

# Knowledge Distillation Methods for Efficient Unsupervised Adaptation Across Multiple Domains

Le Thanh Nguyen-Meidine<sup>1,\*</sup>, Atif Belal<sup>3</sup>, Madhu Kiran<sup>1</sup>, Jose Dolz<sup>1</sup>,  
Louis-Antoine Blais-Morin<sup>2</sup>, Eric Granger<sup>1</sup>

<sup>a</sup>1100 Notre-Dame St W, Montreal, Quebec H3C 1K3, Canada

---

## Abstract

Beyond the complexity of CNNs that require training on large annotated datasets, the domain shift between design and operational data has limited the adoption of CNNs in many real-world applications. For instance, in person re-identification, videos are captured over a distributed set of cameras with non-overlapping view-points. The shift between the source (e.g. lab setting) and target (e.g. cameras) domains may lead to a significant decline in recognition accuracy. Additionally, state-of-the-art CNNs may not be suitable for such real-time applications given their computational requirements. Although several techniques have recently been proposed to address domain shift problems through unsupervised domain adaptation (UDA), or to accelerate/compress CNNs through knowledge distillation (KD), we seek to simultaneously adapt and compress CNNs to generalize well across multiple target domains. In this paper, we propose a progressive KD approach for unsupervised single-target DA (STDA) and multi-target DA (MTDA) of CNNs. Our method for KD-STDA adapts a CNN to a single target domain by distilling from a larger teacher CNN, trained on both target and source domain data in order to maintain its consistency with a common

---

\*Corresponding author

*Email addresses:* [le-thanh.nguyen-meidine.1@ens.etsmtl.ca](mailto:le-thanh.nguyen-meidine.1@ens.etsmtl.ca) (Le Thanh Nguyen-Meidine), [abelal@myamu.ac.in](mailto:abelal@myamu.ac.in) (Atif Belal), [madhu.kiran.1@ens.etsmtl.ca](mailto:madhu.kiran.1@ens.etsmtl.ca) (Madhu Kiran), [jose.dolz@etsmtl.ca](mailto:jose.dolz@etsmtl.ca) (Jose Dolz), [lablaismorin@genetec.com](mailto:lablaismorin@genetec.com) (Louis-Antoine Blais-Morin), [eric.granger@etsmtl.ca](mailto:eric.granger@etsmtl.ca) (Eric Granger)

<sup>1</sup>Laboratoire d'imagerie, de vision et d'intelligence artificielle (LIVIA), Ecole de technologie supérieure, Montreal, Canada.

<sup>2</sup>Genetec Inc., Montreal, Canada.

<sup>3</sup>Aligarh Muslim University, Aligarh, India.

representation. This method is extended to address MTDA problems, where multiple teachers are used to distill multiple target domain knowledge to a common student CNN. A different target domain is assigned to each teacher model for UDA, and they alternatively distill their knowledge to the student model to preserve specificity of each target, instead of directly combining the knowledge from each teacher using fusion methods. Our proposed approach is compared against state-of-the-art methods for compression and STDA of CNNs on the Office31 and ImageClef-DA image classification datasets. It is also compared against state-of-the-art methods for MTDA on Digits, Office31, and OfficeHome. In both settings – KD-STDA and KD-MTDA – results indicate that our approach can achieve the highest level of accuracy across target domains, while requiring a comparable or lower CNN complexity.

*Keywords:* Deep Learning, Convolutional NNs, Knowledge Distillation, Unsupervised Domain Adaptation, CNN Acceleration and Compression

---

## 1. Introduction

Convolutional neural networks (CNNs) have achieved state-of-the-art performance in many visual recognition tasks, i.e., classification, object detection, and segmentation. However, a key limitation that hampers their scalability is their poor generalization to data from new unseen domains, where there exists a shift between data collected for training, source domain, and one/multiple deployment environments, target domains. Particularly, in real-world video surveillance applications like person re-identification, deep Siamese networks are commonly trained to process videos captured over a distributed network of cameras. The variations in data distributions in terms of the different types of the cameras, viewpoints, and capture conditions can cause a significant domain shift, and degrade accuracy [1]. Another drawback of state-of-the-art CNNs is the computational complexity, which impedes its adoption in real-time applications. Although it is possible to address the domain shift problem by annotating target domain data and applying supervised transfer learning, the

cost of collecting and annotating a large image dataset for each target domain can be prohibitively costly. To alleviate this issue when only unlabeled target data is available, several unsupervised domain adaptation (UDA) techniques have been proposed. Currently, most of UDA techniques focus on learning domain invariant features by minimizing a distance or discrepancy between the two data distributions [2], to encourage domain confusion by using adversarial losses [3, 4], or both [5]. Lastly, it is also possible to learn a mapping between source and target images [6, 7], such that images captured in different domains have a similar appearance, which mimics standard supervised learning.

While these techniques are very useful for single-target domain adaptation (STDA), they fail to adapt to a multi-target domain adaptation (MTDA) scenario, a relatively unexplored setting of practical importance, where a single model is trained to generalize well across multiple different target domains. A simple solution would be to solve the problem of MTDA by having one CNN trained on each target domain. Nevertheless, this would be very costly in terms of computational resources. For instance, in person re-identification, MTDA can be addressed by adapting a CNN per target domain, or blending target data to adapt a common CNN, although in many real-world applications, these solutions are too costly, and may lead to a reduction in accuracy. Current literature tackles the MTDA task by either taking advantage of domain labels from each target domain [8] or ignoring these labels and mix all the target domains together [9]. The Figure 1 illustrates the difference between STDA, MTDA and some possible solutions for an MTDA setting.

Current techniques allow to accelerate and compress CNNs while preserving its accuracy. These approaches include: quantization [10, 11, 12], low-rank approximation [13, 14], knowledge distillation (KD) [15, 16, 17, 18] and network pruning [19, 20, 21, 22, 23]. However, the direction of jointly compressing and adapting CNN to a target domain remains largely unexplored. To the best of our knowledge, only the Transfer Channel Pruning (TCP) technique [24] has been proposed to first solve the UDA problem, and then pruning the adapted model. As for using UDA and KD, while recent work [25, 26] resort to KD to

address the UDA task, they focus on reducing the domain gap, rather than the computational cost. Furthermore, most KD techniques rely on labeled dataset whereas in the UDA setting, we do not have access to these labels, which represents an important challenge for the co-joint optimization of KD and UDA.

While it is possible to adopt a compressed CNN to a new target domain, we argue that this scenario will likely degrade the performance of the CNN since over-parametrization is often important for generalization [27]. Another alternative could be to first adapt a model to the target domain and then compress it. However, we experimentally demonstrate that such a scenario results in poor performance due to the lack of label information for knowledge distillation. In this paper, we resolve the aforementioned problem by progressively distilling knowledge from a teacher model that is continuously doing domain adaptation to the student. We argue that this can improve the performance since the student learns how to adapt to a new domain instead of learning directly from a previously targeted domain. In addition, we overcome the problem of unsupervised KD for the student model by jointly optimizing UDA and KD and by keeping the student model consistent with the source domain. Our proposed approach can also be generalized for a multi-target domain adaptation – MTDA, where there exist multiple target domains– by assigning each teacher to a different target domain.

In order to address the problem of CNN complexity and domain shift over one or multiple domains, this paper introduces the following contributions. 1) we propose a new approach for joint KD and UDA that allows training CNNs such that they generalize well on one or multiple target domains; 2) we introduce a consistency loss that ensures that the student models also learn source domain knowledge from the teacher model in order to overcome the challenges of unsupervised KD on target data; 3) this paper extends our previous work [28] in a substantial way by providing new algorithms and extensive experimental results to validate our approach on different UDA and feature distillation techniques; 4) finally, we extend to the MTDA setting, where our KD-MTDA approach relies on multiple teachers to provide a higher capacity to generalize

for multiple target domains, as well to distill their compressed knowledge to a smaller common student CNN that can perform well across multiple target domains.

The rest of this paper is organised as follows. Section 2 provides an overview of state-of-the-art methods for compression and UDA with CNNs. Then, Section 3 introduces our proposed KD approaches for STDA and MTDA of CNNs. The experimental methodology employed for validation, and the results and discussions are described in Sections 4 and 5, respectively.

## 2. Related Work

### 2.1. CNN acceleration and compression techniques:

Time complexity generally depends mostly on the CNN convolutional layers, while memory complexity (number of parameters) depends mostly on the fully connected layers. Currently, one of the most popular ways to reduce CNN complexity is to apply pruning techniques [19, 20, 21, 22, 23, 29, 30], using either weight or filter-level pruning. Another popular strategy to lower the complexity is quantization [10, 11, 12], which reduces the complexity by reducing the representation of weights into a lower precision, i.e., converting 64-bit float precision to an 8-bit integer. Low rank decomposition techniques [13, 14] have also been employed to accelerate CNNs by decomposing weight tensor into lower rank approximation. Lastly, KD [15, 16, 17, 18] aims to transfer knowledge from a larger teacher model to a smaller student model, therefore, preserving teacher accuracy on the student, while reducing complexity.

In this paper, we focus on KD based techniques on the intuition that a CNN can efficiently learn how to adapt to a target domain, improving its accuracy on target domain data. Also, we argue that by leveraging over-parametrization [27], our approach can provide a better generalization capacity than when applying UDA on a pruned CNN which has fewer parameters to optimize. Additionally, KD also provides a natural generalization for an MTDA setting where each teacher is responsible for a target domain. Two main approaches have been pro-

posed for distilling a teacher’s knowledge to a student – either based on network outputs [15] or on intermediate features [17]. For example, a temperature-based softmax was used in [15] to generate softer versions of the teacher network outputs, facilitating the learning of a student model. In contrast, Heo et al. [17] enforce similarities between the teacher and student models at intermediate features layers, by minimizing a partial L2 distance. In particular, if the value measures on the student model is smaller than the value of the teacher, and both are negative, then the result is set to 0. A margin ReLU is also employed since the authors argue that negative value from the network is also useful. Another work [18] proposes to solve the problem of the gap between teacher and student by integrating a teaching assistant that serves as an intermediate student that first learns from the teacher then it will distill the information to the student using the output of its network.

*2.2. Single-target domain adaptation:*

Unsupervised domain adaptation (UDA) techniques alleviate the problem of domain shift of a deep learning model that is trained on a labeled source domain data, using unlabeled data from a target domain. Currently, UDA can be achieved by encouraging domain confusion [3, 4], learning domain-invariant features [2, 31], mapping between source and target domains [32, 33, 34], ensemble learning [35], statistic normalization [36] and target discriminate methods[37]. Domain confusion between source and target domain has been done by employing an adversarial loss[4, 31] or by using a gradient reversal layer like [3] in combination with a domain classifier. Differently, a model can also learn domain-invariant features, for example, by minimizing the Maximum Mean Discrepancy (MMD) between source and target [2]. Other works such as [32, 33] propose to find a mapping from the source to the target domain or vice-versa, typically based on Generative Adversarial Networks (GAN) [32]. This setting is commonly regarded as a zero-sum game between two networks, a discriminator and a generator. While the generator tries to fool the discriminator by generating images with target style, the discriminator tries to distinguish the domain

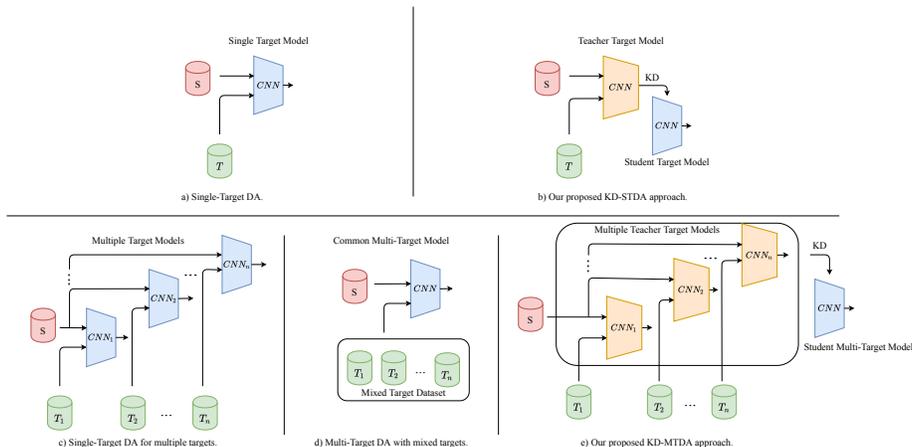


Figure 1: Illustration of the different strategies for STDA and MTDA of CNNs. a) Basic STDA technique given data from a source  $S$  and target  $T$  domain. b) Our proposed KD-STDA approach, where a teacher CNN is adapted to one target domain, and its knowledge is transferred to a smaller student CNN. c) MTDA by applying STDA independently on  $n$  CNNs, one per target domain dataset. d) Applying STDA using a mixture of datasets from  $n$  target domain to adapt a common CNN. e) Our proposed KD-MTDA approach, where multiple teacher CNNs are adapted, one per target, and their knowledge is combined and distilled to a smaller common student CNN.

the images belong to. In [38], the authors further improve mapping by performing the adaptation at the feature level. Another strategy is to employ multiple models as an ensemble or as self-ensembling (multiple models at different times) such that they can produce reliable pseudo-labels on target domains[35]. Last, some other researchers have proposed to adapt the batch norm statistics [36] instead of adapting the network’s layers.

### 2.3. Multi-target domain adaptation

While single target domain adaptation (STDA) has gained significant attention, multi-target domain adaptation (MTDA) remains largely unexplored, since it poses a challenge for many applications. Recent work to tackle the MTDA task [8] proposes to maximize the mutual information between domain labels and domain-specific features while minimizing the mutual information between shared features. Other papers have explored the open multi-domain

adaptation [39] or blended MTDA[9, 40]. The first focuses on doing UDA on multiple target domains on a model so that it can also perform well on another unseen target domain not available during training. The second technique assumes that there’s no information on where the target image comes from and uses a clustering technique to separate them into different pseudo target domains. The Figure 1 shows the difference between our proposed technique and current existing techniques.

*2.4. Joint unsupervised domain adaptation and knowledge distillation:*

Even though joint UDA and KD has already been investigated in the past [41, 42], the aim of these works is on improving the UDA accuracy, neglecting the complexity reduction problem. For example, [41] uses multiple teachers to learn multiple source domains, increasing the accuracy of the student model for sentiment analysis. In [42], KD and DA are combined in a semi-supervised fashion in the context of medical image segmentation. To the best of our knowledge, only the method in [24] currently tackles the problem of compression and domain shift using pruning and UDA. However, unlike the proposed model, the approach in [24], i.e., TCP, operates in several steps. TCP must first be adapted to a target domain, and then the learned model is iteratively pruned while continuously refined using UDA. This contrasts to our model, which progressively compresses and adapts to multiple target domains simultaneously from a single source CNN model. Furthermore, TCP presents some important limitations. First, by performing the UDA and compression as two independent steps, there is a risk of a degrade in performance and a limited compression rate. Second, longer deployment times are needed for obtaining an efficient CNN that is adapted to a target domain, which can hamper the application to many specific problems, for example application that have a lot of changes in environment such as mobile car detection. Lastly, TCP[24] cannot be easily generalized to an MTDA setting, and applying TCP directly on a mixture of multiple target domains can result in poor performance since it has been shown in literature [9] that directly apply an STDA technique on a mixture of multiple targets will

only have limited performance.

### 3. Progressive Domain Adaptation Through Knowledge Distillation

In this paper, we simultaneously address the domain shift and compression problems to obtain an efficient CNN model that is adapted to a target domain. Instead of obtaining a compressed CNN directly from a previously-adapted CNN, we argue that it may be beneficial to learn an efficient CNN by teaching it how to progressively adapt to target domains. This would mean that our efficient model is less constrained to optimize since it does not start from a local minimum. Based on this idea, we propose an approach that relies on a larger teacher model to perform UDA for a target domain, while progressively transferring the knowledge of how to perform UDA to a smaller student. Both student and teacher models begin from a pre-trained CNN, instead of a model previously adapted to a target domain.

The main pipeline of our method is illustrated in Figure 2. The rest of this section provides additional details on UDA techniques used in this paper, KD for both supervised and unsupervised settings, our joint CNN optimization for UDA and KD, and generalization to multiple domains.

#### 3.1. Unsupervised domain adaptation:

Our distillation of domain adaptation is UDA-agnostic, and can be integrated using any UDA approach. We apply our technique on popular discrepancy- and adversarial-based UDA approaches from the literature. First, a Maximum Mean Discrepancy (MMD)-based UDA approach [2], that seeks to minimize the distance between target and source domain distributions in the features space. The choice of this technique for domain adaptation allows for fair comparison of compression of our technique with Transfer Channel Pruning (TCP) [24]. Second, we use a technique called RevGrad [3], which relies on adversarial training in order to find a common representation by encouraging domain confusion between target and source domain. The rest of this section provides additional details on these two approaches.

3.1.1. *MMD approach:*

Let us define the labeled source dataset as  $S = \{x_s, y_s\}$  and the unlabeled target dataset as  $T = \{x_t\}$ . The UDA loss for teacher CNN based on MMD [2, 24] is :

$$\mathcal{L}_{MMD}(\phi, T, S) = \left\| \frac{1}{N_s} \sum_{x_i \in S} \phi(x_i) - \frac{1}{N_t} \sum_{x_j \in T} \phi(x_j) \right\|_{\mathcal{H}}^2 \quad (1)$$

where dataset  $S$  contains  $N_s$  samples (and labels),  $T$  contains  $N_t$  samples,  $\phi$  is the teacher function that maps an input to a feature map, and  $\mathcal{H}$  is the Reproducing Kernel Hilbert Space (RKHS) with Gaussian kernel. As in [2, 24], we add a supervised loss on the source domain dataset, and the overall UDA loss for the teacher CNN is:

$$\mathcal{L}_{TDA}(\Phi, T, S) = \mathcal{L}_{MMD} + \gamma \mathcal{L}_{CE}(\Phi(x_s, 1), y_s) \quad (2)$$

where  $\mathcal{L}_{CE}$  the supervised cross-entropy loss of the teacher model on the source domain data,  $\gamma$  a trade-off hyper-parameter that follows the same variations as [24], and  $\Phi$  the output of the teacher network with a soft-max of temperature value set to 1 (i.e. the regular soft-max).

3.1.2. *RevGrad approach:*

For this approach, UDA of the teacher relies on a domain classifier, a gradient reversal layer (GRL), and the domain confusion loss defined as:

$$\mathcal{L}_{DC}(\phi, T, S) = \frac{1}{N_s + N_t} \sum_{x \in S \cup T} \mathcal{L}_{CE}(C(\phi(x)), d_l) \quad (3)$$

where  $\phi(x)$  is the output from the feature extractor of teacher network  $\Phi$  (before the fully connected layers),  $C$  is the domain classifier for the corresponding teacher network,  $d_l$  the domain label (source or target),  $N_s$  is the number of samples in the source domain  $S$ , and  $N_{t_i}$  is the number of samples in the target domain  $T_i$ . The overall UDA loss is then defined as:

$$\mathcal{L}_{TDA}(\Phi, T, S) = \frac{1}{N_s} \sum_{x_s, y_s \in S} \mathcal{L}_{CE}(\Phi(x_s), y_s) + \alpha \cdot \mathcal{L}_{DC}(\phi, T, S) \quad (4)$$

The cross-entropy loss term in Equ.4 allows the supervised training of the teacher model on the source domain data to ensure the consistency of domain

confusion. The second term is controlled by a hyper-parameter  $\alpha$  that regulates the importance of the domain confusion loss which is maximized using a gradient reversal layer.

### 3.2. Knowledge distillation for domain knowledge transfer:

The next step consists in transferring target domain knowledge from the teacher to student models. Our proposed method is general, and can be adapted to any KD method based on logits and features. We apply our proposed approach on a KD technique based on (1) the work of [15] that distill knowledge from a teacher to student using the output of each network, and (2) the feature distillation method [17] that minimizes the difference between feature maps of intermediate layers. These KD techniques are detailed below.

#### 3.2.1. Logits-based distillation:

Let us defined the temperature based softmax that inputs the logits, and produce a softened output that can represent more information to be distilled:

$$\text{softmax}_i(\tau) = \frac{\exp(z_i \cdot \tau)}{\sum_j \exp(z_j \cdot \tau)} \quad (5)$$

where  $z_i$  the logits of class  $i$  produced by the classification layer of the CNN, the Hinton distillation [15], is then defined as:

$$\mathcal{L}_{LDistill}(\Phi, \Theta, D) = \mathcal{L}_{KL}(\Theta(D, \tau), \Phi(D, \tau)) \quad (6)$$

where  $\Theta(D, \tau)$  and  $\Phi(D, \tau)$  represent the outputs of the student and teacher networks, respectively, with a softmax based on a temperature  $\tau$  in order to soften the output. For simplicity,  $D$  represents a generic dataset that we will replace with the source or target dataset.

#### 3.2.2. Feature-based distillation:

Instead of distilling only on the output of teacher and student, feature KD [17] relies on several intermediate layers inside the network. For simplicity, we present the KD loss for only one layer since the algorithm similar across

all layers. For this technique, the features for distillation is extracted at the end of each residual block and before a ReLU unit since the information before ReLU still carries negative values which are important for KD. Additionally, the authors argue that not all negative values are important, therefore they propose applying a margin ReLU when distilling information from teacher to student. The margin ReLU is defined as:

$$\sigma_m(x) = \max(x, m) \quad (7)$$

where  $m$  is a margin value less than zero. The margin  $m$  can be computed as the expectation value of the negative response, for each input channel. The margin value for channel  $c$  can be computed as:

$$m_c = E[F_{\Phi}^i | F_{\Phi}^i < 0, i \in C] \quad (8)$$

with  $F_{\Phi}^i$  the teacher’s feature of the  $i$  –  $th$  channel of  $C$  channel. The margin value  $m_c$  is computed over all training samples. As for the distillation loss, the authors use a partial L2 distance function, defined as:

$$\mathcal{L}_{PL2}(F_{\Phi}, F_{\Theta}) = \sum_j^{WHC} \begin{cases} 0, & \text{if } F_{\Theta_j} \leq F_{\Phi_j} \leq 0 \\ (F_{\Phi_j} - F_{\Theta_j})^2, & \text{otherwise} \end{cases} \quad (9)$$

In this Equ. 9,  $F_{\Theta}, F_{\Phi} \in \mathcal{R}^{W \times H \times C}$  represents respectively the feature map of the student and teacher networks. Here, the  $j$  –  $th$  component represents an element in the feature map where  $F_{\Theta_j}, F_{\Phi_j} \in \mathcal{R}$ . After the domain adaptation in 3.1, the target domain knowledge is transferred from the teacher to the student, we use a loss function:

$$\mathcal{L}_{FDistill}(\Phi, \Theta, D) = \mathcal{L}_{PL2}(\sigma_{m_c}(F_{\Phi}(D)), r(F_{\Theta}(D))) \quad (10)$$

Here,  $\sigma_{m_c}$  represents margin ReLU and  $r(\cdot)$  represents  $1 \times 1$  convolution layer.  $F^T$  and  $F^S$  represent respectively the features of student network and teacher network. In the next section, we will further explain our contribution that allows employing unsupervised KD for UDA.

### 3.3. Progressive Joint KD and UDA: KD-STDA

The main pipeline of the proposed method is depicted in Figure 2. Our method performs UDA of a teacher model by learning a common domain-invariant feature embedding between the source and target domains, while progressively distilling its knowledge to a student model. Contrary to other works, e.g., [24], the teachers in our setting are not adapted to the target domain, and start from pre-trained weights on ImageNet. Thus, our teachers would train the student CNN to adapt to the target domain, while they also continuously adapt to each target domain. From Equation 6, target distillation is a simple application of the equation on the target dataset:

$$\mathcal{L}_{TKD}(\Phi, \Theta, T) = \mathcal{L}_{Distill}(\Phi, \Theta, T) \quad (11)$$

Equ.11 should allow to transfer information to the student network, since we are only interested in performing correctly on the target domain data. Inspired by current UDA techniques, we propose to add a consistency distillation loss in order to ensure that our student learns a common representation. This consistency is achieved by distilling source domain information into the student CNN:

$$\mathcal{L}_{SKD}(\Phi, \Theta, S) = \mathcal{L}_{distill}(\Phi, \Theta, S) + \alpha \mathcal{L}_{CE}(S) \quad (12)$$

Eq. 12 is the student KD loss, with hyper-parameter  $\alpha$  to balance between the KD and the cross entropy loss of the output of the student model and the ground truth on the source domain. We also propose to add the  $\beta$  hyper-parameter in order to balance out the importance between UDA and KD. Since we are performing jointly KD and DA, in the beginning, the teacher would still be learning from the DA. This means that there is not much to be learned for the student model, besides the source representation which can be learned from the KD loss. In light of this, we propose to start by giving more importance to UDA in the beginning and gradually transfer the importance to KD basing  $\beta$  on an exponential growth function between  $[b, f]$ , with  $b$  the starting value of  $\beta$  and  $f$  the end value. In order to define as exponential growth, we need to

calculate a growth rate based on  $b$  and  $f$ :

$$g = \frac{\log(f/b)}{N_e} \quad (13)$$

where  $N_e$  is the number of epochs and  $g$  the growth rate. Given the growth rate,  $\beta$  at epoch  $t$  can be found as:

$$\beta_t = b \cdot \exp\{g \cdot t\} \quad (14)$$

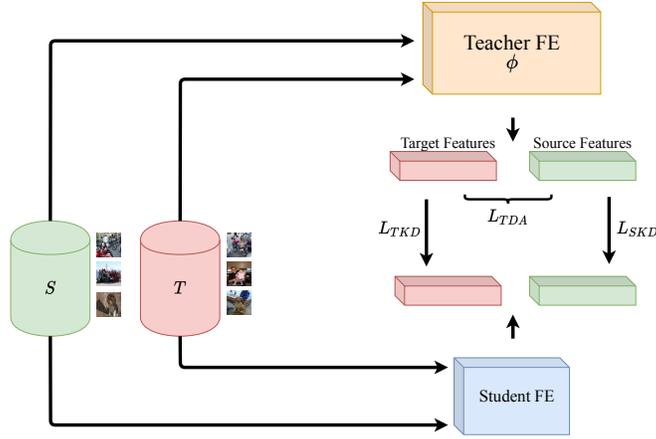


Figure 2: Our proposed progressive KD-STDA, where DA is performed on a teacher CNN, and then a smaller student CNN learns how to adapt on both target and source features through knowledge distillation.

Figure 2 illustrates all these losses and also the proposed techniques. The final loss of our models, is then:

$$\mathcal{L}(\Phi, \Theta, S) = (1 - \beta)\mathcal{L}_{TDA}(\Phi, T, S) + \beta(\mathcal{L}_{TKD}(\Phi, \Theta, T) + \mathcal{L}_{SKD}(\Phi, \Theta, S)) \quad (15)$$

Algorithm 1 described the steps of our approach for the MTDA setting.

### 3.4. Progressive Multi-target KD and UDA: KD-MTDA

A multi-target approach can also be naturally derived from our approach by assigning each teacher a different target domain, and applying the same UDA

---

**Algorithm 1:** KD-STDA: Our method in STDA Setting.

---

**input** : A teacher model  $M_T$ , a student model  $M_S$ , a source dataset  $D_s^{Sup}$ , a target dataset  $D_t^U$

**output** : A target adapted student model

- 1 **for**  $epoch \leftarrow 1$  to  $N_e$  **do**
- 2     **for**  $x_s$  in  $S$  and  $x_t$  in  $T$  **do**
- 3         Optimize the  $(1 - \beta)\mathcal{L}_{DA}$  for the teacher model
- 4         Optimize the  $\beta(\mathcal{L}_{TKD} + \mathcal{L}_{SKD})$  for the student model
- 5         Update  $\beta$  following Eq.14
- 6     **end**
- 7     Evaluate the model
- 8 **end**

---

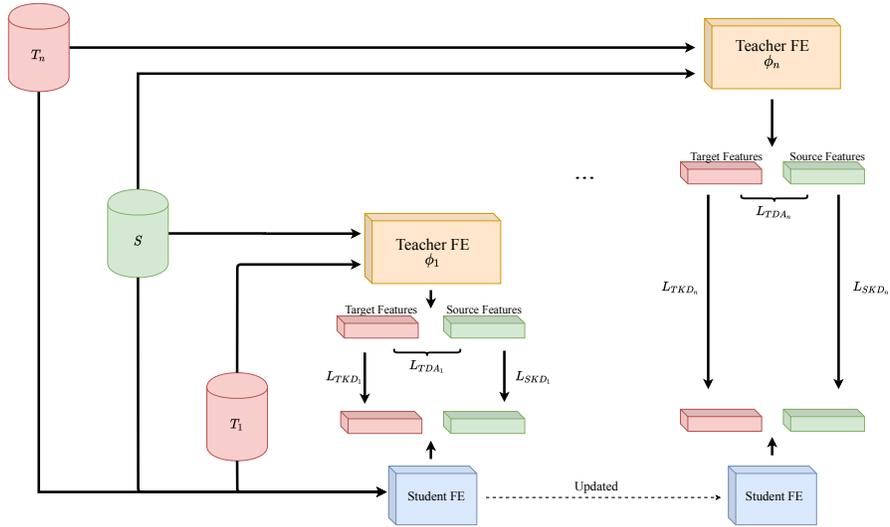


Figure 3: Our proposed method for MTDA with  $n$  target domains. Each teacher CNN is assigned a different target domain, and sequentially teaches the student CNN how to adapt to source and target features for different target domains using knowledge distillation. After each knowledge distillation process with each teacher, features output from the student CNN are recalculated using update weights, for a new distillation with a different teacher.

loss that we previously defined. The progressive KD to the common student is described in Algorithm 2. This approach can be formalized in the same way as in the STDA setting. We define an ensemble of target datasets from  $T_1$  to  $T_n$ , and redefine the loss to a specific target domain dataset for each teacher:

$$\mathcal{L}_{MTDA}(S, T_i) = (1 - \beta)\mathcal{L}_{TDA}(\Phi_i, T_i, S) + \beta(\mathcal{L}_{TKD}(\Phi_i, \Theta, T_i) + \mathcal{L}_{SKD}(\Phi_i, \Theta, S)) \quad (16)$$

where  $\Phi_1$  to  $\Phi_n$  are the respective teachers for target domains  $T_1$  to  $T_n$ . Similar to our STDA setting, we rely on the exponentially growing rate to gradually transfer the importance of domain adaptation to distillation since the student model and teachers are quite similar in the beginning (they all start from pre-trained weights of ImageNet).

Figure 3 illustrated the block diagram of our proposed KD-MTDA method. By assigning each complex teacher a target domain, we can adapt a CNN to generalize well across multiple target domains, thus allowing the student to learn progressively transferable features from all the teachers.

---

**Algorithm 2:** KD-MTDA: Our method in MTDA setting.

---

```

input      : A source domain dataset  $S$ , a set of target dataset  $T_0, T_1, \dots, T_n$ 
output    : A student model adapted to  $n$  targets
1 Initialize a set of teachers models  $\Phi = \{\Phi_0, \Phi_1, \dots, \Phi_n\}$ 
2 Initialize a student model  $\Theta$ 
3 for  $e \leftarrow 1$  to  $N_e$  do
4   for  $x_s \in S$  and  $X_t \in \{T_0, \dots, T_n\}$  do
5     Get the set of data of target domains  $X_t$ 
6     for  $x_t^i \in X_t$  and  $\Phi_i \in \Phi$  do
7       Optimize  $(1 - \beta)\mathcal{L}_{DA}$  for  $\Phi_i$  using  $x_s, x_t^i$ 
8       Optimize the loss of equation  $\beta\mathcal{L}_{SKD}$  for  $\Phi_i$  and  $\Theta$  using  $x_s$ 
9       Optimize the loss of equation  $\beta\mathcal{L}_{TKD}$  for  $\Phi_i$  and  $\Theta$  using  $x_t^i$ 
10    end
11    Update  $\beta = s \cdot \exp^{g \cdot e}$ 
12  end
13  Evaluate the model
14 end

```

---

## 4. Experimental Methodology

### 4.1. Datasets:

*Digits.* This dataset is made of a set of digits datasets: MNIST[43] (**mt**), MNIST-M[44] (**mm**), SVHN[45] (**sv**), and USPS[46] (**up**). Each one has each 10 classes that represent all the digits. For the evaluation on this dataset, we follow the same protocol as in [8] for a fair comparison: where we use 25000 samples for training on **mt,mm,sv,up** and 9000 for testing. On the **up** dataset we use the entire dataset as a domain. We will mainly use this dataset for comparison of MTDA.

*Office31.* [47] This dataset contains three subsets of dataset which are Webcam (**W**), DSLR (**D**) and Amazon (**A**). These subsets contain images from amazon.com (**A**) or office product taken a DSLR camera (**D**) or a webcam (**W**). We evaluate our results based on the same scenario as [24]. These datasets all have 31 common classes and around 4000 images in total and each subset has a different amount of data: Amazon (2800 images), DSLR (500 images), Webcam (800 images).

*ImageClef-DA.* This dataset for UDA has four subsets which are taken from Imagenet (**I**), Pascal-Voc (**P**), Caltech (**C**) and Bing (**B**). Each of these subsets contains 600 images for 12 classes. For this dataset, similar to Office31, we use the same scenario as [24] for fair comparison.

*PACS.* [48] While this dataset is often used for domain generalization, it has been used in [8] for MTDA. This dataset contains 4 subsets – Art painting (**Ap**), Photo (**P**), Cartoon (**Cr**) and Sketch (**S**). With around 2000 images in Art Painting, 2300 for Cartoon, 1700 for Photo and 4000 for Sketch. We will use this dataset for MTDA comparison.

*OfficeCaltech.* This dataset contains 4 subsets of four different domains: Amazon (**A**), Caltech10 (**C**), DSLR(**D**) and Webcam (**W**). They each contain respectively 957, 1123, 295, 157 images of 10 classes.



Figure 4: Samples from the Office31 and ImageClef-DA datasets.

#### 4.2. Baselines methods:

For experiments in the STDA setting, several baselines were selected in order to evaluate our method. We start with the three scenarios presented earlier: (1) UDA then KD, (2) KD then UDA and (3) UDA directly on compact model. For these baselines, our teacher model will always be ResNet50 and our student models ResNet18 or ResNet34. In the first case,  $UDA \rightarrow KD$ , UDA is applied to a teacher model then we start KD to a student model on target domain without labels. For the second baseline  $KD \rightarrow UDA$ , the teacher model is first trained in supervised way on the source dataset then we apply KD with this teacher and a student using labeled data of the source dataset. The student is then applied UDA. As for the last baseline, UDA is applied directly to a student model. We also add RevGrad[3] and CDAN[49] directly on compact model as another baseline. In addition, we also compare to the state-of-the-art method in compression and UDA: TCP[24].

In the MTDA setting, we compare our method with several baselines like a lower-bound where the model is only trained on source domain and then test on multiple target domains. We also introduce baselines from current STDA techniques, RevGrad[3] and ADDA[4] to show that, in their current form, they are

not suitable for an MTDA setting. Additionally, we also compare our method with [8] which is the current state-of-the-art MTDA technique that uses domain labels. As for [39], we do not include them in our comparison since it focuses on the open MTDA with unseen target domain, however we provide a comparison with DADA[40] which does use domain labels.

#### 4.3. Implementation details:

For the STDA setting, we implement our method with two optimizers for domain adaptation and knowledge distillation. For our experiments on Office31 and ImageClef-DA, in order to have a fair comparison to TCP, our teacher model was chosen such that it's the same starting model as TCP and our student models have similar number of FLOPS or less than TCP pruned models. These FLOPS are measured for our models, as for the FLOPS of TCP, we report the number used in their paper. In addition, in order to show that's our framework can work with most domain adaptation and distillation techniques, we provide the result using combination of logits and feature distillation with MMD and RevGrad domain adaptation.

For the implementations of MTDA setting, we use the same number of optimizers as teacher models, with each optimizer responsible for the domain adaptation. Similar to our STDA setting, we have another optimizer for the knowledge distillation of the student. For this approach, we mainly use RevGrad and logits as domain adaptation and distillation technique. In order to have a fair comparison with state-of-the-art method [8], we decided to use LeNet as student backbone and ResNet50 as teachers. Since the authors use a non-standard backbone, with multiple residual layers and deconv layers, we decided to use a standard backbone (LeNet) that's close to their backbone while having less computational power for fair comparison. For the accuracy, we report the average accuracy of our model across all target domains.

For both KD-STDA and KD-MTDA, the details for our hyper-parameters can be found in the supplementary material. Both methods were implemented in PyTorch, and experiments were run on a server equipped with a Nvidia P100

GPU.

## 5. Results and Discussions

### 5.1. STDA with logits- and feature-based distillation:

Table 1 shows the accuracy of UDA methods on the Office31 dataset. Results indicate that our KD-STDA method outperforms all the baselines and state-of-the-art TCP [24]. As for the comparison with our baseline scenarios, the second baseline  $KD \rightarrow UDA$  has the best performance among all the scenarios and better than the state-of-the-art. The difference between the second and third baselines "UDA only" is because the compact model of the second baseline has distilled features which proves to be the better start than the third baseline. The first baseline  $UDA \rightarrow KD$  provides the lowest accuracy due to unsupervised knowledge distillation on the target domain without a consistency loss. In addition, our technique on both student models performs better than state-of-the-art TCP, especially on ResNet18, where our method improved the most upon state-of-the-art. Finally, our progressive method also seems to improve upon UDA since we perform better than the second baseline.

Table 2 shows the accuracy of methods on ImageClef, confirming the tendencies of Table 1. Our third baseline "UDA only" performs better than TCP in this dataset, which can explain why our methods performance better. Both Tables 1 and 2 show that our techniques perform well with our proposed feature-based KD method. Results follow our previous analysis concerning the different baselines. In addition, feature-based distillation can perform better than logits-based distillation, as was shown in [17].

Up until now, we have employed an discrepancy-based (MMD) method for STDA. For this experiment, we replace this method with an adversarial-based method. We choose the popular adversarial method with GRL called RevGrad [3] as since it has been used and adapted to other tasks like object detection, image segmentation. Table 3 shows the difference in accuracy between using RevGrad[3] and MMD[2]. Results indicate that RevGrad provides similar accu-

racy compared to MMD based method with both logits and feature distillation. In Figure 5, the t-SNE visualization [50], computed from the output of the feature extractor of each CNN, of the scenario  $A \rightarrow W$  on ResNet34, further confirms the similarity between both variants.

Table 1: Target domain accuracy of the proposed and baseline UDA methods on the Office31 dataset, when ResNet50 is the teacher, and ResNet34 and ResNet18 are the desired (student) CNN backbones. Results are shown for logits- and feature-based distillation: LB | FB.

Training methods	Source $\rightarrow$ Target							Average
	A $\rightarrow$ W	W $\rightarrow$ A	D $\rightarrow$ W	W $\rightarrow$ D	D $\rightarrow$ A	A $\rightarrow$ D		
Teacher: ResNet50 — Student: ResNet34								
Baseline 1: UDA $\rightarrow$ KD	25.4   9.6	7.1   6.3	28.5   12.6	50.0   13.8	9.7   5.3	30.7   11.2	25.2   9.8	
Baseline 2: KD $\rightarrow$ UDA	75.7   51.8	61.2   22.0	97.8   78.3	99.7   93.2	59.6   17.3	81.1   57.2	79.1   53.3	
Baseline 3: UDA only on ResNet34	67.2	52.3	93.6	96.6	52.2	71.6	72.2	
RevGrad	75.2	59	97.4	99.8	59.1	78.0	78.1	
CDAN	76.2	59.2	98.0	99.6	59.4	78.9	78.5	
TCP: prune rate = 12%	81.8	55.5	98.2	99.8	50.0	77.9	77.2	
KD-STDA MMD (Ours)	85.7   <b>86.0</b>	62.3   <b>67.6</b>	97.1   <b>99.0</b>	<b>100</b>   <b>100</b>	61.8   <b>66.4</b>	82.1   <b>84.7</b>	81.5   <b>83.9</b>	
Teacher: ResNet50 — Student: ResNet18								
Baseline 1: UDA $\rightarrow$ KD	28.8   10.2	5.8   5.3	33.7   11.5	51.1   12.7	7.8   7.8	25.1   9.2	25.3   9.4	
Baseline 2: KD $\rightarrow$ UDA	69.0   53.1	57.3   16.8	96.2   82.3	<b>100</b>   88.6	56.3   19.8	73.6   54.4	75.4   52.5	
Baseline 3: UDA only on ResNet18	60.2	49.2	93.7	97.7	47.6	66.4	69.1	
RevGrad	71.6	53.4	96.5	99.2	53.6	75.5	75.0	
CDAN	72.3	53.9	96.9	99.3	53.9	76.1	75.4	
TCP: prune rate = 45%	77.4	46.3	96.3	<b>100</b>	36.1	72.0	71.3	
KD-STDA MMD (Ours)	78.9   <b>87.3</b>	58.1   <b>64.1</b>	94.2   <b>99.0</b>	<b>100</b>   <b>100</b>	57.2   <b>63.3</b>	81.7   <b>85.7</b>	77.8   <b>83.2</b>	

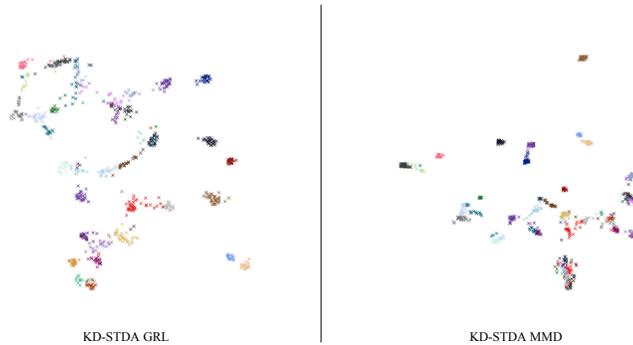


Figure 5: t-SNE visualisation of RevGrad and MMD with KD-STDA for scenario  $A \rightarrow W$  on ResNet34. Best viewed in color. We provide higher resolution figures in Supplementary Material.

Table 2: Same results as in Tables 1 but on the ImageClef-DA dataset.

Training methods	Source → Target						Average
	I → P	P → I	I → C	C → I	C → P	P → C	
Teacher: ResNet50 — Student: ResNet34							
Baseline 1: UDA → KD	48.8   12.2	41.0   14.3	46.0   12.2	39.6   11.5	38.8   9.4	39.0   8.7	42.0   11.4
Baseline 2: KD → UDA	76.6   70.2	87.3   77.3	92.0   88.2	80.0   73.2	65.6   60.6	90.3   84.2	81.9   75.6
Baseline 3: UDA only on ResNet34	73.3	86.3	<b>92.6</b>	79.3	65.8	87.5	80.8
RevGrad	75.9	87.0	92.2	79.3	65.3	90.0	81.6
CDAN	<b>76.8</b>	87.9	93.3	80.2	65.9	90.2	82.4
TCP: prune rate = 12%	75.0	82.6	92.5	80.8	66.2	86.5	80.6
KD-STDA MMD (Ours)	75.6   72.2	<b>89.0</b>   87.5	<b>92.6</b>   92.2	<b>83.8</b>   82.5	66.5   <b>66.8</b>	<b>92.8</b>   89.5	<b>83.3</b>   81.8
Teacher: ResNet50 — Student: ResNet18							
Baseline 1: UDA → KD	45.1   10.3	41.8   12.2	42.5   14.2	43.1   15.5	43.3   7.7	34.5   6.8	41.7   11.1
Baseline 2: KD → UDA	72.1   68.3	86.3   72.3	91.8   84.6	74.6   71.9	61.8   62.1	90.6   85.1	79.5   74.0
Baseline 3: UDA only on ResNet18	70.6	83.8	86.1	75.3	62.0	89.1	77.8
RevGrad	71.2	86.8	92.0	76.7	63.9	89.9	80.1
CDAN	72.1	87.1	<b>92.3</b>	77.8	63.7	89.5	80.4
TCP: prune rate = 45%	67.8	77.5	88.6	71.6	57.7	79.5	73.7
KD-STDA MMD (Ours)	73.1   <b>73.8</b>	88.0   <b>88.2</b>	92.1   91.5	77.3   <b>78.7</b>	<b>65.6</b>   64.4	<b>91.0</b>   89.9	<b>81.1</b>   <b>81.1</b>

Table 3: Same results as in Tables 1 but using MMD and GRL methods.

Training methods	Source → Target						Average
	A → W	W → A	D → W	W → D	D → A	A → D	
Teacher: ResNet50 — Student: ResNet34							
KD-STDA MMD	85.7   <b>86.0</b>	62.3   <b>67.6</b>	97.1   <b>99.0</b>	<b>100</b>   <b>100</b>	61.8   <b>66.4</b>	82.1   <b>84.7</b>	81.5   <b>83.9</b>
KD-STDA RevGrad	<b>82.6</b>   82.1	<b>64.4</b>   63.1	97.8   <b>98.7</b>	<b>100</b>   <b>100</b>	64.0   <b>64.2</b>	<b>86.0</b>   85.3	<b>82.5</b>   82.2

### 5.2. MTDA with logits-based distillation:

Tables 4 shows the accuracy of our generalized algorithm KD-MTDA for a multi-target setting on each different target domains and their average on Digits dataset. Results indicate that our method improves by 3% upon the current state-of-the-art [8] on  $sv \rightarrow mt, mm, up$ . As for the other scenario, we only perform slightly better, this is because we use the same hyper-parameter on different scenarios and on different teachers. One way of improving the current result would be to have a different set of hyper-parameters for each teacher such that it is best optimized for each combination of source and target. In addition, compared to a domain adaptation on a mixed target domains, MTDA focused algorithm perform much better overall.

From Table 5, where we show our accuracy on the PACS dataset, in contrast

Table 4: Target accuracy of the proposed and reference MTDA methods on the Digits dataset.

LeNet	Source → Targets							
	sv → mt, mm, up				mm → sv, mt, up			
	sv → mt	sv → mm	sv → up	Average	mm → sv	mm → mt	mm → up	Average
Teacher: ResNet50 — Student: LeNet								
Source Only	62.1	40.4	39.9	47.5	40.0	84.5	80.4	68.3
ADDA	77.7	64.2	64.1	68.7	40.6	92.8	80.7	71.4
RevGrad	73.8	61.0	62.5	65.8	51.8	61.0	85.3	66.0
MTDA-ITA[8]	84.6	<b>65.3</b>	70	73.3	53.5	<b>98.2</b>	<b>94.1</b>	81.9
KD-MTDA	<b>85.1</b>	63.9	<b>81.5</b>	<b>76.8</b>	<b>61.8</b>	97.8	86.4	<b>82.0</b>

Table 5: Same results as in Tables 4 but on PACS dataset.

LeNet	Source → Targets							
	P → Ap, Cr, S				Ap → Cr, S, P			
	P → Ap	P → Cr	P → S	Average	Ap → Cr	Ap → S	Ap → P	Average
Teacher: ResNet50 — Student: LeNet								
ADDA	24.3	20.1	22.4	22.3	17.8	18.9	32.8	23.2
MTDA-ITA [8]	<b>31.4</b>	23.0	28.2	27.6	27.0	28.9	<b>35.7</b>	30.5
KD-MTDA	24.6	<b>32.2</b>	<b>33.8</b>	<b>30.2</b>	<b>46.6</b>	<b>57.5</b>	<b>35.6</b>	<b>46.6</b>

to the previous comparison, our method perform significantly better than the current state-of-the-art [8]. The results from this dataset also confirm our hypothesis that by applying current STDA technique on multiple target domains, it does not yield good result and the model fail to generalize across multiple target domains. In addition, the paper of [8] use an architecture that includes several Residual block layers whereas, our architecture is based on LeNet5 which only consisted of 5 convolution layers, this meant that our method not only performs better than current method, it can also work with smaller models.

Table 6: Accuracy of KD-MTDA and DADA[40] by using Alexnet (student CNN) and Resnet50 (teacher CNNs) as backbones on the Office-Caltech dataset.

Models	A → C,D,W	C → A,D,W	D → A,C,W	W → A,C,D	Average
Teacher: ResNet50 — Student: AlexNet					
Source only	83.1	88.9	86.7	82.2	85.2
RevGrad[3]	85.9	90.5	88.6	90.4	88.9
DADA[40]	86.3	91.7	89.9	91.3	89.8
KD-MTDA	<b>93.3</b>	<b>93.9</b>	<b>90.1</b>	<b>91.2</b>	<b>92.1</b>

Results in Table6 show the accuracy CNNs trained with KD-MTDA on

Table 7: Accuracy of STDA baselines versus both versions of approach – KD-STDA (single teacher on mixed targets) and KD-MTDA (multiple teachers with multiple targets) on the Office31 dataset

Models	Source → Targets			Average
	A → D,W	D → A,W	W → A,D	
CAT STDA[51] (1 model/target)	78.5	62.9	<b>98.8</b>	80.1
GCAN STDA[52] (1 model/target)	<b>79.5</b>	<b>63.7</b>	98.4	<b>80.5</b>
KD-STDA RevGrad	75.3	64.0	67.0	68.8
KD-MTDA	<b>82.5</b>	<b>74.9</b>	<b>77.6</b>	<b>78.3</b>

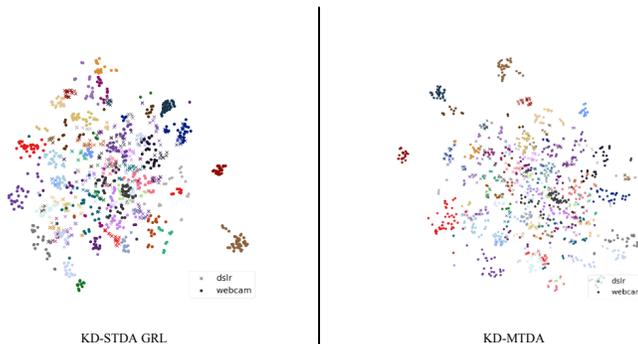


Figure 6: t-SNE of KD-STDA with mixed targets and KD-MTDA. Best viewed in color. We provide higher resolution figures in Supplementary Material.

OfficeCaltech dataset compared with a MTDA technique that does not rely on domain labels. The results in this Table further validate the importance of our method.

Our method is also compared for a single vs. multiple teachers using same hyper-parameters as with the STDA with AlexNet as student CNN backbone. From Table 7, multiple teachers always outperform the single teacher whether using the domain labels or not. In addition, Figure 6, shows that our generalized approach for MTDA separates the features better than our single teacher. Regarding comparisons with the state-of-the-art STDA model with 1 model/target, our MTDA model is capable of reaching almost the same average performance while encoding multiple targets domains inside the same model.

### 5.3. Complexity analysis:

While CNNs pruned using the TCP [24] method require fewer parameters than our student models, Table 8 indicates that our method achieve the same number of FLOPS on ResNet34, and with even fewer FLOPs on ResNet18. This means that while TCP prunes more parameters, it may not have much impact on the number of FLOPS since the pruned filters are ranked and pruned globally across the network, instead of being pruned at each layer. TCP prunes away filters that do not impact the FLOPS but can impact performance. Another important point of having more parameters is that, over-parametrization can help provide better generalization – our student CNNs have a greater chance to provide better generalization than a pruned model with fewer parameters.

Table 8: Computational complexity of proposed methods and TCP

Models	no. operations (GFLOPS)	no. parameters (M)	
		Office31	ImageClef
ResNet50	4.1	25.5	25.5
TCP(12% Pruned)	3.6	15.8	15.9
ResNet34	3.6	21.7	21.7
TCP(46% Pruned)	2.2	10.6	10.9
ResNet18	1.8	11.1	11.1

As for the complexity comparison of an STDA based approach adapted to MTDA setting, i.e. having one model for each target domain instead of having one model for all the target domains. If we assume that there are  $N$  target domains and  $M$  the memory space occupied by a backbone, this would mean that in the scenario of having one model per target domains, we would need  $N \times M$  and in the worst case  $N$  GPUs if it's a complex backbone. This clearly shows that having one model per target domain is not scalable and it is much preferable having one model to handle multiple target domains.

## 6. Conclusion

In this paper, we proposed a joint optimization of KD and UDA that tackles both the problem of domain shift and model compression in both the STDA and MTDA setting. In addition, our method also works with different UDA and KD

techniques, whether it is logits or feature-based. Our results on STDA suggest that our method is capable of adapting and accelerating a model by alleviating the domain shift problem and reduce complexity. The proposed technique is also capable of adapting a model to multiple target domains while keeping high accuracy. In both settings, STDA and MTDA, our method outperforms the current state-of-the-art, especially on compact models. Since UDA is an active area of research, our future research can include a more efficient target domain knowledge transfer method for better compression or better fusion method for combining different target domains.

### Acknowledgements

This research was partially supported by the Mathematics of Information Technology and Complex Systems (MITACS) and the Natural Sciences and Engineering Research Council of Canada (NSERC) organizations.

### References

- [1] D. Mekhazni, A. Bhuiyan, G. Ekladios, E. Granger, Unsupervised domain adaptation in the dissimilarity space for person re-identification, in: ECCV 2020.
- [2] M. Long, Y. Cao, J. Wang, M. I. Jordan, Learning transferable features with deep adaptation networks, in: ICML 2015.
- [3] Y. Ganin, V. Lempitsky, Unsupervised domain adaptation by backpropagation, in: PRLM 2015.
- [4] E. Tzeng, J. Hoffman, K. Saenko, T. Darrell, Adversarial discriminative domain adaptation, CoRR abs/1702.05464.
- [5] J. Shen, Y. Qu, W. Zhang, Y. Yu, Wasserstein distance guided representation learning for domain adaptation (2017). [arXiv:1707.01217](https://arxiv.org/abs/1707.01217).

- [6] H. Fu, M. Gong, C. Wang, K. Batmanghelich, K. Zhang, D. Tao, Geometry-consistent generative adversarial networks for one-sided unsupervised domain mapping, in: CVPR 2019.
- [7] S. Benaim, L. Wolf, One-sided unsupervised domain mapping, in: NIPS 2017.
- [8] B. Gholami, P. Sahu, O. Rudovic, K. Bousmalis, V. Pavlovic, Unsupervised multi-target domain adaptation: An information theoretic approach, CoRR abs/1810.11547.
- [9] Z. Chen, J. Zhuang, X. Liang, L. Lin, Blending-target domain adaptation by adversarial meta-adaptation networks, CoRR abs/1907.03389.
- [10] R. Zhao, Y. Hu, J. Dotzel, C. De Sa, Z. Zhang, Improving neural network quantization without retraining using outlier channel splitting, in: PMLR 2019.
- [11] S. Han, H. Mao, W. J. Dally, Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding, CoRR abs/1510.00149.
- [12] P. Stock, A. Joulin, R. Gribonval, B. Graham, H. Jégou, And the bit goes down: Revisiting the quantization of neural networks, in: Submitted to ICLR 2020.
- [13] M. Jaderberg, A. Vedaldi, A. Zisserman, Speeding up convolutional neural networks with low rank expansions, CoRR abs/1405.3866.
- [14] W. Wen, C. Xu, C. Wu, Y. Wang, Y. Chen, H. Li, Coordinating filters for faster deep neural networks, CoRR abs/1703.09746.
- [15] G. Hinton, O. Vinyals, J. Dean, Distilling the knowledge in a neural network, in: NIPS Deep Learning and Representation Learning Workshop, 2015.

- [16] W. Park, D. Kim, Y. Lu, M. Cho, Relational knowledge distillation, in: CVPR 2019.
- [17] B. Heo, J. Kim, S. Yun, H. Park, N. Kwak, J. Y. Choi, A comprehensive overhaul of feature distillation, in: ICCV 2019.
- [18] S. Mirzadeh, M. Farajtabar, A. Li, H. Ghasemzadeh, Improved knowledge distillation via teacher assistant: Bridging the gap between student and teacher, CoRR abs/1902.03393.
- [19] J. Frankle, M. Carbin, The lottery ticket hypothesis: Finding sparse, trainable neural networks, in: ICLR 2019.
- [20] P. Molchanov, A. Mallya, S. Tyree, I. Frosio, J. Kautz, Importance estimation for neural network pruning, in: CVPR 2019.
- [21] X. Dong, Y. Yang, Network pruning via transformable architecture search, in: NIPS 2019.
- [22] Z. Liu, M. Sun, T. Zhou, G. Huang, T. Darrell, Rethinking the value of network pruning, in: ICLR 2019.
- [23] Y. He, P. Liu, Z. Wang, Z. Hu, Y. Yang, Filter pruning via geometric median for deep convolutional neural networks acceleration, in: CVPR 2019.
- [24] C. Yu, J. Wang, Y. Chen, Z. Wu, Accelerating deep unsupervised domain adaptation with transfer channel pruning, CoRR abs/1904.02654.
- [25] S. Ruder, P. Ghaffari, J. G. Breslin, Knowledge adaptation: Teaching to adapt, CoRR abs/1702.02052.
- [26] M. Orbes-Arteaga, J. Cardoso, L. Sørensen, C. Igel, S. Ourselin, M. Modat, M. Nielsen, A. Pai, Knowledge distillation for semi-supervised domain adaptation (2019). [arXiv:1908.07355](https://arxiv.org/abs/1908.07355).

- [27] Z. Allen-Zhu, Y. Li, Y. Liang, Learning and generalization in overparameterized neural networks, going beyond two layers, in: NIPS 2019.
- [28] L. T. Nguyen-Meidine, E. Granger, M. Kiran, J. Dolz, L.-A. Blais-Morin, Joint progressive knowledge distillation and unsupervised domain adaptation, IJCNN 2020.
- [29] Q. Tian, T. Arbel, J. J. Clark, Structured deep fisher pruning for efficient facial trait classification, *Image and Vision Computing* 77 (2018) 45 – 59. doi:<https://doi.org/10.1016/j.imavis.2018.06.008>.  
URL <http://www.sciencedirect.com/science/article/pii/S0262885618301045>
- [30] P. Singh, V. S. R. Kadi, V. P. Namboodiri, Falf convnets: Fatuous auxiliary loss based filter-pruning for efficient deep cnns, *Image and Vision Computing* 93 (2020) 103857. doi:<https://doi.org/10.1016/j.imavis.2019.103857>.  
URL <http://www.sciencedirect.com/science/article/pii/S0262885619304500>
- [31] Y. Zhang, N. Wang, S. Cai, Adversarial sliced wasserstein domain adaptation networks, *Image and Vision Computing* 102 (2020) 103974. doi:<https://doi.org/10.1016/j.imavis.2020.103974>.  
URL <http://www.sciencedirect.com/science/article/pii/S0262885620301062>
- [32] K. Bousmalis, N. Silberman, D. Dohan, D. Erhan, D. Krishnan, Unsupervised pixel-level domain adaptation with generative adversarial networks, *CoRR* abs/1612.05424.
- [33] W. Hong, Z. Wang, M. Yang, J. Yuan, Conditional generative adversarial network for structured domain adaptation, in: CVPR 2018.
- [34] M. Toldo, U. Michieli, G. Agresti, P. Zanuttigh, Unsupervised domain adaptation for mobile semantic segmentation based on cycle consistency

and feature alignment, *Image and Vision Computing* 95 (2020) 103889.  
doi:<https://doi.org/10.1016/j.imavis.2020.103889>.

URL <http://www.sciencedirect.com/science/article/pii/S0262885620300214>

- [35] K. Saito, Y. Ushiku, T. Harada, Asymmetric tri-training for unsupervised domain adaptation, in: PMLR 2017.
- [36] Y. Li, N. Wang, J. Shi, J. Liu, X. Hou, Revisiting batch normalization for practical domain adaptation, CoRR abs/1603.04779.
- [37] K.-Y. Wei, C.-T. Hsu, Generative adversarial guided learning for domain adaptation, in: BMVC, 2018.
- [38] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, T. Darrell, CyCADA: Cycle-consistent adversarial domain adaptation, in: PMLR 2018.
- [39] Z. Liu, Z. Miao, X. Pan, X. Zhan, D. Lin, S. X. Yu, B. Gong, Open compound domain adaptation, in: ICPR 2020.
- [40] X. Peng, Z. Huang, X. Sun, K. Saenko, Domain agnostic learning with disentangled representations, in: K. Chaudhuri, R. Salakhutdinov (Eds.), *Proceedings of the 36th International Conference on Machine Learning*, Vol. 97 of *Proceedings of Machine Learning Research*, PMLR, Long Beach, California, USA, 2019, pp. 5102–5112.  
URL <http://proceedings.mlr.press/v97/peng19b.html>
- [41] S. Ruder, P. Ghaffari, J. G. Breslin, Knowledge adaptation: Teaching to adapt, CoRR abs/1702.02052.
- [42] M. Orbes-Arteainst, M. J. Cardoso, L. Sørensen, C. Igel, S. Ourselin, M. Modat, M. Nielsen, A. Pai, Knowledge Distillation for Semi-supervised Domain Adaptation, 2019, pp. 68–76. doi:10.1007/978-3-030-32695-1\_8.

- [43] Y. LeCun, C. Cortes, MNIST handwritten digit database [cited 2016-01-14 14:24:11].  
URL <http://yann.lecun.com/exdb/mnist/>
- [44] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, V. Lempitsky, Domain-adversarial training of neural networks (2016). [arXiv:1505.07818](https://arxiv.org/abs/1505.07818).
- [45] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, A. Y. Ng, Reading digits in natural images with unsupervised feature learning, in: NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011, 2011.  
URL [http://ufldl.stanford.edu/housenumbers/nips2011\\_housenumbers.pdf](http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf)
- [46] J. J. Hull, A database for handwritten text recognition research, IEEE Transactions on Pattern Analysis and Machine Intelligence 16 (5) (1994) 550–554. doi:10.1109/34.291440.
- [47] K. Saenko, B. Kulis, M. Fritz, T. Darrell, Adapting visual category models to new domains, in: Proceedings of the 11th European Conference on Computer Vision: Part IV, ECCV'10, Springer-Verlag, 2010, p. 213–226.
- [48] D. Li, Y. Yang, Y.-Z. Song, T. Hospedales, Deeper, broader and artier domain generalization, in: International Conference on Computer Vision, 2017.
- [49] M. Long, Z. CAO, J. Wang, M. I. Jordan, Conditional adversarial domain adaptation, in: S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, R. Garnett (Eds.), Advances in Neural Information Processing Systems, Vol. 31, Curran Associates, Inc., 2018, pp. 1640–1650.  
URL <https://proceedings.neurips.cc/paper/2018/file/ab88b15733f543179858600245108dd8-Paper.pdf>

- [50] L. van der Maaten, G. Hinton, Visualizing data using t-SNE, *Journal of Machine Learning Research* 9 (2008) 2579–2605.  
URL <http://www.jmlr.org/papers/v9/vandermaaten08a.html>
- [51] Z. Deng, Y. Luo, J. Zhu, Cluster alignment with a teacher for unsupervised domain adaptation, in: *The IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [52] X. Ma, T. Zhang, C. Xu, Gcan: Graph convolutional adversarial network for unsupervised domain adaptation, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.

# Supplementary Material

## Appendix A. Hyper-parameters

Table A.9: Hyper-parameters for our algorithms for each backbone and dataset

Hyper parameters	Office31 Resnet34-18	ImageClef-DA ResNet34-18	Digits LeNet	Pacs LeNet
$N_e$	400	400	100	100
$\tau$	20	20	20	20
$\alpha$	0.8	0.8	0.5	0.5
$s$	0.1	0.1	0.1	0.1
$f$	0.8	0.8	0.5	0.5
$\gamma$	-	-	0.5	0.5
UDA Learning Rate	0.001	0.0001	0.01	0.01
KD Learning Rate	0.001	0.001	0.01	0.01
weight decay	0.0005	0.0005	0.0005	0.0005

From Table A.9, we can find the details of our hyper-parameters for both STDA and MTDA setting on different dataset. These hyper-parameters were selected using a standard cross-validation process.

## Appendix B. Additional results

### Appendix B.1. Logits distillation vs. feature distillation for MTDA

In this study, we compare the difference between logits distillation and feature distillation in MTDA setting on Office31 dataset. The Table B.10 shows the accuracy of our MTDA method with either logits or feature distillation on student model.

Table B.10: Accuracy of proposed method using either logits or feature distillation

Types of distillation	Source $\rightarrow$ Targets			Average
	A $\rightarrow$ D,W	D $\rightarrow$ A, W	W $\rightarrow$ A, D	
<b>Teacher: ResNet50 — Student: ResNet18</b>				
KD-MTDA logits distillation	80.8	79.7	78.6	79.7
KD-MTDA feature distillation	<b>81.2</b>	<b>80.5</b>	<b>79.5</b>	<b>80.4</b>

From Table B.10, where we compare the performance of logits-based and feature-based distillation on the setting of MTDA, in contrast to with feature

distillation for STDA, feature distillation in MTDA does not provide a significant improvement in performance. Results indicate that feature distillation in MTDA only perform slightly better than logits distillation.

### Appendix C. t-SNE

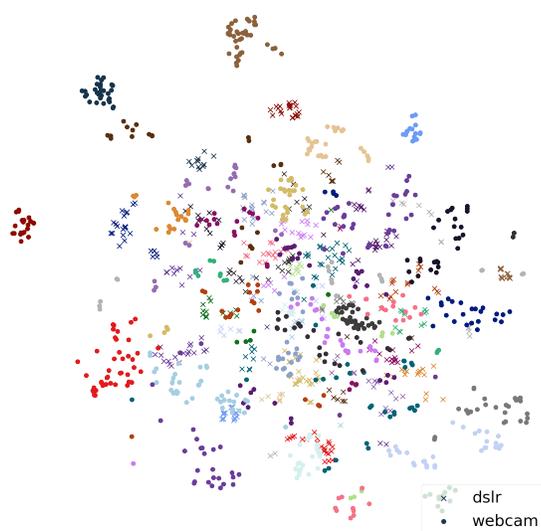


Figure C.7: t-SNE of KD-MTDA. Best viewed in color.

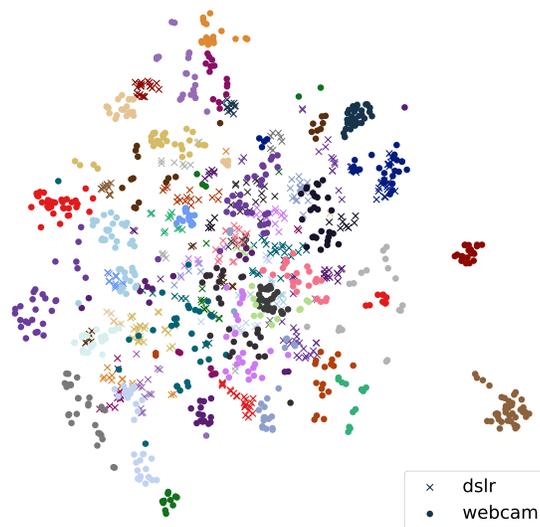


Figure C.8: t-SNE of KD-STDA with mixed targets. Best viewed in color.



Figure C.9: t-SNE of KD-STDA MMD. Best viewed in color.

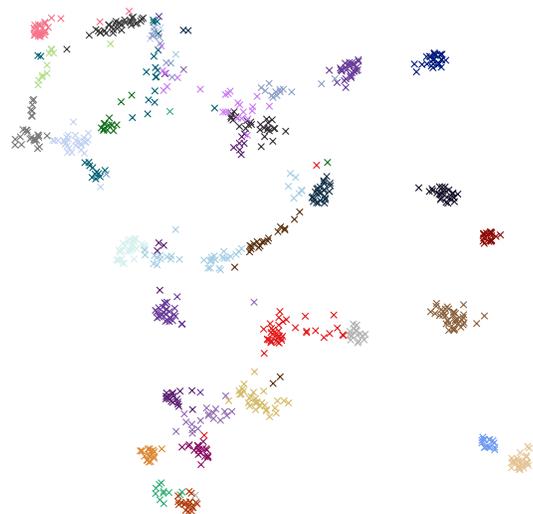


Figure C.10: t-SNE of KD-STDA GRL. Best viewed in color.