012 [28] and Cityscapes [5] datasets.

- CIFAR-10. The CIFAR-10 [27] dataset has 10 classes, 50,000 training samples and 10,000 test samples. Most methods extract 5,000 samples from the training set to use as validation. While we consider using a realistic [29] validation set with 200 samples (see supplementary materials for how we validate our method properly with a limited validation set).
- PASCAL VOC 2012. The original PASCAL VOC 2012 [28] semantic segmentation dataset has 21 classes, 1,464 training samples and 1,449 validation samples (val) featuring common objects. We use the SBD [30] augmented version with 10,582 training samples following common practice.
- Cityscapes. The Cityscapes [5] dataset has 19 classes, 2,975 training samples, and 500 validation samples (val) for urban driving scenes. Following [12], we down-sample the images by half to 512 × 1024.

5.2. Implementation Details

Since DMT in an iterative framework has better performance and overall fast convergence as described in Section 4.2, all DMT experiments are conducted with the default 5 iterations. Our method is implemented on PyTorch with mixed-precision training [31]. All experiments were conducted on a single RTX-2080 Ti GPU. The fully-supervised learning result on the entire dataset is termed as Oracle, i.e. the performance upper-bound for semi-supervised learning. However, this upper-bound only holds when human annotations have good quality across the entire dataset. We show how this supposed upper-bound can be surpassed by our semi-supervised learning in Section 5.3.2.

5.2.1. Network Architectures

- Image classification. We follow MixMatch [19] and use a shallow residual network WideResNet-28-2 (WRN-28-2) [1] as the backbone.
- Semantic segmentation. We follow [12, 8] and use DeepLab-v2 ResNet-101 [3] as the backbone, without multi-scale fusion or CRF post-processing. Our implementation has slightly better performance than previous works, which is better aligned with the original DeepLab-v2 paper (74.75% averaged mean IoU on PASCAL VOC 2012 val set, higher than 73.6% reported in [8]).

5.2.2. Training

• Image classification. Each DMT iteration has 750 epochs with a learning rate of 0.1, a weight decay of 5×10^{-4} , a momentum of 0.9, the cosine annealing technique and a batch size of 512, which is the same as Curriculum Labeling [26]; we do not use SWA [32] for fair comparisons with other methods. Data augmentations are RandAugment [33] with Cutout. We randomly select one augmentation operation with random intensity at each step to avoid hyper-parameter tuning (number of operations n and intensity m are hyper-parameters in RandAugment). We also use mixup [34] by interpolating dynamic weights along with input images.

Table 1: Hyper-parameter settings.

	dataset	labeled ratio	γ_1	γ_2	learning rate	training	epochs	batch size	batch ratio	augmentations
1	PASCAL VOC	1/8	5	5	1×10^{-3}	fine-tuning	5	8	7:1	
2	PASCAL VOC	1/20	5	5	1×10^{-3}	fine-tuning	4	8	7:1	random scale
3	PASCAL VOC	1/50	5	5	1×10^{-3}	fine-tuning	4	8	7:1	
4	PASCAL VOC	1/106	5	5	1×10^{-3}	fine-tuning	4	8	7:1	random crop
5	Cityscapes	1/8	3	3	4×10^{-3}	fine-tuning	10	8	3:1	random horizontal flip
6	Cityscapes	1/30	3	3	4×10^{-3}	fine-tuning	8	8	7:1	
7	CIFAR-10	4k labels	4	4	1×10^{-1}	re-training	750	512	7:1	random augmentation
8	CIFAR-10	1k labels	4	4	1×10^{-1}	re-training	750	512	31:1	with random intensity

• Semantic segmentation. Each DMT iteration has less training steps due to fine-tuning. We use SGD with a momentum of 0.9, the *poly* learning rate schedule and a batch size of 8. Data augmentations include random scaling, random cropping and random flipping. We train and pseudo label at a spatial resolution of 321×321 (PASCAL VOC 2012) and 256×512 (Cityscapes).

To avoid too much hyper-parameter tuning, we set $\gamma_1 = \gamma_2$. Pseudo-labeled data are used along with labeled data in DMT, thus we have a certain ratio (labeled: unlabeled, e.g. 1:7) to combine them in a batch. More details are listed in Tab. 1.

5.2.3. Testing

- Image classification. We report the 5-times averaged *test* set performance with an exponential moving averaged (EMA) network on CIFAR-10 following MixMatch [11].
- Semantic segmentation. We report the three-times averaged val set mean intersection-over-union (mean IoU) in semantic segmentation tasks following common practice, provided the test set labels for these datasets are not publicly available.

5.3. Comparisons

To show the effectiveness and generality of DMT, we compare it with state-of-the-art methods on both image classification and semantic segmentation benchmarks as detailed in Section 5.1: CIFAR-10 [27], PASCAL VOC 2012 [28] and Cityscapes [5]. Standard practice for evaluating semi-supervised learning on these datasets is to treat most of a dataset as the unlabeled subset and use a small portion as the labeled subset.

5.3.1. CIFAR-10

For CIFAR-10, we compare our method with consistency-based Mean Teacher (MT) [9], self-training method Curriculum Labeling (CL) [26], methods explicitly/implicitly using multiple models, Deep Co-Training between two models (DCT) [15] and Dual Student (DS) [35], strong hybrid method MixMatch [11], and the combination of graph-based pseudo label propagation techniques in Density Aware Graph-based framework (DAG) [36]. Methods are evaluated on the commonly adopted 1,000 labels and 4,000 labels splits. Supervised performance with mixup and strong data augmentation on the labeled subset is reported as Baseline. Baseline, CL (our re-implementation, without SWA) and DMT are implemented in the same codebase, while for other methods with WRN-28-2 we take the reported numbers from MixMatch. The remaining numbers are taken from the original papers.

As shown in Tab. 2, DMT steadily improves CL with the dynamic loss, larger performance gain is shown in harder setting (smaller labeled subset). Note that CL is already a very high-performance method (only 2.33% less than Oracle performance using 4,000 labels), it is rather difficult to gain further improvements in a strictly controlled comparison such as ours. While our method still achieves a slight improvement. We also compare DMT with more methods in Tab. 3, where DMT shows better performance than state-of-the-art methods. We are aware that recently ReMixMatch [19] has obtained an accuracy of 94.86%, which is certainly benefited from using multiple forward passes and multiple losses, e.g. rotation loss, thus takes much more computing to train.

Table 2: Results (%) between DMT and CL on CIFAR-10 test set.

	Baseline	CL [26]	DMT
4000 labels	86.08	94.02	94.21 (+0.19)
1000 labels	75.14	90.61	$91.51 \; (\mathbf{+0.90})$

Table 3: Results (%) for DMT and other methods on CIFAR-10 test set using 4,000 labels. Oracle performance is 96.35%.

method	Baseline	MT [9]	DCT [15]	DS [35]	MixMatch [11]	DAG [36]	CL [26]	Ours (DMT)
network	WRN-28-2	WRN-28-2	CNN-13	CNN-13	WRN-28-2	CNN-13	WRN-28-2	WRN-28-2
accuracy	86.08	89.64	90.97	91.11	93.76	93.87	94.02	94.21

5.3.2. PASCAL VOC 2012

PASCAL VOC 2012 is the most commonly used benchmark for semi-supervised semantic segmentation. We compare our method with consistency-based Mean Teacher directly adapted to semantic segmentation (MT-Seg), Mean Teacher with strong CutMix augmentation [10], feature-level consistency-based method Cross-Consistency Training (CCT) [37], adapted Dual Student for semantic segmentation with an auxiliary flaw detector (GCT) [38], GAN-based method [8] that pre-trains a discriminator to select pseudo labels, and hybrid method s4GAN + MLMT [12] that adds consistency regularization upon [8] by an extra classification branch. Methods are evaluated on 4 challenging splits: 1/106 (100 labels), 1/50, 1/20 and 1/8. We do not use more than 1/8 data which is becoming easier and pose less challenge to state-of-the-art methods. Supervised performance on the labeled subset is reported as Baseline. MT-Seg, CCT and GCT performance are the re-evaluated results in the GCT codebase²: others are taken from the original papers. All methods use the same network architecture as ours except for CCT in which a slightly superior architecture PSPNet-ResNet-101 [39] is used for evaluation.

As shown in Tab. 4, DMT outperforms other methods with a clear margin.

²https://github.com/ZHKKKe/PixelSSL/tree/master/task/sseg

Table 4: Mean IoU (%) results for DMT and other methods on PASCAL VOC 2012 val set. Performance gap to Oracle is shown in brackets. † $Updated\ numbers\ from\ s4GAN\ +\ MLMT$. * $ImageNet\ pre-training$.

method	network	1/106	1/50	1/20	1/8	Oracle
Baseline	DeepLab-v2	46.66 (-28.09)	55.62 (-19.13)	62.29 (-12.46)	67.37 (-7.38)	74.75
MT-Seg [9]	DeepLab-v2	-	-	-	67.65 (-5.94)	73.59
Hung et al. [8]	DeepLab-v2	-	$57.2^{\dagger} \ (-17.7)$	$64.7^{\dagger} \ (-10.2)$	69.5 (-5.4)	74.9
${\rm s4GAN+MLMT[12]}$	DeepLab-v2	-	63.3 (-12.3)	67.2 (-8.4)	71.4 (-4.2)	75.6
CutMix [10]*	DeepLab-v2	53.79 (-18.75)	64.81 (-7.73)	66.48 (-6.06)	67.60 (-4.94)	72.54
CCT [37]	PSPNet	-	-	-	70.45 (-4.80)	75.25
GCT [38]	DeepLab-v2	-	-	-	70.57 (-3.49)	74.06
Ours (DMT)	DeepLab-v2	63.04 (-11.71)	67.15 (-7.60)	69.92 (-4.83)	$72.70 \ (-2.05)$	74.75

However, some methods use GAN and extra network branch, or have implementation flaws, resulting in different Oracle and Baseline performances. Thus, we further show performance gaps to Oracle in brackets for fair comparisons, where DMT's performance is also the closest to Oracle. Our proposed DMT is the only method showing a clear improvement over Baseline on the challenging 100 labels split other than CutMix (1/106 in Tab. 4). In addition, DMT shows more stable performance across different labeled ratios (Fig. 8).

Comparing with human supervision. We design an interesting experiment to show how DMT is even superior to human annotators. Specifically, the original PASCAL VOC 2012 dataset only labels 1,464 training images, called the *train* set. While the commonly used 10,582 training set *trainaug* contains 9,118 images from SBD [30]. The SBD dataset uses the same set of images as PASCAL VOC and annotates object outlines by Amazon Mechanical Turk (AMT), which can be filled as segmentation masks. However, unprofessional annotators from AMT tend to draw coarse outlines. Thus, *trainaug* has worse label quality than *train*. Therefore, we use *train* as the labeled subset and the 9,118 images from SBD as the unlabeled subset (by removing the SBD labels), and experiment DMT with the same hyper-parameters in the 1/8 split experiment. Surprisingly, as shown in Tab. 5, DMT exceeds the performance of Oracle (fully-supervised training on *trainaug*). This suggests that DMT renders human

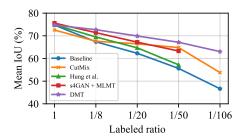


Figure 8: Comparison of semi-supervised semantic segmentation performance on PASCAL VOC 2012 with multiple data splits. The proposed DMT has more stable performance across different labeled ratios.

Table 5: Mean IoU (%) comparisons between Oracle and DMT on PASCAL VOC 2012 1464/9118 split. val mean IoU (%) reported.

	number of images	number of labels	mean IoU
Baseline	1464	1464	72.10 ± 0.53
DMT	10582	1464	74.85 ± 0.29
Oracle	10582	10582	74.75 ± 0.25

supervision at this quality (AMT) unnecessary for semantic segmentation, except for somewhat faster training since DMT needs two models (roughly twice the training budget of one model).

Qualitative results. We provide qualitative comparisons in segmentation results among Baseline, DMT and ground truth. As shown in Fig. 9, there is confusion between similar classes in Baseline predictions (column 3), such as horse and cow (row 1, 2), dog and bird (row 3), train and motorbike (row 4). After dynamic mutual training (column 4), class confusion is mostly resolved. Also finer details are recovered (e.g. distant people in row 5 and chairs in row 6). Moreover, Baseline entirely fails to detect the dining table in row 6.

5.3.3. Cityscapes

Cityscapes features complex street scenes which are less ventured by semisupervised learning methods. Of the six methods evaluated in Section 5.3.2, three methods have chosen to report performance on this dataset. We follow

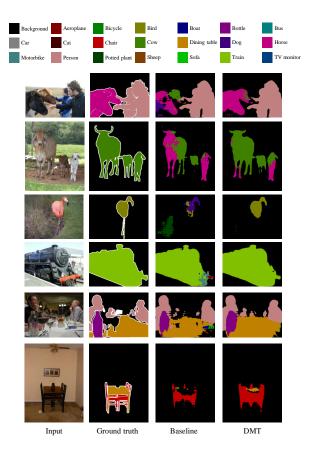


Figure 9: Qualitative results on PASCAL VOC 2012. Models are trained with 1/20 pixel-level labeled data in the training set. White regions are ignored in ground truth.

the same evaluation protocols as PASCAL VOC 2012, to evaluate each method on 1/30 (100 labels) and 1/8 splits. All reported numbers are taken from the original papers.

As shown in Tab. 6, Hung et al. [8] is merely comparable to our fully-tuned Baselines. Although stronger methods s4GAN + MLMT [12] and CutMix [10] still obtain good results, DMT outperforms them by $2 \sim 3\%$. Refer to supplementary materials for proper baseline training and qualitative results.

Table 6: Mean IoU (%) results for DMT and other methods on Cityscapes val set. Performance gap to Oracle is shown in brackets. * ImageNet pre-training.

method	1/30	1/8	Oracle
Baseline	49.54 (-18.62)	59.65 (-8.51)	68.16
Hung et al. [8]	-	58.8 (-8.9)	67.7
$s4GAN + MLMT \ [12]*$	-	59.3 (-6.5)	65.8
CutMix[10]*	51.20 (-16.48)	60.34 (-7.34)	67.68
Ours (DMT)	54.80 (-13.36)	63.03 (-5.13)	68.16

5.4. Ablations

To further validate our method choice and provide additional insights for online and offline self-training, we conduct the following ablations.

- Online ST. Instead of 5 iterations of DMT, we perform online self-training with fixed confidence threshold 0.9 for 20 epochs.
- CBST. The CBST algorithm [13] is modified to direct ranking, i.e. iterative class-balanced self-training without dynamic loss, similar to a class-balanced version of CL [26].
- **DST.** Same as DMT, except using only one model to provide pseudo labels for itself, i.e. DMT with only one model F_A fine-tuning itself.
- **DMT-Naive.** We directly re-weight the loss by confidence without distinguishing the three cases in Eq. 3.
- **DMT-Flip.** In the third condition of Eq. 3, since the pseudo label is likely incorrect, instead of setting loss to 0, we flip the pseudo label to the current model's prediction and weight the loss by $(1 c_A)^{\gamma_2}$, acting as an estimate of disagreement between models, given the pseudo label is flipped.

To show clear differences between setups, the ablations are carried out on PASCAL VOC 2012 in Tab. 7. This dataset is sufficiently complex and experiments run faster due to fine-tuning.

Table 7: Ablations on PASCAL VOC 2012 (one random 1/20 split and one random 1/50 split). val mean IoU (%) is reported.

ablations	Baseline	Online ST	CBST	DST	DMT-Naive	DMT	DMT-Flip
1/20	61.90	63.12	65.09	69.43	70.00	70.16	70.17
1/50	56.29	53.52	62.29	66.50	64.95	68.37	68.35

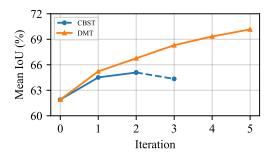


Figure 10: Comparison of DMT and CBST (both COCO initialized) on one PASCAL VOC $2012\ 1/20$ split. Iteration 0 represents training on the labeled subset only. CBST performance starts to degrade at iteration three due to too much pseudo label noise.

Online ST has very limited performance due to continuously fitting on selferrors, where it even performs worse than Baseline when labels are extremely scarce (1/50). We also observe that its performance remains similar without data augmentation on the unlabeled subset, echos our analysis that consistency regularization cannot be integrated for online self-training without multiple forward passes (Section 2.2).

CBST is an offline method, thus it is less sensitive to self-error than Online ST and consistently improves over Baseline by a clear margin. However, it only conducts self-training without considering pseudo label noise. In our experiments, performance increase stops at iteration three for CBST while DMT benefits from all 5 iterations. Because pseudo label noise prevents further improvements on CBST when using more unlabeled data.

DST requires half the computing budget compared to DMT and performs well. As fine-tuning goes on in each iteration, a relatively large learning rate and

data augmentations drive the model to deviate from its previous-self³ rapidly. Thus, sufficient inter-model disagreement is provided for dynamic weighting to take effects. While in DMT, larger inter-model disagreement by starting from different model initializations naturally enables better final results, especially on the more challenging 1/50 split.

DMT-Naive is a simpler formulation to integrate inter-model disagreement, the performance is even comparable to DMT on the 1/20 split. But its performance degrades significantly when label noise is severe (1/50 split). Although this naive policy outperforms CBST on this task, its performance is similar or worse than CL on CIFAR-10, indicates poor generalization ability.

DMT-Flip is more complex than DMT. But its performance is similar to DMT, differences are at the level of random variations. We suspect by flipping labels, there is a similar drawback as online self-training: fitting newly made self-errors (Section 3). Thus the results are no better than ignoring those labels (usually fewer than 5% in training). Moreover, considering simplifying the policy also brings notable degradation (DMT-Naive), the three-cases setup in DMT (Eq. 3) offers a reasonable trade-off between performance and complexity.

6. Related Work

6.1. Semi-Supervised Learning

Since popular semi-supervised learning methods are described and compared in Section 2 and 5.3. We focus on methods that are most relevant to DMT. Among typical deep learning methods that use more than one models to explicitly/implicitly exploit the inter-model disagreement, the most similar to DMT are Dual Student [35] and Deep Co-training [15]. In Dual Student, two models are trained in parallel online to select stable examples for each other to learn. However, the stable examples are decided by one model alone and the other

 $^{^3}$ Previous-self denotes the model state that produced the pseudo labels before the current mutual training iteration.

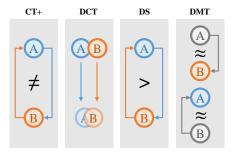


Figure 11: Comparison of inter-model disagreement usages among Co-Teaching+ (CT+) [40], Deep Co-Training (DCT) [15], Dual Student (DS) [35] and DMT. Left: CT+ let models decide training examples for each other, but it only updates when models disagree (\neq). Middle-Left: DCT maximizes the agreement between two models (A, B). Middle-Right: DS let one model teach another when it is more certain on an example (>). Right: DMT updates one model with fixed (grey) pseudo labels from another, depending on how much they disagree (\approx).

model cannot dispute that supervision. In Deep Co-training, two models are also trained in parallel online, and the learning objective is to minimize their disagreement on the unlabeled subset. Since disagreement is minimized, the models can rapidly converge to the same set of weights online, thus explicit weight distance constrains have to be employed to avoid collapse. Largely different from them, the dynamic loss in DMT is determined by both models and does not require special constrains to avoid collapse. We illustrate the major differences between DMT and these two methods in Fig. 11.

Other than methodological differences to previous methods, DMT is also generally applicable to both image classification and semantic segmentation tasks. While most semi-supervised learning methods are ad-hoc and only work well in a limited range of tasks (e.g. only work on either image classification [9] or pixel-wise task [38]). To the best of our knowledge, no previous methods have been shown to reach state-of-the-art performance on both image classification and semantic segmentation without additional efforts. For instance, consistency-based methods such as Mean Teacher [9] works well in image classification but performs only comparable to Baseline on PASCAL VOC 2012

semantic segmentation (Tab. 4). High-dimensional perturbations have to be imposed to work in semantic segmentation, e.g. CutMix [10] and CCT [37].

6.2. Learning with Noisy Labels

Learning with noisy labels is a well-studied topic. Most researches on this topic tackle random noise that can be modeled by a noise transition matrix, with each matrix entry as the label random flip probability from one class to another [21, 23]. We focus on general methods that have no explicit noise source modeling. Decoupling [24] trains two models simultaneously online, and only performs gradient descent when two models disagree, to decouple "when" and "how" to update model parameters, i.e. not allowing the noisy labels to control when to learn. Co-teaching [41] also trains two models online, while each model selects low loss examples for the other to train, and a similar policy has been exploited in semi-supervised learning by Dual Student [35]. Recently, Co-teaching+ [40] combines the idea of Decoupling and Co-teaching. However, these methods still mostly show good performance regarding only random noise.

Contrary to how disagreement is leveraged in Decoupling/Co-teaching+ [24, 40], in semi-supervised learning we deal with pseudo label noise made by deep nets similar to the model in training, where the disagreement between models signifies possible errors instead of an parameter update opportunity. Besides, Co-teaching [41] uses two models to select examples for each other, and the training targets are still determined by one model only. While in DMT the two models collaborate explicitly to re-weigh loss. Major differences between Co-teaching+ and DMT are illustrated in Fig. 11.

There is one approach that deals particularly with pseudo labels [36]. Specifically, given a feature embedding graph with every data point (image) as a node, a node's pseudo label can thus be rectified by its neighbors. It is feasible for relatively small datasets such as CIFAR-10, but rather unrealistic for segmentation datasets, where one image has millions of data points (pixels). Besides, the graph-based method is applied at the pseudo-labeling phase, making it complementary to DMT, which is adopted at training phase.

7. Conclusions and Discussions

In this paper, we have proposed Dynamic Mutual Training (DMT) to counter the pseudo supervision noise by a re-weighted loss function based on the intermodel disagreement. Furthermore, we have adapted DMT to an iterative framework for better performance in both image classification and semantic segmentation. DMT is flexible and easy to implement. We have evaluated the proposed method on different datasets including CIFAR-10, PASCAL VOC 2012 and Cityscapes. The experiments (comparisons and ablations) clearly demonstrate the effectiveness of the proposed DMT, and show its state-of-the-art outcomes in classification and segmentation.

We find that DMT is more promising in semantic segmentation than image classification, probably because dynamic weighting exploits pixels with high-quality pseudo labels in an image and provides better pseudo supervision on each image overall. In addition, image classification on CIFAR-10 requires retraining for each iteration and the two models do not have equal classification ability throughout training, thus making it difficult to exploit inter-model disagreement. Besides, it is hard to estimate confidence when recent image classification models require heavy data augmentation in training, confidence distribution is quite different compared to the generated pseudo labels. Thus, a better confidence estimation process could bring further gains, e.g. multiple forward pass statistics (at the cost of computing). Also advances in the learning with noisy labels community may be potentially useful, which is worthy of future investigations.

Our work shares the common limitation of most offline self-training methods: the initial model learned on the labeled subset may be insufficient if the labels are too few, e.g. 100 labels on Cityscapes. Better pre-trained weights, e.g. COCO pre-trained weights for PASCAL VOC 2012, can potentially alleviate this issue, given a small gap to Oracle on PASCAL VOC 2012 (11.71%, Tab. 4). If off-the-shelf pre-trained weights are unavailable, self-supervised learning [42] is a good way to initiate learning. We intend to investigate how self-supervised learning

can help semi-supervised learning in the future, especially for structured tasks like semantic segmentation.

8. Acknowledgements

This work was partially supported by National Key Research and Development Program of China (No. 2019YFC1521104), Art major project of National Social Science Fund (I8ZD22). The author Qianyu Zhou is supported by Wu Wenjun Honorary Doctoral Scholarship, AI Institute, Shanghai Jiao Tong University. Also, the authors would like to thank Shikun Liu (Imperial College London) for insightful discussions and identifying an incorrect data augmentation policy in the earlier semantic segmentation code.

Appendix A. Extra Qualitative Results

Qualitative comparisons among Baseline (supervised only), our method (DMT) and ground truth on Cityscapes are shown in Fig. A.12. Specifically, there is confusion among similar classes from Baseline (column 3), such as bus and car (row 1), wall, fence and building (row 2), road and sidewalk (row 3). In contrast, DMT does a better job (column 4). DMT also detects the existence of small objects better, e.g. motorcycle in row 4, pole in row 5. While Baseline claims existence of non-existent class sky in row 6. In row 7, where the real-world fence is more complex, Baseline produces chaotic results which become consistent after dynamic mutual training.

Appendix B. Proper Baseline Training

A concern of non-fully training baselines in semi-supervised learning has been raised in [29]. We find that other than unifying data augmentation schemes and tricks (e.g. using the same strong augmentations and mixup in CIFAR-10 baselines), one important factor is the number of epochs. For example, we have a labeled subset that is 1/8 the entire dataset, 8 times more epochs (i.e. keep the number of steps unchanged) are too many, and the same number of epochs are too few. Thus, we make a compromise between them and train for $\sqrt{\frac{1}{labeled\ ratio}}N$ epochs, where N is the number of epochs used in Oracle training, which is 300, 30, 60 for CIFAR-10, PASCAL VOC 2012 and Cityscapes, respectively. As a result, the reported baseline performance in this paper is noticeably higher than previous works, sometimes even comparable to some previous state-of-the-art methods on Cityscapes (Tab. 5).

In our fully-supervised baselines, the learning rate is set to 0.2 (CIFAR-10), 2×10^{-3} (PASCAL VOC 2012), and 4×10^{-3} (Cityscapes).

Appendix C. Realistic Validation

Oliver et al. [29] has raised a concern of using unrealistically large validation sets in semi-supervised image classification. For instance, most prior arts split

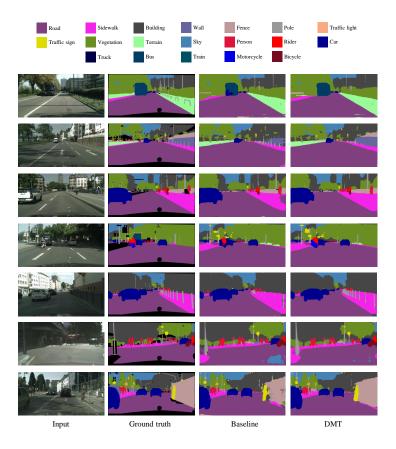


Figure A.12: Qualitative results on Cityscapes. Models are trained with 1/8 pixel-level labeled data in the training set. Black regions are ignored in ground truth.

a validation set of 5,000 images for hyper-parameter tuning on CIFAR-10, even larger than their labeled subset. While we use only 200 images (we call this validation set *valtiny* and will release it along with our source codes). Intuitively, 200 images can hardly tell the difference between two models with 100%/200 = 0.5% accuracy gap. Different from semantic segmentation, since the mean IoU is a fine-grained metric that can work with a small validation set.

Fine-grained testing. To test with fewer images, we propose fine-grained testing for image classification, where we count the probability of being correct. Concretely, if a model makes the right decision on an image, we count it as the probability that model assigns for the correct class instead of 1. In this way,

we assess not only whether a prediction is correct, but also how correct it is. We observe it helpful when normal testing can not tell the difference, especially when comparing similar setups (a few different hyper-parameter values). But fine-grained testing does not work well when comparing mixup methods and non-mixup methods, since mixup is better class-calibrated [43].

References

- [1] S. Zagoruyko, N. Komodakis, Wide residual networks, in: British Machine Vision Conference, 2016, pp. 1–1.
- [2] Z. Wu, C. Shen, A. Van Den Hengel, Wider or deeper: Revisiting the resnet model for visual recognition, Pattern Recognition 90 (2019) 119–133.
- [3] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs, Pattern Analysis and Machine Intelligence 40 (2017) 834–848.
- [4] J. Fu, J. Liu, Y. Li, Y. Bao, W. Yan, Z. Fang, H. Lu, Contextual deconvolution network for semantic segmentation, Pattern Recognition 101 (2020) 107152.
- [5] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, B. Schiele, The cityscapes dataset for semantic urban scene understanding, in: Computer Vision and Pattern Recognition, 2016, pp. 3213–3223.
- [6] D. Yarowsky, Unsupervised word sense disambiguation rivaling supervised methods, in: Annual Meeting of the Association for Computational Linguistics, 1995, pp. 189–196.
- [7] D.-H. Lee, Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks, in: International Conference on Machine Learning workshop, 2013, pp. 1–1.

- [8] W. Hung, Y. Tsai, Y. Liou, Y. Lin, M. Yang, Adversarial learning for semisupervised semantic segmentation, in: British Machine Vision Conference, 2018, p. 65.
- [9] A. Tarvainen, H. Valpola, Mean teachers are better role models: Weightaveraged consistency targets improve semi-supervised deep learning results, in: Neural Information Processing Systems, 2017, pp. 1195–1204.
- [10] G. French, S. Laine, T. Aila, M. Mackiewicz, G. Finlayson, Semi-supervised semantic segmentation needs strong, varied perturbations, in: British Machine Vision Conference, 2020, pp. 1–1.
- [11] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, C. A. Raffel, Mixmatch: A holistic approach to semi-supervised learning, in: Neural Information Processing Systems, 2019, pp. 5050–5060.
- [12] S. Mittal, M. Tatarchenko, T. Brox, Semi-supervised semantic segmentation with high- and low-level consistency, Pattern Analysis and Machine Intelligence (2019) 1–1.
- [13] Y. Zou, Z. Yu, B. Vijaya Kumar, J. Wang, Unsupervised domain adaptation for semantic segmentation via class-balanced self-training, in: European Conference on Computer Vision, 2018, pp. 289–305.
- [14] Y. Bengio, J. Louradour, R. Collobert, J. Weston, Curriculum learning, in: International Conference on Machine Learning, 2009, pp. 41–48.
- [15] S. Qiao, W. Shen, Z. Zhang, B. Wang, A. Yuille, Deep co-training for semi-supervised image recognition, in: European Conference on Computer Vision, 2018, pp. 135–152.
- [16] J. Peng, G. Estrada, M. Pedersoli, C. Desrosiers, Deep co-training for semisupervised image segmentation, Pattern Recognition 107 (2020) 107269.
- [17] Z.-H. Zhou, M. Li, Semi-supervised learning by disagreement, Knowledge and Information Systems 24 (2010) 415–439.

- [18] Y. Grandvalet, Y. Bengio, Semi-supervised learning by entropy minimization, in: Neural Information Processing Systems, 2005, pp. 529–536.
- [19] D. Berthelot, N. Carlini, E. D. Cubuk, A. Kurakin, K. Sohn, H. Zhang, C. Raffel, Remixmatch: Semi-supervised learning with distribution alignment and augmentation anchoring, in: International Conference on Learning Representations, 2020, pp. 1–1.
- [20] S. Laine, T. Aila, Temporal ensembling for semi-supervised learning, in: International Conference on Learning Representations, 2017, pp. 1–1.
- [21] D. Angluin, P. Laird, Learning from noisy examples, Machine Learning 2 (1988) 343–370.
- [22] C. Zhang, S. Bengio, M. Hardt, B. Recht, O. Vinyals, Understanding deep learning requires rethinking generalization, in: International Conference on Learning Representations, 2017, pp. 1–1.
- [23] J. Goldberger, E. Ben-Reuven, Training deep neural-networks using a noise adaptation layer, in: International Conference on Learning Representations, 2017, pp. 1–1.
- [24] E. Malach, S. Shalev-Shwartz, Decoupling "when to update" from "how to update", in: Neural Information Processing Systems, 2017, pp. 960–970.
- [25] Y. Zou, Z. Yu, X. Liu, B. Kumar, J. Wang, Confidence regularized self-training, in: International Conference on Computer Vision, 2019, pp. 5982–5991.
- [26] P. Cascante-Bonilla, F. Tan, Y. Qi, V. Ordonez, Curriculum labeling: Revisiting pseudo-labeling for semi-supervised learning, in: AAAI Conference on Artificial Intelligence, 2021, pp. 1–1.
- [27] A. Krizhevsky, Learning multiple layers of features from tiny images, Technical Report, University of Toronto, 2009.

- [28] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, A. Zisserman, The pascal visual object classes challenge: A retrospective, International Journal of Computer Vision 111 (2015) 98–136.
- [29] A. Oliver, A. Odena, C. A. Raffel, E. D. Cubuk, I. Goodfellow, Realistic evaluation of deep semi-supervised learning algorithms, in: Neural Information Processing Systems, 2018, pp. 3235–3246.
- [30] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, J. Malik, Semantic contours from inverse detectors, in: International Conference on Computer Vision, 2011, pp. 991–998.
- [31] P. Micikevicius, S. Narang, J. Alben, G. F. Diamos, E. Elsen, D. García, B. Ginsburg, M. Houston, O. Kuchaiev, G. Venkatesh, H. Wu, Mixed precision training, in: International Conference on Learning Representations, 2018, pp. 1–1.
- [32] P. Izmailov, D. Podoprikhin, T. Garipov, D. Vetrov, A. G. Wilson, Averaging weights leads to wider optima and better generalization, in: Conference on Uncertainty in Artificial Intelligence, 2018, pp. 1–1.
- [33] E. D. Cubuk, B. Zoph, J. Shlens, Q. V. Le, Randaugment: Practical automated data augmentation with a reduced search space, in: Computer Vision and Pattern Recognition workshop, 2020, pp. 702–703.
- [34] H. Zhang, M. Cisse, Y. N. Dauphin, D. Lopez-Paz, mixup: Beyond empirical risk minimization, in: International Conference on Learning Representations, 2018, pp. 1–1.
- [35] Z. Ke, D. Wang, Q. Yan, J. Ren, R. W. Lau, Dual student: Breaking the limits of the teacher in semi-supervised learning, in: International Conference on Computer Vision, 2019, pp. 6728–6736.
- [36] S. Li, B. Liu, D. Chen, Q. Chu, L. Yuan, N. Yu, Density-aware graph for deep semi-supervised visual recognition, in: Computer Vision and Pattern Recognition, 2020, pp. 13400–13409.

- [37] Y. Ouali, C. Hudelot, M. Tami, Semi-supervised semantic segmentation with cross-consistency training, in: Computer Vision and Pattern Recognition, 2020, pp. 1–1.
- [38] Z. Ke, D. Qiu, K. Li, Q. Yan, R. W. Lau, Guided collaborative training for pixel-wise semi-supervised learning, in: European Conference on Computer Vision, volume 2, 2020, p. 6.
- [39] H. Zhao, J. Shi, X. Qi, X. Wang, J. Jia, Pyramid scene parsing network, in: Computer Vision and Pattern Recognition, 2017, pp. 2881–2890.
- [40] X. Yu, B. Han, J. Yao, G. Niu, I. W. Tsang, M. Sugiyama, How does disagreement help generalization against label corruption?, in: International Conference on Machine Learning, 2019, pp. 1–1.
- [41] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. Tsang, M. Sugiyama, Co-teaching: Robust training of deep neural networks with extremely noisy labels, in: Neural Information Processing Systems, 2018, pp. 8527–8537.
- [42] X. Zhai, A. Oliver, A. Kolesnikov, L. Beyer, S4l: Self-supervised semisupervised learning, in: International Conference on Computer Vision, 2019, pp. 1476–1485.
- [43] S. Thulasidasan, G. Chennupati, J. A. Bilmes, T. Bhattacharya, S. Michalak, On mixup training: Improved calibration and predictive uncertainty for deep neural networks, in: Neural Information Processing Systems, 2019, pp. 13888–13899.