

# Semi-supervised Domain Adaptation via Minimax Entropy

Kuniaki Saito<sup>1</sup>, Donghyun Kim<sup>1</sup>, Stan Sclaroff<sup>1</sup>, Trevor Darrell<sup>2</sup> and Kate Saenko<sup>1</sup>

<sup>1</sup>Boston University, <sup>2</sup>University of California, Berkeley

<sup>1</sup>{keisaito, donhk, sclaroff, saenko}@bu.edu, <sup>2</sup>trevor@eecs.berkeley.edu

## Abstract

Contemporary domain adaptation methods are very effective at aligning feature distributions of source and target domains without any target supervision. However, we show that these techniques perform poorly when even a few labeled examples are available in the target domain. To address this semi-supervised domain adaptation (SSDA) setting, we propose a novel Minimax Entropy (MME) approach that adversarially optimizes an adaptive few-shot model. Our base model consists of a feature encoding network, followed by a classification layer that computes the features’ similarity to estimated prototypes (representatives of each class). Adaptation is achieved by alternately maximizing the conditional entropy of unlabeled target data with respect to the classifier and minimizing it with respect to the feature encoder. We empirically demonstrate the superiority of our method over many baselines, including conventional feature alignment and few-shot methods, setting a new state of the art for SSDA.

## 1. Introduction

Deep convolutional neural networks [15] have significantly improved image classification accuracy with the help of large quantities of labeled training data, but often generalize poorly to new domains. Recent unsupervised domain adaptation (UDA) methods [10, 18, 19, 27, 35] improve generalization on unlabeled target data by aligning distributions, but can fail to learn discriminative class boundaries on target domains (see Fig. 1.) We show that in the Semi-Supervised Domain Adaptation (SSDA) setting where a few target labels are available, such methods often do not improve performance relative to just training on labeled source and target examples, and can even make it worse.

We propose a novel approach for SSDA that overcomes the limitations of previous methods and significantly improves the accuracy of deep classifiers on novel domains with only a few labels per class. Our approach, which we call Minimax Entropy (MME), is based on optimizing a minimax loss on the conditional entropy of unlabeled data,

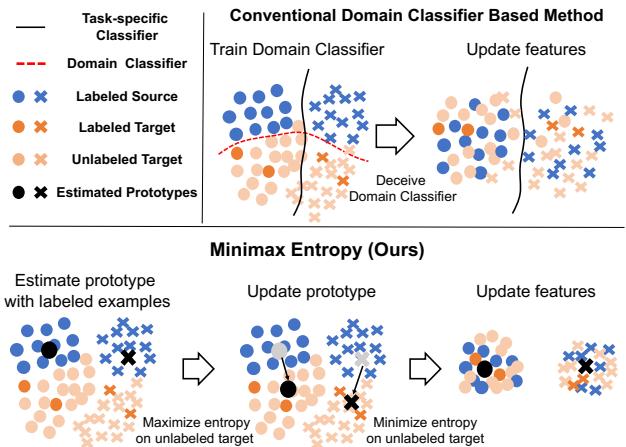


Figure 1: We address the task of semi-supervised domain adaptation. Top: Existing domain-classifier based methods align source and target distributions but can fail by generating ambiguous features near the task decision boundary. Bottom: Our method estimates a representative point of each class (prototype) and extracts discriminative features using a novel minimax entropy technique.

as well as the task loss; this reduces the distribution gap while learning discriminative features for the task.

We exploit a cosine similarity-based classifier architecture recently proposed for few-shot learning [11, 4]. The classifier (top layer) predicts a K-way class probability vector by computing cosine similarity between K class-specific weight vectors and the output of a feature extractor (lower layers), followed by a softmax. Each class weight vector is an estimated “prototype” that can be regarded as a representative point of that class. While this approach outperformed more advanced methods in few-shot learning and we confirmed its effectiveness in our setting, as we show below it is still quite limited. In particular, it does not leverage unlabeled data in the target domain.

Our key idea is to minimize the distance between the class prototypes and neighboring unlabeled target samples, thereby extracting discriminative features. The problem is how to estimate domain-invariant prototypes without many labeled target examples. The prototypes are dominated by the source domain, as shown in the leftmost side of Fig. 2

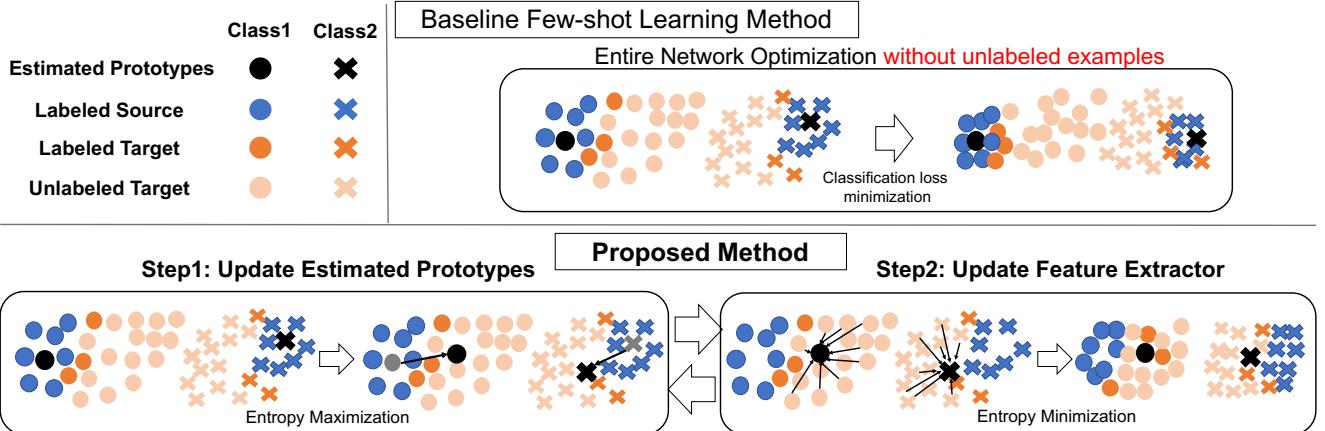


Figure 2: Top: baseline few-shot learning method, which estimates class prototypes by weight vectors, yet does not consider unlabeled data. Bottom: our model extracts discriminative and domain-invariant features using unlabeled data through a domain-invariant prototype estimation. Step 1: we update the estimated prototypes in the classifier to maximize the entropy on the unlabeled target domain. Step 2: we minimize the entropy with respect to the feature extractor to cluster features around the estimated prototype.

(bottom), as the vast majority of labeled examples come from the source. To estimate domain-invariant prototypes, we move weight vectors toward the target feature distribution. Entropy on target examples represents the similarity between the estimated prototypes and target features. A uniform output distribution with high entropy indicates that the examples are similar to all prototype weight vectors. Therefore, we move the weight vectors towards target by maximizing the entropy of unlabeled target examples in the first adversarial step. Second, we update the feature extractor to minimize the entropy of the unlabeled examples, to make them better clustered around the prototypes. This process is formulated as a mini-max game between the weight vectors and the feature extractor and applied over the unlabeled target examples.

Our method offers a new state-of-the-art in performance on SSDA; as reported below, we reduce the error relative to baseline few-shot methods which ignore unlabeled data by 8.5%, relative to current best-performing alignment methods by 8.8%, and relative to a simple model jointly trained on source and target by 11.3% in one adaptation scenario. Our contributions are summarized as follows:

- We highlight the limitations of state-of-the-art domain adaptation methods in the SSDA setting;
- We propose a novel adversarial method, Minimax Entropy (MME), designed for the SSDA task;
- We show our method’s superiority to existing methods on benchmark datasets for domain adaptation.

## 2. Related Work

**Domain Adaptation.** Semi-supervised domain adaptation (SSDA) is a very important task [7, 38], however it has not been fully explored, especially with regard to deep learning based methods. We revisit this task and compare our approach to recent semi-supervised learning or unsupervised domain adaptation methods. The main challenge in domain adaptation (DA) is the gap in feature distributions between domains, which degrades the source classifier’s performance. Most recent work has focused on unsupervised domain adaptation (UDA) and, in particular, feature distribution alignment. The basic approach measures the distance between feature distributions in source and target, then trains a model to minimize this distance. Many UDA methods utilize a domain classifier to measure the distance [10, 35, 18, 19]. The domain-classifier is trained to discriminate whether input features come from the source or target, whereas the feature extractor is trained to deceive the domain classifier to match feature distributions. UDA has been applied to various applications such as image classification [26], semantic segmentation [31], and object detection [5, 28]. Some methods minimize task-specific decision boundaries’ disagreement on target examples [29, 27] to push target features far from decision boundaries. In this respect, they increase between-class variance of target features; on the other hand, we propose to make target features well-clustered around estimated prototypes. Our MME approach can reduce within-class variance as well as increasing between-class variance, which results in more discriminative features. Interestingly, we empirically observe that UDA methods [10, 19, 27] often fail in improving accuracy in SSDA.

**Semi-supervised learning (SSL).** Generative [6, 30],

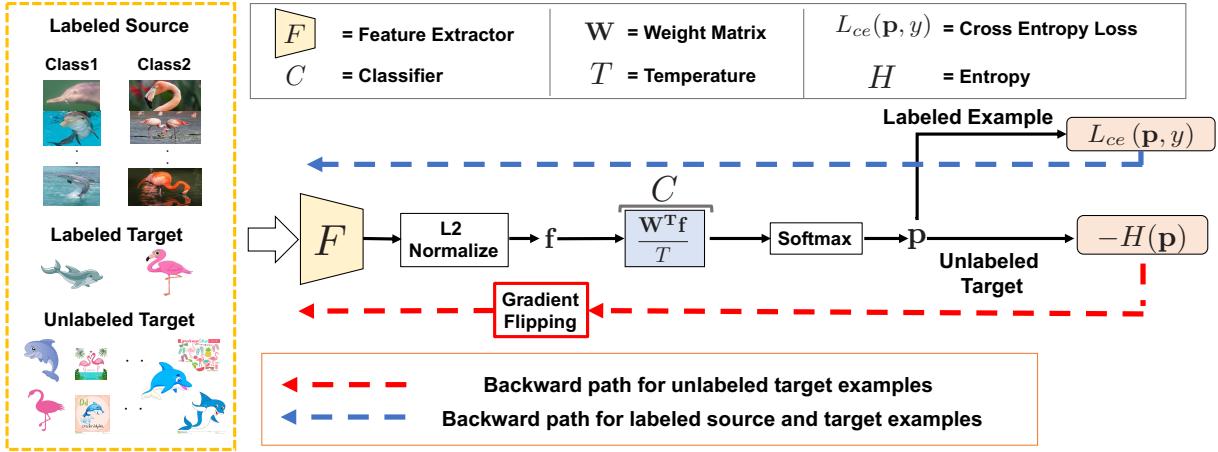


Figure 3: An overview of the model architecture and MME. The inputs to network are labeled source examples, a few labeled target examples, and unlabeled target examples. Our model consists of the feature extractor  $F$  and the classifier  $C$  which has weight vectors ( $\mathbf{W}$ ) and temperature  $T$ .  $\mathbf{W}$  is trained to maximize entropy on unlabeled target (Step 1 in Fig. 2) whereas  $F$  is trained to minimize it (Step 2 in Fig. 2). To achieve the adversarial learning, the sign of gradients for entropy loss on unlabeled target examples is flipped by gradient reversal layer [10, 35].

model-ensemble [16], and adversarial approaches [21] have boosted performance in semi-supervised learning, but do not address domain shift. Conditional entropy minimization (CEM) is a widely used method in SSL [12, 9]. However, we found that CEM fails to improve performance when there is a large domain gap between the source and target domains (see experimental section.) MME can be regarded as a variant of entropy minimization which overcomes the limitation of CEM in domain adaptation.

**Few-shot learning (FSL).** Few shot learning [33, 37, 25] aims to learn novel classes given a few labeled examples and labeled “base” classes. SSDA and FSL make different assumptions: FSL does not use unlabeled examples and aims to acquire knowledge of *novel* classes, while SSDA aims to adapt to the *same* classes in a new domain. However both tasks aim to extract discriminative features given a few labeled examples from a novel domain or novel classes. We employ a network architecture with  $\ell_2$  normalization on features before the last linear layer and a temperature parameter  $T$ , which was proposed for face verification [24] and applied to few-shot learning [11, 4]. Generally, classification of a feature vector with a large norm results in confident output. To make the output more confident, networks can try to increase the norm of features. However, this does not necessarily increase the between-class variance because increasing the norm does not change the direction of vectors.  $\ell_2$  normalization on feature vectors can solve this issue. To make the output more confident, the network focuses on making the direction of the features from the same class closer to each other and separating different classes. This simple architecture was shown to be very effective for few-shot learning [4] and we build our method on it in our

work.

### 3. Minimax Entropy Domain Adaptation

In semi-supervised domain adaptation, we are given source images and the corresponding labels in the source domain  $\mathcal{D}_s = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^{m_s}$ . In the target domain, we are also given a limited number of labeled target images  $\mathcal{D}_t = \{(\mathbf{x}_i^t, y_i^t)\}_{i=1}^{m_t}$ , as well as unlabeled target images  $\mathcal{D}_u = \{(\mathbf{x}_i^u)\}_{i=1}^{m_u}$ . Our goal is to train the model on  $\mathcal{D}_s$ ,  $\mathcal{D}_t$ , and  $\mathcal{D}_u$  and evaluate on  $\mathcal{D}_u$ .

#### 3.1. Similarity based Network Architecture

Inspired by [4], our base model consists of a feature extractor  $F$  and a classifier  $C$ . For the feature extractor  $F$ , we employ a deep convolutional neural network and perform  $\ell_2$  normalization on the output of the network. Then, the normalized feature vector is used as an input to  $C$  which consists of weight vectors  $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_K]$  where  $K$  represents the number of classes and a temperature parameter  $T$ .  $C$  takes  $\frac{F(\mathbf{x})}{\|F(\mathbf{x})\|}$  as an input and outputs  $\frac{1}{T} \frac{\mathbf{W}^T F(\mathbf{x})}{\|F(\mathbf{x})\|}$ . The output of  $C$  is fed into a softmax-layer to obtain the probabilistic output  $\mathbf{p} \in R^n$ . We denote  $\mathbf{p}(\mathbf{x}) = \sigma(\frac{1}{T} \frac{\mathbf{W}^T F(\mathbf{x})}{\|F(\mathbf{x})\|})$ , where  $\sigma$  indicates a softmax function. In order to classify examples correctly, the direction of a weight vector has to be representative to the normalized features of the corresponding class. In this respect, the weight vectors can be regarded as estimated prototypes for each class. An architecture of our method is shown in Fig. 3.

### 3.2. Training Objectives

We estimate domain-invariant prototypes by performing entropy maximization with respect to the estimated prototype. Then, we extract discriminative features by performing entropy minimization with respect to feature extractor. In our method, the prototypes are parameterized by the weight vectors of the last linear layer. First, we train  $F$  and  $C$  to classify labeled source and target examples correctly and utilize an entropy minimization objective to extract discriminative features for the target domain. We use a standard cross-entropy loss to train  $F$  and  $C$  for classification:

$$\mathcal{L} = \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}_s, \mathcal{D}_t} \mathcal{L}_{ce}(\mathbf{p}(\mathbf{x}), y). \quad (1)$$

With this classification loss, we ensure that the feature extractor generates discriminative features with respect to the source and a few target labeled examples. However, the model is trained on the source domain and a small fraction of target examples for classification. This does not learn discriminative features for the entire target domain. Therefore, we propose minimax entropy training using unlabeled target examples.

A conceptual overview of our proposed adversarial learning is illustrated in Fig. 2. We assume that there exists a single domain-invariant prototype for each class, which can be a representative point for both domains. The estimated prototype will be near source distributions because source labels are dominant. Then, we propose to estimate the position of the prototype by moving each  $\mathbf{w}_i$  toward target features using unlabeled data in the target domain. To achieve this, we increase the entropy measured by the similarity between  $\mathbf{W}$  and unlabeled target features. Entropy is calculated as follows,

$$H = -\mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}_u} \sum_{i=1}^K p(y=i|\mathbf{x}) \log p(y=i|\mathbf{x}) \quad (2)$$

where  $K$  is the number of classes and  $p(y=i|\mathbf{x})$  represents the probability of prediction to class  $i$ , namely  $i$  th dimension of  $\mathbf{p}(\mathbf{x}) = \sigma(\frac{1}{T} \frac{\mathbf{W}^T F(\mathbf{x})}{\|F(\mathbf{x})\|})$ . To have higher entropy, that is, to have uniform output probability, each  $\mathbf{w}_i$  should be similar to all target features. Thus, increasing the entropy encourages the model to estimate the domain-invariant prototypes as shown in Fig. 2.

To obtain discriminative features on unlabeled target examples, we need to cluster unlabeled target features around the estimated prototypes. We propose to decrease the entropy on unlabeled target examples by the feature extractor  $F$ . The features should be assigned to one of the prototypes to decrease the entropy, resulting in the desired discriminative features. Repeating this prototype estimation (entropy maximization) and entropy minimization process yields discriminative features.

To summarize, our method can be formulated as adversarial learning between  $C$  and  $F$ . The task classifier  $C$  is

trained to maximize the entropy, whereas the feature extractor  $F$  is trained to minimize it. Both  $C$  and  $F$  are also trained to classify labeled examples correctly. The overall adversarial learning objective functions are:

$$\begin{aligned} \hat{\theta}_F &= \operatorname{argmin}_{\theta_F} \mathcal{L} + \lambda H \\ \hat{\theta}_C &= \operatorname{argmin}_{\theta_C} \mathcal{L} - \lambda H \end{aligned} \quad (3)$$

where  $\lambda$  is a hyper-parameter to control a trade-off between minimax entropy training and classification on labeled examples. Our method can be formulated as the iterative minimax training. To simplify training process, we use a gradient reversal layer [10] to flip the gradient between  $C$  and  $F$  with respect to  $H$ . With this layer, we can perform the minimax training with one forward and back-propagation, which is illustrated in Fig. 3.

### 3.3. Theoretical Insights

As shown in [1], we can measure domain-divergence by using a domain classifier. Let  $h \in \mathcal{H}$  be a hypothesis,  $\epsilon_s(h)$  and  $\epsilon_t(h)$  be the expected risk of source and target respectively, then  $\epsilon_t(h) \leq \epsilon_s(h) + d_{\mathcal{H}}(p, q) + C_0$  where  $C_0$  is a constant for the complexity of hypothesis space and the risk of an ideal hypothesis for both domains and  $d_{\mathcal{H}}(p, q)$  is the  $\mathcal{H}$ -divergence between  $p$  and  $q$ .

$$d_{\mathcal{H}}(p, q) \triangleq 2 \sup_{h \in \mathcal{H}} \left| \Pr_{\mathbf{x}^s \sim p} [h(\mathbf{f}^s) = 1] - \Pr_{\mathbf{x}^t \sim q} [h(\mathbf{f}^t) = 1] \right|. \quad (4)$$

where  $\mathbf{f}^s$  and  $\mathbf{f}^t$  denote the features in the source and target domain respectively. In our case the features are outputs of the feature extractor. The  $\mathcal{H}$ -divergence relies on the capacity of the hypothesis space  $\mathcal{H}$  to distinguish distributions  $p$  and  $q$ . This theory states that the divergence between domains can be measured by training a domain classifier and features with low divergence are the key to having a well-performing task-specific classifier. Inspired by this, many methods [10, 2, 35, 34] train a domain classifier to discriminate different domains while also optimizing the feature extractor to minimize the divergence.

Our proposed method is also connected to Eq. 4. Although we do not have a domain classifier or a domain classification loss, our method can be considered as minimizing domain-divergence through minimax training on unlabeled target examples. We choose  $h$  to be a classifier that decides a binary domain label of a feature by the value of the entropy, namely,

$$h(\mathbf{f}) = \begin{cases} 1, & \text{if } H(C(\mathbf{f})) \geq \gamma, \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

where  $C$  denotes our classifier,  $H$  denotes entropy, and  $\gamma$  is a threshold to determine a domain label. Here, we assume  $C$  outputs the probability of the class prediction for simplicity. Eq. 4 can be rewritten as follows,

$$\begin{aligned}
d_{\mathcal{H}}(p, q) &\triangleq 2 \sup_{h \in \mathcal{H}} \left| \Pr_{\mathbf{f}^s \sim p} [h(\mathbf{f}^s) = 1] - \Pr_{\mathbf{f}^t \sim q} [h(\mathbf{f}^t) = 1] \right| \\
&= 2 \sup_{C \in \mathcal{C}} \left| \Pr_{\mathbf{f}^s \sim p} [H(C(\mathbf{f}^s)) \geq \gamma] - \Pr_{\mathbf{f}^t \sim q} [H(C(\mathbf{f}^t)) \geq \gamma] \right| \\
&\leq 2 \sup_{C \in \mathcal{C}} \Pr_{\mathbf{f}^t \sim q} [H(C(\mathbf{f}^t)) \geq \gamma].
\end{aligned}$$

In the last inequality, we assume that  $\Pr_{\mathbf{f}^s \sim p} [H(C(\mathbf{f}^s)) \geq \gamma] \leq \Pr_{\mathbf{f}^t \sim p} [H(C(\mathbf{f}^t)) \geq \gamma]$ . This assumption should be realistic because we have access to many labeled source examples and train entire networks to minimize the classification loss. Minimizing the cross-entropy loss (Eq. 1) on source examples ensures that the entropy on a source example is very small. Intuitively, this inequality states that the divergence can be bounded by the ratio of target examples having entropy greater than  $\gamma$ . Therefore, we can have the upper bound by finding the  $C$  that achieves maximum entropy for all target features. Our objective is finding features that achieve lowest divergence. We suppose there exists a  $C$  that achieves the maximum in the inequality above, then the objective can be rewritten as,

$$\min_{\mathbf{f}^t} \max_{C \in \mathcal{C}} \Pr_{\mathbf{f}^t \sim q} [H(C(\mathbf{f}^t)) \geq \gamma] \quad (6)$$

Finding the minimum with respect to  $\mathbf{f}^t$  is equivalent to find a feature extractor  $F$  that achieves that minimum. Thus, we derive the minimax objective of our proposed learning method in Eq. 3. To sum up, our maximum entropy process can be regarded as measuring the divergence between domains, whereas our entropy minimization process can be regarded as minimizing the divergence. In our experimental section, we observe that our method actually reduces domain-divergence (Fig. 6c). In addition, target features produced by our method look aligned with source features and are just as discriminative. These come from the effect of the domain-divergence minimization.

## 4. Experiments

### 4.1. Setup

We randomly selected one or three labeled examples per one class as the labeled training target examples. We call each one-shot and three-shot setting respectively. We selected three other labeled examples as the validation for the target domain. The validation examples are used for early stopping, deciding hyper-parameter  $\lambda$ , and training scheduling. The other unlabeled target examples are used for training and evaluating classification accuracy (%). All examples of the source are used for training.

**Datasets.** LSDAC [23] is a benchmark dataset for large-scale domain adaptation. It has 345 classes and 6 domains. Since this dataset has various classes and many examples, we mainly used this dataset in the experiment. As labels of some domains and classes are very noisy, we pick 4 domains (Real, Clipart, Painting, Sketch) and 126 classes. We focus on the adaptation scenario where the target domain is

different from real and picked up 7 scenarios from the domains. See our supplemental material for the detail. **Office-Home** [36] contains 4 domains (Real, Clipart, Art, Product) with 65 classes. This dataset is one of the benchmark datasets for unsupervised domain adaptation. We evaluated our method on 12 scenarios in total. **Office** [26] contains 3 domains (Amazon, Webcam, DSLR) with 31 classes. Webcam and DSLR are small domains and some classes do not have a lot of examples while Amazon has many examples. To evaluate on the domain with enough examples, we have 2 scenarios where we set Amazon as the target domain and DSLR and Webcam as the source domain.

**Implementation Details.** All experiments are implemented in Pytorch [22]. We employ AlexNet [15] and VGG16 [32] in the experiments. To investigate the effect of deeper architectures, we use ResNet34 [13] in experiments on LSDAC. We remove the last linear layer of these networks to build  $F$ . we add a K-way linear classification layer  $C$  with randomly initialized weight matrix  $W$ . The value of temperature  $T$  is set 0.05 following the results of [24] in all settings. Every iteration, we prepared two mini-batches. One of them consists of labeled examples while the other one consists of unlabeled target examples. The half of labeled examples comes from source and the other half comes from labeled target. Using the two mini-batches, we calculated the objective in Eq. 3. To implement the adversarial learning in Eq. 3, we use a gradient reversal layer [10, 35] to flip the gradient with respect to entropy loss. The sign of the gradient is flipped between  $C$  and  $F$  during backpropagation. In all experiments, we set the trade-off parameter  $\lambda$  in Eq. 3 as 0.1. This is decided by the validation performance on Real to Clipart experiments. We will also show the performance sensitivity to this parameter in our supplemental material. We adopt SGD with the momentum of 0.9. Please see our supplemental material for more details including learning rate scheduling. We will publicize our implementation and dataset splits we used.

**Baselines.** **S+T** [4, 24] is a model trained with the labeled source and unlabeled target examples without using unlabeled target examples. **DANN** [10] employs a domain classifier to match feature distributions. This is the most popular method of UDA. For fair comparison, we modify this method so that it is trained with the labeled source, labeled target, and unlabeled target examples. **ADR** [27] utilizes a task-specific decision boundary to align features. This method is proposed to extract discriminative target features. **CDAN** [19] is one of the state-of-the art methods on UDA. This method performs domain alignment on features that are conditioned on the output of classifiers. In addition, it utilizes entropy minimization on target examples. CDAN integrates domain-classifier based alignment and entropy minimization. Comparison with the UDA methods (DANN, ADR, CDAN) reveals how much gain will be ob-

Net	Method	R to C		R to P		P to C		C to S		S to P		R to S		P to R		MEAN	
		1-shot	3-shot														
AlexNet	S+T	43.3	47.1	42.4	45.0	40.1	44.9	33.6	36.4	35.7	38.4	29.1	33.3	55.8	58.7	40.0	43.4
	DANN	43.3	46.1	41.6	43.8	39.1	41.0	35.9	36.5	36.9	38.9	32.5	33.4	53.6	57.3	40.4	42.4
	ADR	43.1	46.2	41.4	44.4	39.3	43.6	32.8	36.4	33.1	38.9	29.1	32.4	55.9	57.3	39.2	42.7
	CDAN	46.3	46.8	45.7	45.0	38.3	42.3	27.5	29.5	30.2	33.7	28.8	31.3	56.7	58.7	39.1	41.0
	ENT	37.0	45.5	35.6	42.6	26.8	40.4	18.9	31.1	15.1	29.6	18.0	29.6	52.2	60.0	29.1	39.8
	MME	<b>48.9</b>	<b>55.6</b>	<b>48.0</b>	<b>49.0</b>	<b>46.7</b>	<b>51.7</b>	<b>36.3</b>	<b>39.4</b>	<b>39.4</b>	<b>43.0</b>	<b>33.3</b>	<b>37.9</b>	<b>56.8</b>	<b>60.7</b>	<b>44.2</b>	<b>48.2</b>
VGG	S+T	49.0	52.3	55.4	56.7	47.7	51.0	43.9	48.5	50.8	55.1	37.9	45.0	69.0	71.7	50.5	54.3
	DANN	43.9	56.8	42.0	57.5	37.3	49.2	46.7	48.2	51.9	55.6	30.2	45.6	65.8	70.1	45.4	54.7
	ADR	48.3	50.2	54.6	56.1	47.3	51.5	44.0	49.0	50.7	53.5	38.6	44.7	67.6	70.9	50.2	53.7
	CDAN	57.8	58.1	57.8	59.1	51.0	57.4	42.5	47.2	51.2	54.5	42.6	49.3	71.7	74.6	53.5	57.2
	ENT	39.6	50.3	43.9	54.6	26.4	47.4	27.0	41.9	29.1	51.0	19.3	39.7	68.2	72.5	36.2	51.1
	MME	<b>60.6</b>	<b>64.1</b>	<b>63.3</b>	<b>63.5</b>	<b>57.0</b>	<b>60.7</b>	<b>50.9</b>	<b>55.4</b>	<b>60.5</b>	<b>60.9</b>	<b>50.2</b>	<b>54.8</b>	<b>72.2</b>	<b>75.3</b>	<b>59.2</b>	<b>62.1</b>
ResNet	S+T	55.6	60.0	60.6	62.2	56.8	59.4	50.8	55.0	56.0	59.5	46.3	50.1	71.8	73.9	56.9	60.0
	DANN	58.2	59.8	61.4	62.8	56.3	59.6	52.8	55.4	57.4	59.9	52.2	54.9	70.3	72.2	58.4	60.7
	ADR	57.1	60.7	61.3	61.9	57.0	60.7	51.0	54.4	56.0	59.9	49.0	51.1	72.0	74.2	57.6	60.4
	CDAN	65.0	69.0	64.9	67.3	63.7	68.4	53.1	57.8	63.4	65.3	54.5	59.0	73.2	78.5	62.5	66.5
	ENT	65.2	71.0	65.9	69.2	65.4	71.1	54.6	60.0	59.7	62.1	52.1	61.1	75.0	<b>78.6</b>	62.6	67.6
	MME	<b>70.0</b>	<b>72.2</b>	<b>67.7</b>	<b>69.7</b>	<b>69.0</b>	<b>71.7</b>	<b>56.3</b>	<b>61.8</b>	<b>64.8</b>	<b>66.8</b>	<b>61.0</b>	<b>61.9</b>	<b>76.1</b>	78.5	<b>66.4</b>	<b>68.9</b>

Table 1: Results on LSDAC dataset (%). The table have results of one-shot and three-shot setting on 4 domains, R: Real, C: Clipart, P: Clipart, S: Sketch. Our method outperformed other baselines for all adaptation scenarios and all three networks except for only one case where our method showed equivalent performance to ENT.

Net	Method	Office-Home		Office	
		1-shot	3-shot	1-shot	3-shot
AlexNet	S+T	44.1	50.0	50.2	61.8
	DANN	45.1	50.3	<b>55.8</b>	64.8
	ADR	44.5	49.5	50.6	61.3
	CDAN	41.2	46.2	49.4	60.8
	ENT	38.8	50.9	48.1	65.1
	MME	<b>49.2</b>	<b>55.2</b>	<b>56.5</b>	<b>67.6</b>
VGG	S+T	57.4	62.9	68.7	73.3
	DANN	60.0	63.9	69.8	75.0
	ADR	57.4	63.0	69.4	73.7
	CDAN	55.8	61.8	65.9	72.9
	ENT	51.6	64.8	70.6	75.3
	MME	<b>62.7</b>	<b>67.6</b>	<b>73.4</b>	<b>77.0</b>

Table 2: Results on Office-Home and Office dataset (%). The value is the accuracy averaged over all adaptation scenarios. Performance on each setting is summarized in supplementary material.

tained compared to the existing domain alignment-based methods. ENT [12] is a model trained with labeled source and target and unlabeled target using entropy minimization. Entropy is calculated on unlabeled target examples and entire network is trained to minimize it. The difference from MME is that ENT does not have a maximization process. Comparison with this baseline clarifies the importance of the entropy maximization process.

Please note that all methods except for CDAN are trained with exactly the same architecture used in our method in this experiment. In case of CDAN, we could not find advantage of using the architecture. The detail of baseline implementations is in our supplemental material.

## 4.2. Results

**Overview.** The main results on the LSDAC dataset are shown in Table 1. First of all, our method outperformed other baselines for all adaptation scenarios and all three networks except for one case where our method showed equivalent performance to ENT. On average, our method outperformed S+T with 9.5 % and 8.9 % in ResNet one-shot and three-shot setting respectively. The results on Office-Home and Office are summarized in Table 2. Due to the limited space, we show the results averaged on all adaptation scenarios. Our method outperformed all baselines on these datasets too.

**Comparison with UDA Methods.** Generally, baseline UDA methods (DANN, ADR and CDAN) need strong base networks such as VGG or ResNet to perform better than S+T. Interestingly, these methods cannot improve the performance in some cases. The superiority of MME over existing UDA methods is supported by Tables 1 and 2. Since CDAN uses entropy minimization and ENT significantly hurts the performance for AlexNet and VGG, CDAN does not consistently improve the performance for AlexNet and VGG.

**Comparison with Entropy Minimization.** ENT does

Method	R - C	R - P	P - C	C - S	S - P	R - S	P - R	Avg
Source	41.1	42.6	37.4	30.6	30.0	26.3	52.3	37.2
DANN	44.7	36.1	35.8	33.8	<b>35.9</b>	27.6	49.3	37.6
ADR	40.2	40.1	36.7	29.9	30.6	25.9	51.5	36.4
CDAN	44.2	39.1	37.8	26.2	24.8	24.3	<b>54.6</b>	35.9
ENT	33.8	43.0	23.0	22.9	13.9	12.0	51.2	28.5
MME	<b>47.6</b>	<b>44.7</b>	<b>39.9</b>	<b>34.0</b>	33.0	<b>29.0</b>	53.5	<b>40.2</b>

Table 3: Results on the LSDAC dataset in the unsupervised domain adaptation setting (%).

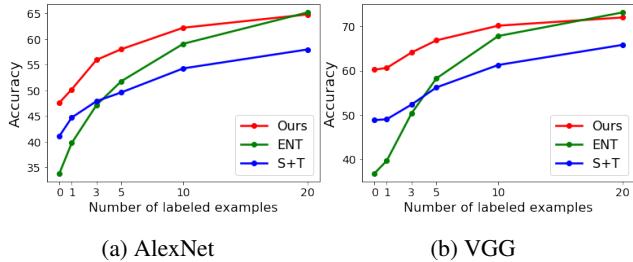


Figure 4: Accuracy changes with respect to the number of labeled target examples. The results shows the ENT method needs more labeled examples to get the similar performance to our method.

not improve performances in some cases because it does not account for domain gap. Comparing results on one-shot and three-shot, entropy minimization gains performance with the help of labeled examples. As we have more labeled target examples, the estimation of prototypes will be more accurate without any adaptation methods. In case of ResNet, entropy minimization often improves the performance. There are two reasons. First, ResNet pre-trained on ImageNet have discriminative representations. Therefore, given a few labeled target examples, the model can extract discriminative features, which contributed to the performance gain in entropy minimization. Second, ResNet has batch-normalization (BN) layers [14]. It is reported that BN has the effect of aligning feature distributions [3, 17]. Hence, entropy minimization was done on aligned feature representations, which improved the performance. When there is a large domain gap such as C to S, S to P, and R to S in Table 1, BN is not enough to handle the domain gap. Therefore, our proposed method performs much better than entropy minimization in such cases. We show the analysis on BN in our supplemental material. This analysis reveals the effectiveness of BN for entropy minimization.

### 4.3. Analysis

**Varying Number of Labeled Examples.** First, we show the results on unsupervised domain adaptaton setting in Table 3. Our method performed better than other methods on average. In addition, only our method improved performance compared to source only model in all settings.

Method	R to C		R to S	
	1-shot	3-shot	1-shot	3-shot
S+T (Standard Linear)	41.4	44.3	26.5	28.7
S+T (Few-shot [4, 24])	<b>43.3</b>	<b>47.1</b>	<b>29.1</b>	<b>33.3</b>
MME (Standard Linear)	44.9	47.7	30.0	32.2
MME (Few-shot [4, 24])	<b>48.9</b>	<b>55.6</b>	<b>33.3</b>	<b>37.9</b>

Table 4: Comparison of architecture in LSDAC dataset using AlexNet. We observe the effectiveness of architecture proposed in [4, 24].

Furthermore, we observe the behavior of our method when the number of labeled examples in the target domain varies. The number of labeled examples is varied from 0 to 20. The results are shown in Fig. 4. Our method works much better than S+T given a few labeled examples. On the other hand, ENT needs 5 labeled examples per class to improve the performance. As we have more labeled examples, the performance gap between ENT and ours is reduced. This result is quite reasonable because prototype estimation will become more accurate without any adaptation method as we have more labeled examples.

**Effect of Network Architecture.** We introduce an ablation study on the network architecture proposed in [4, 24] with AlexNet on LSDAC. As shown in Fig. 3, we employ  $\ell_2$  normalization and temperature scaling. In this experiment, we compared it with a model having a standard linear layer without  $\ell_2$  normalization and temperature. The result is shown in Table 4. By using the network architecture proposed in [4, 24], we can improve the performance of both our method and baseline S+T model. S+T model is trained only on source examples and a few labeled target examples. Therefore, we can argue that the network architecture is an effective technique to improve the performance when we are given a few labeled examples from the target domain.

**Feature Visualization.** In addition, we plot the learned features with t-SNE [20] in Fig. 5. We employ the scenario Real to Clipart of LSDAC using AlexNet as pre-trained model. Fig 5 (a-d) visualizes the target features and estimated prototypes. The color of the cross represents its class and black points are the prototypes. With our method, the target features are clustered to their prototypes and do not have a large variance within the class. We visualize features on the source domain (red cross) and target domain (blue cross) in Fig. 5 (e-h). As we discussed in the method section, our method is supposed to minimize domain-divergence. Then, target features are well-aligned with source features with our method. From Fig. 5f, entropy minimization (ENT) also tries to extract discriminative features, but it fails to find domain-invariant prototypes.

**Quantitative Feature Analysis.** We quantitatively investigate the characteristics of the features we obtain using the same adaptation scenario. First, we perform the analysis on the eigenvalues of the covariance matrix of target

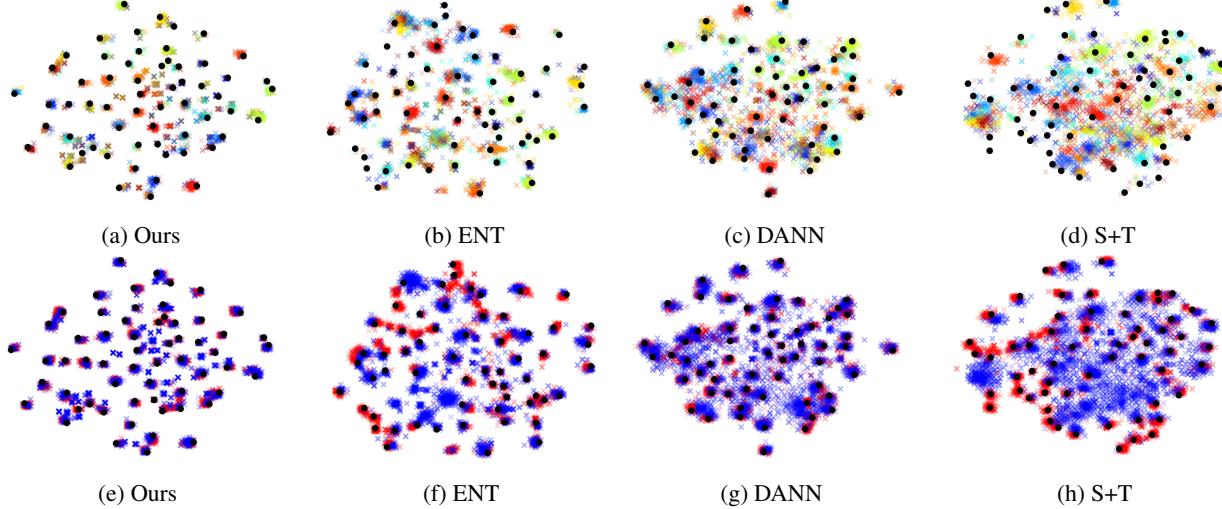


Figure 5: Feature visualization with t-SNE. (a-d) We plot the class prototypes (black circles) and features on the target domain (crosses). The color of a cross represents its class. We observed that features on our method show more discriminative features than other methods. (e-h) Red: Features on the source domain. Blue: Features on the target domain. Our method shows well-aligned features between domains compared to other methods.

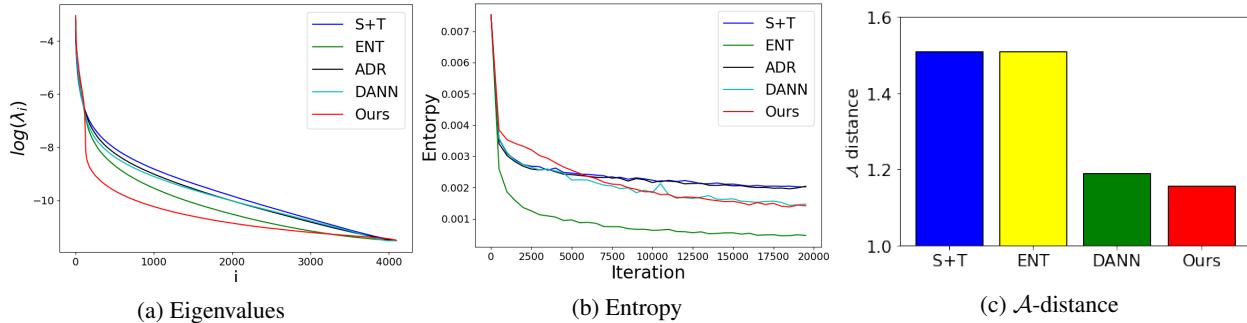


Figure 6: (a) Eigenvalues of the covariance matrix of the features on the target domain. Eigenvalues reduce quickly in our method, which shows that features are more discriminative than other methods. (b) Our method achieves lower entropy than baselines except ENT. (c) Our method clearly reduces domain-divergence compared to S+T.

features. We follow the analysis done in [8]. Eigenvectors represent the components of the features and eigenvalues represent their contributions. If the features are highly discriminative, only a few components are needed to summarize them. Therefore, in such a case, the first few eigenvalues are expected to be large, and the rest to be small. Obviously, the features are summarized by fewer components in our method as shown in Fig. 9a. Second, we show the change of entropy value on the target in Fig. 9b. ENT diminishes the entropy quickly, but results in poor performance. This indicates that the method increases the confidence of predictions incorrectly while our method achieves higher accuracy at the same time. Finally, in Fig. 6c, we calculated  $\mathcal{A}$ -distance by training a SVM as a domain classifier as proposed in [1]. Our method greatly reduces the distance compared to S+T. The claim that our method reduces a domain divergence is empirically supported with

this result.

## 5. Conclusion

We propose a novel Minimax Entropy (MME) approach that adversarially optimizes an adaptive few-shot model for semi-supervised domain adaptation (SSDA). Our model consists of a feature encoding network, followed by a classification layer that computes the features’ similarity to a set of estimated prototypes (representatives of each class). Adaptation is achieved by alternately maximizing the conditional entropy of unlabeled target data with respect to the classifier and minimizing it with respect to the feature encoder. We empirically demonstrate the superiority of our method over many baselines, including conventional feature alignment and few-shot methods, setting a new state of the art for SSDA.

## 6. Acknowledgement

This work was supported by Honda, DARPA and NSF Award No. 1535797.

## References

- [1] S. Ben-David, J. Blitzer, K. Crammer, and F. Pereira. Analysis of representations for domain adaptation. In *NIPS*, 2007. 4, 8
- [2] K. Bousmalis, G. Trigeorgis, N. Silberman, D. Krishnan, and D. Erhan. Domain separation networks. In *NIPS*, 2016. 4
- [3] F. M. Cariucci, L. Porzi, B. Caputo, E. Ricci, and S. R. Bulò. Autodial: Automatic domain alignment layers. In *ICCV*, 2017. 7, 10
- [4] W.-Y. Chen, Y.-C. Liu, Z. Kira, Y.-C. F. Wang, and J.-B. Huang. A closer look at few-shot classification. *arXiv*, 2018. 1, 3, 5, 7, 10
- [5] Y. Chen, W. Li, C. Sakaridis, D. Dai, and L. Van Gool. Domain adaptive faster r-cnn for object detection in the wild. In *CVPR*, 2018. 2
- [6] Z. Dai, Z. Yang, F. Yang, W. W. Cohen, and R. R. Salakhutdinov. Good semi-supervised learning that requires a bad gan. In *NIPS*, 2017. 2
- [7] J. Donahue, J. Hoffman, E. Rodner, K. Saenko, and T. Darrell. Semi-supervised domain adaptation with instance constraints. In *CVPR*, 2013. 2
- [8] A. Dubey, O. Gupta, R. Raskar, and N. Naik. Maximum-entropy fine grained classification. In *NIPS*, 2018. 8
- [9] A. Erkan and Y. Altun. Semi-supervised learning via generalized maximum entropy. In *AISTATS*, 2010. 3
- [10] Y. Ganin and V. Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*, 2014. 1, 2, 3, 4, 5, 10
- [11] S. Gidaris and N. Komodakis. Dynamic few-shot visual learning without forgetting. In *CVPR*, 2018. 1, 3
- [12] Y. Grandvalet and Y. Bengio. Semi-supervised learning by entropy minimization. In *NIPS*, 2005. 3, 6
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5
- [14] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv*, 2015. 7
- [15] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 1, 5
- [16] S. Laine and T. Aila. Temporal ensembling for semi-supervised learning. *arXiv*, 2016. 3
- [17] Y. Li, N. Wang, J. Shi, J. Liu, and X. Hou. Revisiting batch normalization for practical domain adaptation. *arXiv*, 2016. 7, 10
- [18] M. Long, Y. Cao, J. Wang, and M. I. Jordan. Learning transferable features with deep adaptation networks. In *ICML*, 2015. 1, 2
- [19] M. Long, Z. Cao, J. Wang, and M. I. Jordan. Conditional adversarial domain adaptation. In *NIPS*, 2018. 1, 2, 5, 10
- [20] L. v. d. Maaten and G. Hinton. Visualizing data using t-sne. *JMLR*, 9(11):2579–2605, 2008. 7
- [21] T. Miyato, S.-i. Maeda, M. Koyama, K. Nakae, and S. Ishii. Distributional smoothing with virtual adversarial training. *arXiv*, 2015. 3, 10, 12
- [22] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. 2017. 5
- [23] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang. Moment matching for multi-source domain adaptation. *arXiv*, 2018. 5
- [24] R. Ranjan, C. D. Castillo, and R. Chellappa. L2-constrained softmax loss for discriminative face verification. *arXiv*, 2017. 3, 5, 7
- [25] S. Ravi and H. Larochelle. Optimization as a model for few-shot learning. *arXiv*, 2016. 3
- [26] K. Saenko, B. Kulis, M. Fritz, and T. Darrell. Adapting visual category models to new domains. In *ECCV*, 2010. 2, 5
- [27] K. Saito, Y. Ushiku, T. Harada, and K. Saenko. Adversarial dropout regularization. In *ICLR*, 2018. 1, 2, 5, 10
- [28] K. Saito, Y. Ushiku, T. Harada, and K. Saenko. Strong-weak distribution alignment for adaptive object detection. *arXiv*, 2018. 2
- [29] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*, 2018. 2
- [30] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *NIPS*, 2016. 2
- [31] S. Sankaranarayanan, Y. Balaji, A. Jain, S. N. Lim, and R. Chellappa. Learning from synthetic data: Addressing domain shift for semantic segmentation. In *CVPR*, 2018. 2
- [32] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv*, 2014. 5
- [33] J. Snell, K. Swersky, and R. Zemel. Prototypical networks for few-shot learning. In *NIPS*, 2017. 3
- [34] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell. Adversarial discriminative domain adaptation. In *CVPR*, 2017. 4
- [35] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv*, 2014. 1, 2, 3, 4, 5
- [36] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Panchanathan. Deep hashing network for unsupervised domain adaptation. In *CVPR*, 2017. 5
- [37] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, et al. Matching networks for one shot learning. In *NIPS*, 2016. 3
- [38] T. Yao, Y. Pan, C.-W. Ngo, H. Li, and T. Mei. Semi-supervised domain adaptation with subspace learning for visual recognition. In *CVPR*, 2015. 2

## Supplemental Material

### 1. Datasets

First, we show the examples of datasets we employ in the experiments in Fig 7. We also attach a list of classes used in our experiments on LSDAC with this material.

## 2. Implementation Detail

We provide details of our implementation. **We will publish our implementation upon acceptance.** The reported performance in the main paper is obtained by one-time training. In this material, we also report both average and variance on multiple runs and results on different dataset splits (i.e., different train/val split).

**Implementation of MME.** For VGG and AlexNet, we replace the last linear layer with randomly initialized linear layer. With regard to ResNet34, we remove the last linear layer and add two fully-connected layers following [19]. We use the momentum optimizer where the initial learning rate is set 0.01 for all fully-connected layers whereas it is set 0.001 for other layers including convolution layers and batch-normalization layers. We employ learning rate annealing strategy proposed in [10]. Each mini-batch consists of labeled source, labeled target and unlabeled target images. Labeled examples and unlabeled examples are separately forwarded. We sample  $s$  labeled source and labeled target images and  $2s$  unlabeled target images.  $s$  is set to be 32 for AlexNet, but 24 for VGG and ResNet due to GPU memory constraints. We use horizontal-flipping and random-cropping based data augmentation for all training images.

### 2.1. Baseline Implementation

Except for CDAN, we implemented all baselines by ourselves. **S+T** [4]. This approach only uses labeled source and target examples with the cross-entropy loss for training.

**DANN** [10]. We train a domain classifier on the output of the feature extractor. It has three fully-connected layers with relu activation. The dimension of the hidden layer is set 512. We use a sigmoid activation only for the final layer. The domain classifier is trained to distinguish source examples and unlabeled target examples.

**ADR** [27]. We put dropout layer with 0.1 dropout rate after l2-normalization layer. For unlabeled target examples, we calculate sensitivity loss and trained  $C$  to maximize it whereas trained  $F$  to minimize it. We also implemented  $C$  with deeper layers, but could not find improvement.

**ENT.** The difference from MME is that the entire network is trained to minimize entropy loss for unlabeled examples in addition to classification loss.

**CDAN** [19]. We used the official implementation of CDAN provided in <https://github.com/thuml/CDAN>. For brevity, CDAN in our paper denotes CDAN+E in their paper. We changed their implementation so that the model is trained with labeled target examples. Similar to DANN, the domain classifier of CDAN is trained to distinguish source examples and unlabeled target examples.

## 3. Additional Results Analysis

**Results on Office-Home and Office.** In Table 5 and Table 6, we report all results on Office-Home and Office. In almost all settings, our method outperformed baseline methods.

**Sensitivity to hyper-parameter  $\lambda$ .** In Fig. 8, we show our method’s performance when varying the hyper-parameter  $\lambda$  which is the trade-off parameter between classification loss on labeled examples and entropy on unlabeled target examples. The best validation result is obtained when  $\lambda$  is 0.1. From the result on validation, we set  $\lambda$  0.1 in all experiments.

**Changes in accuracy during training.** We show the learning curve during training in Fig 9. Our method gradually increases the performance whereas others quickly converges.

**Comparison with virtual adversarial training.** Here, we present the comparison with general semi-supervised learning algorithm. We select virtual adversarial training (VAT) [21] as the baseline because the method is one of the state-of-the art algorithms on semi-supervised learning and works well on various settings. The work proposes a loss called virtual adversarial loss. The loss is defined as the robustness of the conditional label distribution around each input data point against local perturbation. We add the virtual adversarial loss for unlabeled target examples in addition to classification loss. We employ hyper-parameters used in the original implementation because we could not see improvement in changing the parameters. We show the results in Table 7. We do not observe the effectiveness of VAT in SSDA. This could be due to the fact that the method does not consider the domain-gap between labeled and unlabeled examples. In order to boost the performance, it should be better to account for the gap.

**Analysis of Batch Normalization.** We investigate the effect of BN and analyze the behavior of entropy minimization and our method with ResNet. When training all models, unlabeled target examples and labeled examples are forwarded separately. Thus, the BN stats are calculated separately between unlabeled target and labeled ones. Some previous work [3, 17] have demonstrated that this operation can reduce domain-gap. We call this batch strategy as a “Separate BN”. To analyze the effect of Separate BN, we compared this with a “Joint BN” where we forwarded unlabeled and labeled examples at once. BN stats are calculated jointly and Joint BN will not help to reduce domain-gap. We compare ours with entropy minimization on both Separate BN and Joint BN. Entropy minimization with Joint BN performs much worse than Separate BN as shown in Table 8. This results show that entropy minimization does not reduce domain-gap by itself. On the other hand, our method works well even in case of Joint BN. This is because our training method is designed to reduce domain-gap.



Figure 7: Example images in LSDAC, Office-Home, and Office.

Network	Method	R to C	R to P	R to A	P to R	P to C	P to A	A to P	A to C	A to R	C to R	C to A	C to P	Mean
<b>One-shot</b>														
AlexNet	S+T	37.5	63.1	44.8	54.3	31.7	31.5	48.8	31.1	53.3	48.5	33.9	50.8	44.1
	DANN	<b>42.5</b>	64.2	45.1	56.4	36.6	32.7	43.5	34.4	51.9	51.0	33.8	49.4	45.1
	ADR	37.8	63.5	45.4	53.5	32.5	32.2	49.5	31.8	53.4	49.7	34.2	50.4	44.5
	CDAN	36.1	62.3	42.2	52.7	28.0	27.8	48.7	28.0	51.3	41.0	26.8	49.9	41.2
	ENT	26.8	65.8	45.8	56.3	23.5	21.9	47.4	22.1	53.4	30.8	18.1	53.6	38.8
	MME	42.0	<b>69.6</b>	<b>48.3</b>	<b>58.7</b>	<b>37.8</b>	<b>34.9</b>	<b>52.5</b>	<b>36.4</b>	<b>57.0</b>	<b>54.1</b>	<b>39.5</b>	<b>59.1</b>	<b>49.2</b>
VGG	S+T	39.5	75.3	61.2	71.6	37.0	52.0	63.6	37.5	69.5	64.5	51.4	65.9	57.4
	DANN	<b>52.0</b>	75.7	62.7	72.7	45.9	51.3	64.3	44.4	68.9	64.2	52.3	65.3	60.0
	ADR	39.7	76.2	60.2	71.8	37.2	51.4	63.9	39.0	68.7	64.8	50.0	65.2	57.4
	CDAN	43.3	75.7	60.9	69.6	37.4	44.5	67.7	39.8	64.8	58.7	41.6	66.2	55.8
	ENT	23.7	77.5	64.0	<b>74.6</b>	21.3	44.6	66.0	22.4	70.6	62.1	25.1	67.7	51.6
	MME	49.1	<b>78.7</b>	<b>65.1</b>	74.4	<b>46.2</b>	<b>56.0</b>	<b>68.6</b>	<b>45.8</b>	<b>72.2</b>	<b>68.0</b>	<b>57.5</b>	<b>71.3</b>	<b>62.7</b>
<b>Three-shot</b>														
AlexNet	S+T	44.6	66.7	47.7	57.8	44.4	36.1	57.6	38.8	57.0	54.3	37.5	57.9	50.0
	DANN	47.2	66.7	46.6	58.1	44.4	36.1	57.2	39.8	56.6	54.3	38.6	57.9	50.3
	ADR	45.0	66.2	46.9	57.3	38.9	36.3	57.5	40.0	57.8	53.4	37.3	57.7	49.5
	CDAN	41.8	69.9	43.2	53.6	35.8	32.0	56.3	34.5	53.5	49.3	27.9	56.2	46.2
	ENT	44.9	70.4	47.1	60.3	41.2	34.6	60.7	37.8	60.5	58.0	31.8	63.4	50.9
	MME	<b>51.2</b>	<b>73.0</b>	<b>50.3</b>	<b>61.6</b>	<b>47.2</b>	<b>40.7</b>	<b>63.9</b>	<b>43.8</b>	<b>61.4</b>	<b>59.9</b>	<b>44.7</b>	<b>64.7</b>	<b>55.2</b>
VGG	S+T	49.6	78.6	63.6	72.7	47.2	55.9	69.4	47.5	73.4	69.7	56.2	70.4	62.9
	DANN	56.1	77.9	63.7	73.6	52.4	56.3	69.5	50.0	72.3	68.7	56.4	69.8	63.9
	ADR	49.0	78.1	62.8	73.6	47.8	55.8	69.9	49.3	73.3	69.3	56.3	71.4	63.0
	CDAN	50.2	80.9	62.1	70.8	45.1	50.3	74.7	46.0	71.4	65.9	52.9	71.2	61.8
	ENT	48.3	81.6	65.5	76.6	46.8	56.9	73.0	44.8	<b>75.3</b>	<b>72.9</b>	59.1	<b>77.0</b>	64.8
	MME	<b>56.9</b>	<b>82.9</b>	<b>65.7</b>	<b>76.7</b>	<b>53.6</b>	<b>59.2</b>	<b>75.7</b>	<b>54.9</b>	<b>75.3</b>	<b>72.9</b>	<b>61.1</b>	76.3	<b>67.6</b>

Table 5: Results on Office-Home. Our method performs better than baselines in most settings.

**Results on Multiple Runs.** We investigate the stability of our method and several baselines. Table 10 shows results averaged accuracy and standard deviation of three runs. The deviation is not large and we can say that our method is stable.

**Results on Different Splits.** We investigate the stability

of our method for labeled target examples. Table 9 shows results on different splits.  $sp0$  corresponds to the split we use in the experiment on our paper. For each split, we randomly picked up labeled training examples and validation examples. Our method consistently performs better than other methods.

Network	Method	W to A		D to A	
		1-shot	3-shot	1-shot	3-shot
AlexNet	S+T	50.4	61.2	62.4	50.0
	DANN	57.0	64.4	65.2	54.5
	ADR	50.2	61.2	61.4	50.9
	CDAN	50.4	60.3	61.4	48.5
	ENT	50.7	64.0	66.2	50.0
	MME	<b>57.2</b>	<b>67.3</b>	<b>67.8</b>	<b>55.8</b>
VGG	S+T	69.2	73.2	68.2	73.3
	DANN	69.3	75.4	70.4	74.6
	ADR	69.7	73.3	69.2	74.1
	CDAN	65.9	74.4	64.4	71.4
	ENT	69.1	75.4	72.1	75.1
	MME	<b>73.1</b>	<b>76.3</b>	<b>73.6</b>	<b>77.6</b>

Table 6: Results on Office. Our method outperformed other baselines in all settings.

Method	R to C	R to P	P to C	C to P	C to S	S to P	R-S	P to R
S+T	47.1	45.0	44.9	35.9	36.4	38.4	33.3	58.7
VAT	46.1	43.8	44.3	35.8	35.6	38.2	31.8	57.7
MME	55.6	49.0	51.7	40.2	39.4	43.0	37.9	60.7

Table 7: Comparison with VAT [21] using AlexNet on LSDAC. VAT does not perform better than S+T

Method	Joint BN	Separate BN
ENT	63.6	68.9
MME	69.5	69.6

Table 8: Ablation study of batch-normalization. The performance of the ENT method highly depends on the choice of BN while our method shows consistent behavior.

Method	1-shot			3-shot		
	sp0	sp1	sp2	sp0	sp1	sp2
S+T	43.3	43.8	43.8	47.1	45.9	48.8
DANN	43.3	44.0	45.4	46.1	43.1	45.3
ENT	37.0	32.9	38.2	45.5	45.4	47.8
MME	<b>48.9</b>	<b>51.2</b>	<b>51.4</b>	<b>55.6</b>	<b>55.0</b>	<b>55.8</b>

Table 9: Results on different training splits on LSDAC, Real to Clipart adaptation scenario using AlexNet.

Method	1-shot	3-shot
CDAN	$62.9 \pm 1.5$	$65.3 \pm 0.1$
ENT	$59.5 \pm 1.5$	$63.6 \pm 1.3$
MME	<b><math>64.3 \pm 0.8</math></b>	<b><math>66.8 \pm 0.4</math></b>

Table 10: Results on three runs on LSDAC, Sketch to Painting adaptation scenario using ResNet.

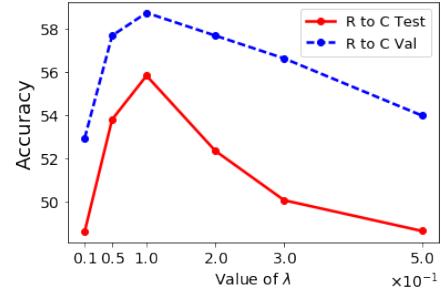


Figure 8: Sensitivity to hyper-parameter  $\lambda$ . The result is obtained when we use AlexNet on LSDAC, Real to Clipart.

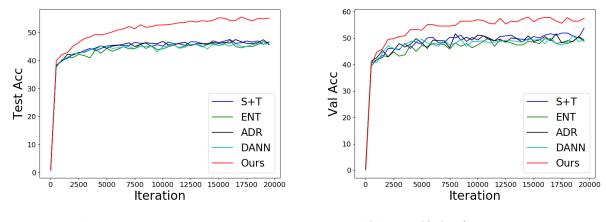


Figure 9: Test and validation accuracy over iterations. Our method increases performances over iterations while others quickly converges. The result is obtained on Real to Clipart adaptation on LSDAC using AlexNet.