# Source-Free Domain Adaptation with Frozen Multimodal Foundation Model

Song Tang[1,2,3], Wenxin Su[1], Mao Ye[*4], and Xiatian Zhu[*5]

[1]University of Shanghai for Science and Technology [2]Universität Hamburg [3]ComOriginMat Inc.
[4]University of Electronic Science and Technology of China [5]University of Surrey

tangs@usst.edu.cn, {suwenxin43, cvlab.uestc}@gmail.com , xiatian.zhu@surrey.ac.uk

## Abstract

*Source-Free Domain Adaptation (SFDA) aims to adapt a source model for a target domain, with only access to unlabeled target training data and the source model pretrained on a supervised source domain. Relying on pseudo labeling and/or auxiliary supervision, conventional methods are inevitably error-prone. To mitigate this limitation, in this work we for the first time explore the potentials of off-the-shelf vision-language (ViL) multimodal models (e.g., CLIP) with rich whilst heterogeneous knowledge. We find that directly applying the ViL model to the target domain in a zero-shot fashion is unsatisfactory, as it is not specialized for this particular task but largely generic. To make it task specific, we propose a novel **D**istilling mult**I**modal **F**oundation m**O**del (**DIFO**) approach. Specifically, DIFO alternates between two steps during adaptation: (i) Customizing the ViL model by maximizing the mutual information with the target model in a prompt learning manner, (ii) Distilling the knowledge of this customized ViL model to the target model. For more fine-grained and reliable distillation, we further introduce two effective regularization terms, namely most-likely category encouragement and predictive consistency. Extensive experiments show that DIFO significantly outperforms the state-of-the-art alternatives. Code is here.*

## 1. Introduction

Unsupervised Domain Adaptation (UDA) relies on both well-annotated source data and unannotated target data. However, due to heightened safety and privacy concerns, accessing source data freely has become difficult [19, 24]. In response, Source-Free Domain Adaptation (SFDA) has gained attention as a more practical solution, aiming to transfer a pretrained source model to the target domain using only unlabeled target data.
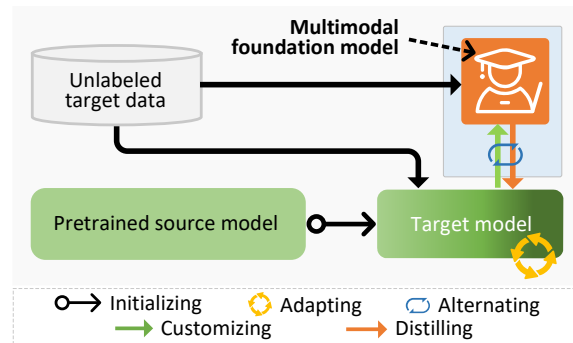


Figure 1. We expand beyond traditional SFDA methods that rely solely on a pretrained source model and unlabeled target data. Instead, we innovate by exploring off-the-shelf multimodal foundation models, such as CLIP, in an unsupervised manner (marked by the box with blue background).

Due to the absence of source samples, traditional distribution matching approaches are no longer viable [8, 15]. The predominant alternative is self-supervised learning, which generates or mines auxiliary information to facilitate unsupervised adaptation. Two main approaches exist: constructing a pseudo source domain to leverage established UDA methods such as adversarial learning [17, 47] or domain shift minimization based on distribution measurement [5, 16, 43] and mining extra supervision from the source model [12, 19, 46] or target data [40, 49, 50]. In the presence of domain distribution shift, applying the source model to the target domain introduces inevitable errors in pseudo-labeling or auxiliary supervision, thereby limiting adaptation performance.

To address identified limitations, we pioneer the exploration of off-the-shelf multimodal foundation models, such as the vision-language (ViL) model CLIP [24], transcending the constraints of both the source model and target data knowledge. However, direct application of the ViL model proves unsatisfactory, lacking specialization for specific tasks. To overcome this, we propose a novel task-specific distillation approach named **D**istilling mult**I**modal

---
*Corresponding author

*Foundation mOdel (DIFO)*. Initially, we customize the ViL model through *unsupervised* prompt learning for imposing task-specific information. Subsequently, we distil the knowledge from this customized ViL model to the target model, with joint supervision through two designed regularization terms: (1) most-likely category encouragement for coarse-grained distillation and (2) predictive consistency for fine-grained distillation.

Our **contributions** are summarized as follows. **(1)** Pioneering the use of generic but heterogeneous knowledge sources (e.g., the off-the-shelf ViL model) for the SFDA problem, transcending the limited knowledge boundary of a pretrained source model and unlabeled target data. **(2)** Development of the novel DIFO approach to effectively distill useful task-specific knowledge from the general-purpose ViL model. **(3)** Extensive evaluation on standard benchmarks, demonstrating the significant superiority of our DIFO over previous state-of-the-art alternatives under conventional closed-set settings, as well as more challenging partial-set and open-set settings.

## 2. Related Work

**Source-free domain adaptation.** Existing SFDA approaches fall into three distinct categories. The first explicitly aligns the pseudo source domain with the target domain, treating SFDA as a specialized case of unsupervised domain adaptation. This alignment is achieved by constructing the pseudo source domain through a generative model [22, 44] or by splitting the target domain based on prior source hypotheses [6].

The second group extracts cross-domain factors from the source domain and transfers them in successive model adaptation for aligning feature distributions across the two domains. For example, [37] establishes a mapping relationship from a sample and its exemplar Support Vector Machine (SVM) (an individual classifier) on the source domain to ensure individual classification on the target domain. Some approaches leverage pre-trained source models to generate auxiliary factors, such as multi-hypothesis [19], prototypes [42], source distribution estimation [5], or hard samples [21] to aid in feature alignment.

The third group incorporates auxiliary information refined from the unlabeled target domain. In addition to widely used pseudo-labels [3, 25], geometry information, such as intrinsic neighborhood structure [39] and target data manifold [40], has also been exploited.

Despite continual advancements, these methods are limited by the knowledge derived solely from the pretrained source model and unlabeled target data. We break this limitation by tapping into the rich knowledge encoded in off-the-shelf multimodal foundation models.

**Large multimodal model.** Multimodal vision-language (ViL) models, such as CLIP [31] and ALIGN [14], have shown promise across various mono-modal and multimodal tasks by capturing modality-invariant features. Approaches in this domain can be broadly categorized into two lines.

The first line focuses on enhancing ViL model performance. For instance, in [9, 53], prompt learning optimizes the text encoder through the use of tailored, learnable prompts designed for specific scenarios. Other efforts aim to improve data efficiency by repurposing noisy data [1].

The second line utilizes ViL models as external knowledge to enhance downstream tasks, as demonstrated in this paper. Previous work in knowledge transfer primarily falls into two frameworks. For the first scheme, where the ViL model is directly applied to the target task in a zero-shot fashion [23], domain generality is leveraged without task-specific refinement. The second scheme does not focus on source model adaptation. Instead, it fine-tunes the ViL model to the target domain through prompt or adaptor learning with an amount of manal labels [4].

A relevant method to our DIFO is the UDA method DAPL [9]. Although both adopt CLIP, they differ significantly in problem setting and methodology. DAPL employs CLIP to learn domain-specific prompts, aiming to disentangle domain and category information in CLIP's visual features. In contrast, DIFO aligns target features to a progressively customized vision-language latent space in a memory-aware fashion. Importantly, DAPL requires labeled source data, making it inapplicable in SFDA.

## 3. Methodology

**Problem statement.** In the context of two distinct yet inter-related domains—namely, the labeled source domain and the unlabeled target domain—both characterized by the same set of $C$ categories, the following notation is employed. The source samples and their corresponding labels are represented as $\mathcal{X}_s$ and $\mathcal{Y}_s$ respectively. Similarly, the target samples and their true labels are denoted as $\mathcal{X}_t = \{x_i\}_{i=1}^n$ and $\mathcal{Y}_t = \{y_i\}_{i=1}^n$, where $n$ signifies the number of samples.

We aim to learn a target model $\theta_t : \mathcal{X}_t \to \mathcal{Y}_t$. This involves utilizing (1) a pre-trained source model $\theta_s : \mathcal{X}_s \to \mathcal{Y}_s$, (2) unlabeled target data, and (3) a Visual-Language (ViL) model denoted as $\theta_v$.

**Overview.** As depicted in Fig. 2, the proposed DIFO framework alternates between two distinct steps to customize and distill the off-the-shelf ViL knowledge.

*In the first step*, we engage in prompt learning on the ViL model for the purpose of task-specific customization. This serves to mitigate the guidance error within the ViL model. In particular, we adopt a mutual information-based alignment approach. This approach is characterized by its richness in context and interaction between the target model and the ViL model, as opposed to placing blind trust in either model alone as conventional methods.
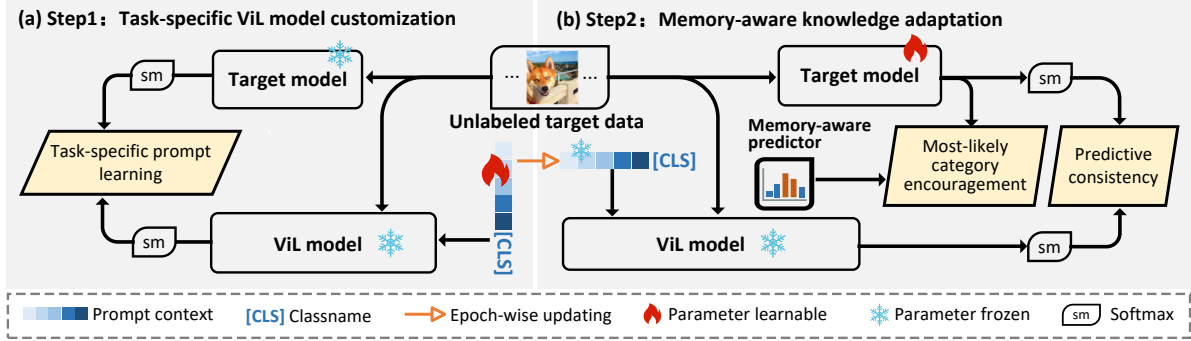
Figure 2. Overview of our DIFO: The process involves two alternating steps. First, we perform (a) *task-specific customization* of a ViL model through task-specific prompt learning ($L_{\text{TSC}}$). This is achieved under soft predictive guidance using mutual information maximization. Second, we undertake (b) *memory-aware knowledge adaptation*, incorporating two regularizations: most-likely category encouragement ($L_{\text{MCE}}$) predicted by our dynamic memory-aware predictor, along with the tupical predictive consistency ($L_{\text{PC}}$). These regularizations are designed to facilitate a coarse-to-fine adaptation.

*In the second step*, knowledge adaptation takes place within a unique constraint that encourages the identification of the most probable category labels in the logit space, while concurrently maintaining the typical predictive consistency. The most likely category labels are determined by a carefully designed memory-aware predictor, which dynamically integrates knowledge from both the target model and the ViL model in a cumulative fashion.

### 3.1. Task-Specific ViL Model Customization

We adopt the prompt learning framework for ViL model customization, with all the parameters of the ViL model frozen throughout. The only learnable part in customization is the prompts each assigned for a specific class. To optimize these prompts, we need a useful supervision. In SFDA, however, it is challenging to customize such a domain-generic ViL model towards to the target domain, at the absence of a well-trained target domain model. This is because, none of them can reasonably make predictions. That means there is no clearly good supervision signals available.

To address this challenge, we propose to explore the wisdom of the crowd by leveraging their predictive interaction as the supervision. Formally, we denote the predictions by the target model and the ViL model as $\theta_t(\boldsymbol{x}_i)$ and $\theta_v(\boldsymbol{x}_i)$, respectively, given an unlabeled target sample $\boldsymbol{x}_k$. We conduct the customization by maximizing the mutual information of their predictions as:

$$L_{\text{TSC}} = -\min_{\boldsymbol{v}} \mathbb{E}_{\boldsymbol{x}_i \in \mathcal{X}_t} \mathrm{I}\left(\theta_t\left(\boldsymbol{x}_i\right), \theta_v\left(\boldsymbol{x}_i, \boldsymbol{v}\right)\right) \quad (1)$$

where $\boldsymbol{v}$ is the prompt context to be learned and the function $\mathrm{I}(\cdot, \cdot)$ measures the mutual information [13].

This alignment design differs significantly from the conventional adoption of the Kullback–Leibler (KL) divergence. First of all, the mutual information is a lower optimization bound than KL divergence, facilitating deeper alignment (see Theorem 6 with the proof provided in `Supplementary`).

**Theorem 1** *Given two random variables $X, Y$. Their mutual information* $\mathrm{I}(X, Y)$ *and KL divergence* $D_{\text{KL}}(X||Y)$ *satisfy the unequal relationship as follows.*

$$-\mathrm{I}(X, Y) \leq D_{\text{KL}}(X, Y). \quad (2)$$

Crucially, the KL divergence exhibits an inherent bias towards a specific prediction, making it less suitable for our context where none of the predictions holds a significant advantage. On the contrary, mutual information considers the joint distribution or correlation between the two predictions. This distinction arises from their respective definitions: $-\mathrm{I}(X, Y) = -H(X) + H(X|Y)$ and $D_{\text{KL}}(X, Y) = -H(X) + H(X : Y)$, where

$$\begin{aligned} H(X \mid Y) &= -\sum p(\boldsymbol{x}, \boldsymbol{y}) \log p(\boldsymbol{x}|\boldsymbol{y}) \\ H(X : Y) &= -\sum p(\boldsymbol{x}) \log p(\boldsymbol{y}). \end{aligned} \quad (3)$$

The conditional entropy component $H(X|Y)$ of mutual information explicitly captures the joint distributions, a feature absent in KL divergence. Empirically, we also confirm the significance of incorporating this joint distribution-based interaction between the two predictions during the customization of the ViL model (see ablation study in `Tab. 6` and task-specific knowledge adaptation analysis in `Section 4.3`).

### 3.2. Memory-Aware Knowledge Adaptation

As previously mentioned, even with customization for the target domain, the ViL model may not be fully adapted due to no robust target model available in prior. This limitation hinders effective knowledge adaptation at this stage. To address this issue, we propose the incorporation of a specialized memory-aware predictor to provide additional learning guidance – most-likely category encouragement, complementing the conventional predictive consistency constraint.

**Most-likely category encouragement.** The rationale behind incorporating this learning constraint is to harness the collective knowledge of both the target model and the ViL model in order to enhance the discernment of probable category labels for each sample. Given the sluggish nature of this search process, it has been devised to function as a form of learning regularization. An illustration of this regularization process is presented in Fig. 3. Specifically, it is realized through two distinct steps as detailed below.

*(I) Memory-aware predictor.* We initiate the process by generating pseudo-labels that represent the most likely category distribution, utilizing historical information stored in a prediction bank. The prediction bank archives two types of historical data for all samples in the target domain: (1) predictions from the target model denoted by $\{\boldsymbol{p}_i\}_{i=1}^n$ and (2) predictions from the ViL model denoted by $\{\boldsymbol{p}_i'\}_{i=1}^n$.

Throughout the adaptation process, the predictions from the target model are updated iteratively. At the end of each training iteration, the newly predicted labels for the training batch from the target model replace their counterparts in the prediction bank. In contrast, predictions from the ViL model are updated collectively in an epoch-wise manner, triggering updates every $M$ iterations. This mixed-update strategy is designed to strike a balance between maintaining the stability of the customized ViL model's guidance and capturing the task-specific dynamics inherent in the adaptation process.

Based on the provided prediction bank, the creation of a pseudo-label for the most probable category involves a historical prediction fusion process as:

$$\bar{\boldsymbol{p}}_i = \omega\, \boldsymbol{p}_i + (1-\omega)\, \boldsymbol{p}_i'. \qquad (4)$$

Here, the weight $\omega$, drawn from an Exponential distribution with parameter $\lambda$, is a crucial factor. This fusion introduces dynamic bias rectification (represented by $\boldsymbol{p}_i$) based on the guidance from the customized ViL model ($\boldsymbol{p}_i'$). The role of $\boldsymbol{p}_i$ is to provide adjustments, leading us to adopt an asymmetric random weighting approach represented by $\omega$.

*(II) Category attention calibration.* Subsequently, we formulate a regularization technique employing pseudo-labels acquired through category attention calibration. Specifically, we begin by identifying the *top-N* most probable categories using $\bar{\boldsymbol{p}}_i$. The indices of these identified categories are denoted by $\mathcal{M}_i = \{m_k\}_{k=1}^N$. With $\mathcal{M}_i$, the target model's logit of a target domain sample $x_i$, denoted as $\boldsymbol{l}_i$, is segregated into positive and negative category groups. We define this regularization as:

$$L_{\text{MCE}} = \min_{\theta t} \mathbb{E}_{\boldsymbol{x}_i \in \mathcal{X}_t} \log \frac{\exp\left(a_i/\tau\right)}{\sum_{j \neq \mathcal{M}_i} \exp\left(b_i \cdot \boldsymbol{l}_{i,j}/\tau\right)}$$
$$a_i = \prod_{k=1}^N \boldsymbol{l}_{i,m_k}, \quad b_i = \sum_{k=1}^N \boldsymbol{l}_{i,m_k} \qquad (5)$$
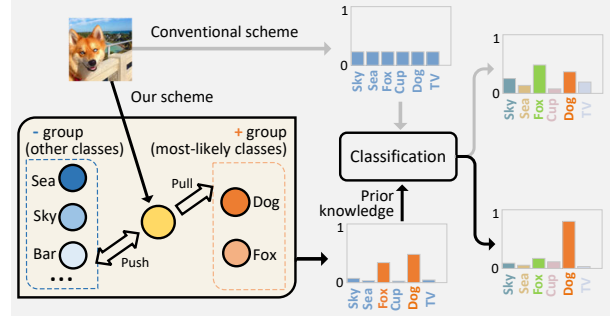


Figure 3. Illustration of most-likely category encouragement. In contrast to the conventional approach that assigns equal importance to all categories (depicted by the gray line), our approach (represented by the black line) introduces additional supervision by incorporating extra knowledge about the two most likely categories.

where $\boldsymbol{l}_{i,a}$ denotes the $a$-th element of $\boldsymbol{l}_i$ and $\tau$ is the temperature parameter.

In Eq. (5), we note that the product operation with $a_i$ in the numerator amplifies penalties for the probability decrease on the most likely categories compared to the sum form. Similarly, the sum with $b_i$ in the denominator serves as an increasing weighting parameter to enhance suppression of values at other locations. Moreover, $a_i$ is more sensitive to changes than $b_i$ due to $\frac{\partial a_i}{\partial m_k} \propto O(n^{N-1})$ and $\frac{\partial b_i}{\partial m_k} \propto O(1)$. By combining the use of $a_i$ and $b_i$, we globally impose a calibration effect on the elements corresponding to the most likely categories within the logit $\boldsymbol{l}_i$. Essentially, attention is introduced to these potential categories, as illustrated in the box with a yellow background in Fig. 3.

**Predictive Consistency.** For the purpose of knowledge adaptation, we incorporate the conventional predictive consistency loss as:

$$L_{\text{PC}} = \min_{\theta_t}\left[-\mathbb{E}_{\boldsymbol{x}_i \in \mathcal{X}_t}\text{I}\left(\theta_t\left(\boldsymbol{x}_i\right), \theta_v\left(\boldsymbol{x}_i, \boldsymbol{v}*\right)\right) + \alpha L_B\right], \quad (6)$$

where $\theta_t(\boldsymbol{x}_i)$ represents the target prediction, $\theta_v(\boldsymbol{x}_i, \boldsymbol{v})$ denotes the ViL prediction, and $\boldsymbol{v}$ is the prompt context learned during the initial phase of task-specific customization. The function $\text{I}(\cdot, \cdot)$ corresponds to the mutual information function. The parameter $\alpha$ serves as a trade-off parameter, and the category balance term $L_B = \text{KL}(\bar{\boldsymbol{q}}||\frac{1}{C})$ aligns with previous approaches [41, 49], preventing solution collapse by ensuring the empirical label distribution $\bar{\boldsymbol{q}}$ matches the uniform distribution $\frac{1}{C}$. For the reasons elaborated in `Section 3.1`, we employ mutual information for alignment.

### 3.3. Model training

To systematically distill and leverage task-specific knowledge from the ViL model, we adopt an epoch-wise training approach for DIFO. The training process is divided into $T$ epochs, each comprising two stages aligned with the two

**Algorithm 1** Training of DIFO

**Input**: Pre-trained source model $\theta_s$, target model $\theta_t$, ViL model $\theta_v$, unlabelled target domain $\mathcal{X}_t$, learnable prompt context $v$, #epoch $T$, #iteration per epoch $M$.
**Output**: The adapted target model $\theta_t$.
**Procedure**:

1: **Initialisation**: Set $\theta_t = \theta_s$ and $v$='a photo of a [CLS].'
2: **for** $t = 1:T$ **do**
3:     Update ViL predictions in the prediction bank.
4:     =========== *Step1* ===========
5:     **for** $m = 1:M$ **do**
6:         Sample a batch from $\mathcal{X}_t$;
7:         Forward prompt $v$ and this batch $\mathcal{X}_t^b$ through $\theta_v$;
8:         Forward this batch data through $\theta_t$;
9:         Customize $\theta_v$ by optimizing $L_{\text{TSC}}$ (Eq. (1)) and obtain task-specific prompt context $v^*$.
10:    **end for**
11:    =========== *Step2* ===========
12:    **for** $m = 1:M$ **do**
13:        Sample a batch from $\mathcal{X}_t$;
14:        Forward the $v^*$ and this batch through $\theta_v$;
15:        Forward this batch data through $\theta_t$;
16:        Discover most-likely category (Eq. (4));
17:        Update model $\theta_t$ by optimizing $L_{\text{MKA}}$ (Eq. (7)).
18:        Update target predictions in the prediction bank.
19:    **end for**
20:    Set $v = v^*$.
21: **end for**
22: **return** Adapted model $\theta_t$.

steps in the DIFO framework (Fig. 2). During the first stage, training is governed by the objective $L_{\text{TSC}}$, and in the subsequent second stage, the objective function transitions to

$$L_{\text{MKA}} = L_{\text{PC}} + \beta L_{\text{MCE}}, \qquad (7)$$

where $\beta$ is a trade-off parameter. We summarize the whole training procedure of DIFO in Algorithm 1.

## 4. Experiments

**Datasets.** We evaluate four standard benchmarks: **Office-31** [33], **Office-Home** [45], **VisDA** [29] and **DomainNet-126** [30]. Among them, **Office-31** is a small-scaled dataset; **Office-Home** is a medium-scale dataset; **VisDA** and **DomainNet-126** are both large-scale dataset. The details of the four datasets are provided in Supplementary.

**Competitor.** We compare DIFO with 18 existing top-performing methods into three groups. (1) *The first group* contains Source (the source model's results), CLIP [31] and Source+CLIP where Source+CLIP directly average the results of the source model and CLIP. (2) *The second group* includes three state-of-the-art UDA methods

Table 1. Closed-set SFDA on **Office-31** (%)

| Method | Venue | A→D | A→W | D→A | D→W | W→A | W→D | Avg. |
|---|---|---|---|---|---|---|---|---|
| Source | – | 79.1 | 76.6 | 59.9 | 95.5 | 61.4 | 98.8 | 78.6 |
| SHOT [25] | ICML20 | 93.7 | 91.1 | 74.2 | 98.2 | 74.6 | **100.** | 88.6 |
| NRC [49] | NIPS21 | 96.0 | 90.8 | 75.3 | 99.0 | 75.0 | **100.** | 89.4 |
| GKD [38] | IROS21 | 94.6 | 91.6 | 75.1 | 98.7 | 75.1 | **100.** | 89.2 |
| HCL [12] | NIPS21 | 94.7 | 92.5 | 75.9 | 98.2 | 77.7 | **100.** | 89.8 |
| AaD [50] | NIPS22 | 96.4 | 92.1 | 75.0 | **99.1** | 76.5 | **100.** | 89.9 |
| AdaCon [2] | CVPR22 | 87.7 | 83.1 | 73.7 | 91.3 | 77.6 | 72.8 | 81.0 |
| CoWA [20] | ICML22 | 94.4 | 95.2 | 76.2 | 98.5 | 77.6 | 99.8 | 90.3 |
| SCLM [40] | NN22 | 95.8 | 90.0 | 75.5 | 98.9 | 75.5 | 99.8 | 89.4 |
| ELR [51] | ICLR23 | 93.8 | 93.3 | 76.2 | 98.0 | 76.9 | **100.** | 89.6 |
| PLUE [26] | CVPR23 | 89.2 | 88.4 | 72.8 | 97.1 | 69.6 | 97.9 | 85.8 |
| TPDS [41] | IJCV23 | 97.1 | 94.5 | 75.7 | 98.7 | 75.5 | 99.8 | 90.2 |
| **DIFO**-C-RN | – | 93.6 | 92.1 | 78.5 | 95.7 | 78.8 | 97.0 | 89.3 |
| **DIFO**-C-B32 | – | **97.2** | **95.5** | **83.0** | 97.2 | **83.2** | 98.8 | **92.5** |

DAPL [9], PADCLIP [18] and ADCLIP [36] that are also multimodal guiding-based. (3) *The third group* comprises 13 current state-of-the-art SFDA models: SHOT [25], NRC [49], GKD [38], HCL [12], AaD [50], AdaCon [2], CoWA [20], SCLM [40], ELR [51], PLUE [26], TPDS [41] and CRS [52].

For comprehensive comparisons, we implement DIFO in two variants: (1) DIFO-C-RN (weak version) and (2) DIFO-C-B32 (strong version). The key distinction lies in the backbone of the CLIP image-encoder. Specifically, for DIFO-C-RN, ResNet101 [11] is employed on the VisDA dataset, while ResNet50 [11] is used on the other three datasets. On the other hand, DIFO-C-B32 adopts ViT-B/32 [10] as the backbone across all datasets.

**SFDA settings.** We consider three distinct settings: the conventional closed-set SFDA setting, the partial-set and the open-set SFDA settings. The experiment implementation details are provided in Supplementary.

### 4.1. Comparison Results

**Comparison on Closed-set SFDA setting.** The comparisons of the four evaluation datasets are listed in Tab. 1∼3. DIFO-C-B32 surpasses the previous best method CoWA (on Office-31), TPDS (on Office-Home) and PLUE (on VisDA) and GKD (on DomainNet-126) by **2.2%**, **9.6% 2.0%** and **11.3%** in average accuracy respectively. Specifically, DIFO-C-B32 obtains the best results on 4 out of 6 tasks on Office-31 while surpassing previous methods on all tasks of the other three datasets. As for DIFO-C-RN, besides Office-31, it obtains the second-best results and beat the previous best methods by **5.9%**, **0.5%** and **8.0%** on Office-Home, VisDA and DomainNet-126 in average accuracy. The comparison of DIFO-C-RN shows that our method can still perform well despite using a weaker CLIP. Based on a strong CLIP (see results of DIFO-C-B32), our method's performance can improve further as we expected. All of the results indicate that the DIFO can boost the cross-domain performance in closed-set SFDA setting.

Table 2. Closed-set SFDA on **Office-Home** and **VisDA** (%). **SF** and **M** means source-free and multimodal, respectively; the full results on **VisDA** are in `Supplementary`.

| Method | Venue | SF | M | Ar→Cl | Ar→Pr | Ar→Rw | Cl→Ar | Cl→Pr | Cl→Rw | Pr→Ar | Pr→Cl | Pr→Rw | Rw→Ar | Rw→Cl | Rw→Pr | Avg. | VisDA Sy→Re |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source | – | – | – | 43.7 | 67.0 | 73.9 | 49.9 | 60.1 | 62.5 | 51.7 | 40.9 | 72.6 | 64.2 | 46.3 | 78.1 | 59.2 | 49.2 |
| DAPL-RN [9] | TNNLS23 | ✗ | ✓ | 54.1 | 84.3 | 84.8 | 74.4 | 83.7 | 85.0 | 74.5 | 54.6 | 84.8 | 75.2 | 54.7 | 83.8 | 74.5 | 86.9 |
| PADCLIP-RN [18] | ICCV23 | ✗ | ✓ | 57.5 | 84.0 | 83.8 | 77.8 | 85.5 | 84.7 | 76.3 | 59.2 | 85.4 | 78.1 | 60.2 | 86.7 | 76.6 | 88.5 |
| ADCLIP-RN [36] | ICCVW23 | ✗ | ✓ | 55.4 | 85.2 | 85.6 | 76.1 | 85.8 | 86.2 | 76.7 | 56.1 | 85.4 | 76.8 | 56.1 | 85.5 | 75.9 | 87.7 |
| SHOT [25] | ICML20 | ✓ | ✗ | 56.7 | 77.9 | 80.6 | 68.0 | 78.0 | 79.4 | 67.9 | 54.5 | 82.3 | 74.2 | 58.6 | 84.5 | 71.9 | 82.7 |
| NRC [49] | NIPS21 | ✓ | ✗ | 57.7 | 80.3 | 82.0 | 68.1 | 79.8 | 78.6 | 65.3 | 56.4 | 83.0 | 71.0 | 58.6 | 85.6 | 72.2 | 85.9 |
| GKD [38] | IROS21 | ✓ | ✗ | 56.5 | 78.2 | 81.8 | 68.7 | 78.9 | 79.1 | 67.6 | 54.8 | 82.6 | 74.4 | 58.5 | 84.8 | 72.2 | 83.0 |
| AaD [50] | NIPS22 | ✓ | ✗ | 59.3 | 79.3 | 82.1 | 68.9 | 79.8 | 79.5 | 67.2 | 57.4 | 83.1 | 72.1 | 58.5 | 85.4 | 72.7 | 88.0 |
| AdaCon [2] | CVPR22 | ✓ | ✗ | 47.2 | 75.1 | 75.5 | 60.7 | 73.3 | 73.2 | 60.2 | 45.2 | 76.6 | 65.6 | 48.3 | 79.1 | 65.0 | 86.8 |
| CoWA [20] | ICML22 | ✓ | ✗ | 56.9 | 78.4 | 81.0 | 69.1 | 80.0 | 79.9 | 67.7 | 57.2 | 82.4 | 72.8 | 60.5 | 84.5 | 72.5 | 86.9 |
| SCLM [40] | NN22 | ✓ | ✗ | 58.2 | 80.3 | 81.5 | 69.3 | 79.0 | 80.7 | 69.0 | 56.8 | 82.7 | 74.7 | 60.6 | 85.0 | 73.0 | 85.3 |
| ELR [51] | ICLR23 | ✓ | ✗ | 58.4 | 78.7 | 81.5 | 69.2 | 79.5 | 79.3 | 66.3 | 58.0 | 82.6 | 73.4 | 59.8 | 85.1 | 72.6 | 85.8 |
| PLUE [26] | CVPR23 | ✓ | ✗ | 49.1 | 73.5 | 78.2 | 62.9 | 73.5 | 74.5 | 62.2 | 48.3 | 78.6 | 68.6 | 51.8 | 81.5 | 66.9 | 88.3 |
| TPDS [41] | IJCV23 | ✓ | ✗ | 59.3 | 80.3 | 82.1 | 70.6 | 79.4 | 80.9 | 69.8 | 56.8 | 82.1 | 74.5 | 61.2 | 85.3 | 73.5 | 87.6 |
| **DIFO**-C-RN | – | ✓ | ✓ | 62.6 | 87.5 | 87.1 | 79.5 | 87.9 | 87.4 | 78.3 | 63.4 | 88.1 | 80.0 | 63.3 | 87.7 | 79.4 | 88.8 |
| **DIFO**-C-B32 | – | ✓ | ✓ | 70.6 | 90.6 | 88.8 | 82.5 | 90.6 | 88.8 | 80.9 | 70.1 | 88.9 | 83.4 | 70.5 | 91.2 | 83.1 | 90.3 |

Table 3. Closed-set SFDA on **DomainNet-126** (%). **SF** and **M** means source-free and multimodal, respectively.

| Method | Venue | SF | M | C→P | C→R | C→S | P→C | P→R | P→S | R→C | R→P | R→S | S→C | S→P | S→R | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source | – | – | – | 44.6 | 59.8 | 47.5 | 53.3 | 75.3 | 46.2 | 55.3 | 62.7 | 46.4 | 55.1 | 50.7 | 59.5 | 54.7 |
| DAPL-RN [9] | TNNLS23 | ✗ | ✓ | 72.4 | 87.6 | 65.9 | 72.7 | 87.6 | 65.6 | 73.2 | 72.4 | 66.2 | 73.8 | 72.9 | 87.8 | 74.8 |
| ADCLIP-RN [36] | ICCVW23 | ✗ | ✓ | 71.7 | 88.1 | 66.0 | 73.2 | 86.9 | 65.2 | 73.6 | 73.0 | 68.4 | 72.3 | 74.2 | 89.3 | 75.2 |
| SHOT [25] | ICML20 | ✓ | ✗ | 63.5 | 78.2 | 59.5 | 67.9 | 81.3 | 61.7 | 67.7 | 67.6 | 57.8 | 70.2 | 64.0 | 78.0 | 68.1 |
| GKD [38] | IROS21 | ✓ | ✗ | 61.4 | 77.4 | 60.3 | 69.6 | 81.4 | 63.2 | 68.3 | 68.4 | 59.5 | 71.5 | 65.2 | 77.6 | 68.7 |
| NRC [49] | NIPS21 | ✓ | ✗ | 62.6 | 77.1 | 58.3 | 62.9 | 81.3 | 60.7 | 64.7 | 69.4 | 58.7 | 69.4 | 65.8 | 78.7 | 67.5 |
| AdaCon [2] | CVPR22 | ✓ | ✗ | 60.8 | 74.8 | 55.9 | 62.2 | 78.3 | 58.2 | 63.1 | 68.1 | 55.6 | 67.1 | 66.0 | 75.4 | 65.4 |
| CoWA [20] | ICML22 | ✓ | ✗ | 64.6 | 80.6 | 60.6 | 66.2 | 79.8 | 60.8 | 69.0 | 67.2 | 60.0 | 69.0 | 65.8 | 79.9 | 68.6 |
| PLUE [26] | CVPR23 | ✓ | ✗ | 59.8 | 74.0 | 56.0 | 61.6 | 78.5 | 57.9 | 61.6 | 65.9 | 53.8 | 67.5 | 64.3 | 76.0 | 64.7 |
| TPDS [41] | IJCV23 | ✓ | ✗ | 62.9 | 77.1 | 59.8 | 65.6 | 79.0 | 61.5 | 66.4 | 67.0 | 58.2 | 68.6 | 64.3 | 75.3 | 67.1 |
| **DIFO**-C-RN | – | ✓ | ✓ | 73.8 | 89.0 | 69.4 | 74.0 | 88.7 | 70.1 | 74.8 | 74.6 | 69.6 | 74.7 | 74.3 | 88.0 | 76.7 |
| **DIFO**-C-B32 | – | ✓ | ✓ | 76.6 | 87.2 | 74.9 | 80.0 | 87.4 | 75.6 | 80.8 | 77.3 | 75.5 | 80.5 | 76.7 | 87.3 | 80.0 |

Table 4. Results (%) of CLIP and Source+CLIP on the four evaluation datasets. The backbone of CLIP image-encoder in CLP-C-RN and CLP-C-B32 are the same as **DIFO**-C-RN and **DIFO**-C-B32, respectively. The full results are provided in `Supplementary`.

| Method | Venue | Office-31 →A | →D | →W | →Avg. | Office-Home →Ar | →Cl | →Pr | →Rw | →Avg. | VisDA Sy→Re | DomainNet-126 →C | →P | →R | →S | →Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CLIP-RN [31] | ICML21 | 73.1 | 73.9 | 67.0 | 71.4 | 72.5 | 51.9 | 81.5 | 82.5 | 72.1 | 83.7 | 67.9 | 70.2 | 87.1 | 65.4 | 72.7 |
| Source+CLIP-RN | – | 76.3 | 90.4 | 84.0 | 83.6 | 75.4 | 57.4 | 84.4 | 85.7 | 75.7 | 82.0 | 71.8 | 71.4 | 87.3 | 66.5 | 74.3 |
| **DIFO**-C-RN | – | 78.6 | 95.3 | 93.9 | 89.3 | 79.3 | 63.1 | 87.7 | 87.5 | 79.4 | 88.8 | 74.5 | 74.2 | 88.5 | 69.7 | 76.7 |
| CLIP-B32 [31] | ICML21 | 76.0 | 82.7 | 80.6 | 79.8 | 74.6 | 59.8 | 84.3 | 85.5 | 76.1 | 82.9 | 74.7 | 73.5 | 85.7 | 71.2 | 76.3 |
| Source+CLIP-B32 | – | 78.5 | 93.0 | 89.6 | 87.0 | 78.9 | 62.5 | 86.1 | 87.7 | 78.8 | 82.0 | 76.8 | 73.7 | 86.0 | 70.8 | 76.8 |
| **DIFO**-C-B32 | – | 83.1 | 98.0 | 96.4 | 92.5 | 82.3 | 70.4 | 90.8 | 88.8 | 83.1 | 90.3 | 80.4 | 76.9 | 87.3 | 75.3 | 80.0 |

**Comparison to CLIP based prediction results.** The original CLIP model can conduct general image classification. We carry out a quantitative comparison between DIFO's adaptation performance and CLIP's performance on the four datasets, averaging the adaptation results of DIFO grouped by the target domain name.

As presented in the bottom of Tab. 4, DIFO-C-B32 outperforms CLIP-B32 on all tasks. On average accuracy, DIFO-C-B32 increases the performance by **12.7**%, **7.0**%, **7.4**% and **3.7**% in Office-31, Office-Home, VisDA and DomainNet-126, respectively. Regarding the weak version, as reported in the top, DIFO-C-RN maintains similar advantages with the increase of **17.9**%, **7.3**%, **5.1**% and **4.0**%. The result shows that *the domain generality of the original CLIP model cannot fully excel to the target domain, and task-specific customization is needed.*

Interestingly, compared with CLIP-B32, except for VisDA with a tiny gap of **0.9**%, Source+CLIP-B32 averagely improve by **7.2**% at most on the other datasets. Meanwhile, Source+CLIP-B32 is beaten by DIFO-C-B32 with an increase of **3.2**% at least. In the group of DIFO-C-RN, we have the same observation. These results imply that directly weighting the source model and CLIP is an intuitive knowledge adaptation scheme, but it is hard to perform adaptation deeply. Considering Source+CLIP is an average version, we conduct a comprehensive comparison with the weighting strategy where the weighting coefficient of CLIP prediction varies from 0.0 to 1.0. Here, we conduct this experiment
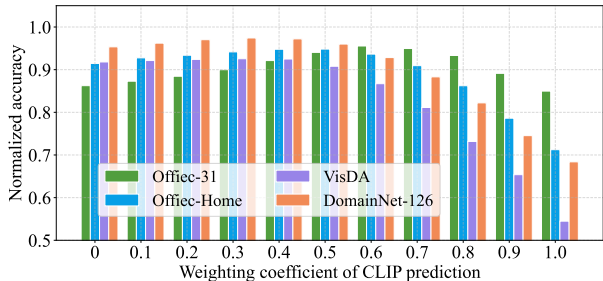
Figure 4. The performance of the scheme directly weighting the source model and CLIP-B32. All results are normalized by corresponding DIFO-C-B32 accuracies for a clear view.

based on more challenging CLIP-B32 due to its large performance gap with Source (see the first row in `Tab.1∼3`). For a clear view, all weighted accuracies are normalized by the corresponding DIFO-C-B32 accuracies, respectively. As shown in Fig. 4, no result can exceed the value of 1.0. This indicates that *weighting the source model and CLIP in a zero-shot manner cannot obtain desirable task-specific fusion, and a carefully designed distilling is necessary.*

**Comparison on Partial-set and Open-set SFDA settings.** These are the variations of traditional Closed-set SFDA setting, following the same as SHOT [25] (the detailed setting introduction is provided in `Supplementary`). As reported in Tab. 5, compared with previous best method CoWA (Partial-set) and CRS (Open-set), our DIFO-C-B32 improves by **2.4**% and **2.7**%, respectively.

## 4.2. Model Analysis

**Feature distribution visualization.** Taking task Ar→Cl in Office-Home as a toy experiment, we visualize feature distribution using t-SNE tool. Meanwhile, we choose 5 comparisons, including the source model (termed Source), CLIP-B32's zero shot (termed CLIP), SHOT, TPDS and Oracle (trained on domain Cl with the real labels). As shown at the top of Fig. 5, from Souce to DIFO-C-B32, category aliasing gradually relieves. Compared with Oracle, DIFO-C-B32 has the most similar distribution shape. To verify this point, we also give the 3D Density chart results arranged at the bottom of Fig. 5. These results confirm the effectiveness of our DIFO-C-B32 in terms of Feature distribution.

**Ablataion study.** We evaluate the (1) effect of $L_{\text{TSC}}$, $L_{\text{MCE}}$ and $L_{\text{PC}}$, (2) effect of optimization of mutual information, (3) effect of task-specific customization and (4) effect of historical prediction fusion.

For this first issue, we conduct a progressive experiment to isolate the loss's effects. The top four rows of Tab. 6 list the ablation study results. For convenience comparison, the baseline (the first row) is the source model results. When single $L_{\text{TSP}}$ or $L_{\text{MCE}}$ works, the accuracy largely increases on the three datasets with an improvement of about **20**% in average accuracy compared with the baseline. As both of them

Table 5. Partial-set SFDA and Open-set SFDA on **Office-Home** (%). The full results are provided in `Supplementary`.

| Partial-set SFDA | Venue | Avg. | Open-set SFDA | Venue | Avg. |
|---|---|---|---|---|---|
| Source | – | 62.8 | Source | – | 46.6 |
| SHOT [25] | ICML20 | 79.3 | SHOT [25] | ICML20 | 72.8 |
| HCL [12] | NIPS21 | 79.6 | HCL [12] | NIPS21 | 72.6 |
| CoWA [20] | ICML22 | 83.2 | CoWA [20] | ICML22 | 73.2 |
| AaD [50] | NIPS22 | 79.7 | AaD [50] | NIPS22 | 71.8 |
| CRS [52] | CVPR23 | 80.6 | CRS [52] | CVPR23 | 73.2 |
| **DIFO**-C-B32 | – | **85.6** | **DIFO**-C-B32 | – | **75.9** |

Table 6. Classification results of ablation study (%) on **Office-31 Office-Home** and **VisDA**.

| $L_{\text{TSC}}$ | $L_{\text{MCE}}$ | $L_{\text{PC}}$ | **Office-31** | **Office-Home** | **VisDA** | Avg. |
|---|---|---|---|---|---|---|
| ✗ | ✗ | ✗ | 78.6 | 59.2 | 49.2 | 62.3 |
| ✓ | ✗ | ✗ | 82.4 | 77.4 | 84.4 | 81.4 |
| ✗ | ✓ | ✗ | 82.1 | 76.5 | 88.6 | 82.4 |
| ✓ | ✓ | ✗ | 87.0 | 80.0 | 88.3 | 85.1 |
| ✓ | ✓ | ✓ | **92.5** | **83.1** | **90.3** | **88.6** |
| **DIFO**-C-B32 w/ KL | | | 90.4 | 81.5 | 89.0 | 87.0 |
| **DIFO**-C-B32 w/ CLIP | | | 90.7 | 81.1 | 88.8 | 86.8 |
| **DIFO**-C-B32 w/o $\boldsymbol{p}_i'$ | | | 89.8 | 73.5 | 87.0 | 83.4 |
| **DIFO**-C-B32 w/o $\boldsymbol{p}_i$ | | | 88.9 | 82.2 | 88.9 | 86.7 |

are adopted, the accuracy evident increase (**3.7**% in average, the fourth row) on the top of the case of only $L_{\text{TSC}}$ and further enhanced by adopting of item $L_{\text{PC}}$ (**3.5**% in average, the fifth row). The results indicate: (1) all objective components positively affect the final performance, (2) $L_{\text{MCE}}$, $L_{\text{PC}}$ is crucial due to providing a new soft supervision for coarse-to-fine adaptation.

For the second and third issues, we propose two variation methods of DIFO-C-B32 to evaluate the effect. One is DIFO-C-B32 w/ KL where the mutual information maximization loss in $L_{\text{TSC}}$, $L_{\text{PC}}$ are replace by KL divergence loss. The other one is DIFO-C-B32 w/ CLIP where the prompt learning-based customization for CLIP is cancelled, and the inputted prompt is set to the fixed template of *"a photo of a [CLS]."* during the entire adaptation. As presented in the last two rows in Tab. 6, DIFO-C-B32 (the fifth row) beats DIFO-C-B32 w/ KL and DIFO-C-B32 w/ CLIP with average improvement of **1.6**% at least, respectively confirming the effect of adopting mutual information optimization and task-specific customization. As for the fourth issue, its effect is verified by the performance decreases (**3.4**% in average at most) in the variation methods (the last two rows), which remove $\boldsymbol{p}_i'$ and $\boldsymbol{p}_i$ from the fusion respectively.

## 4.3. Task-Specific Knowledge Adaptation Analysis

In this part, we give a feature space shift analysis using the measure of MMD (maximum mean discrepancy) distance [54] to verify whether the proposed method ensures a task-specific knowledge adaptation.

In this experiment, we first train a domain-invariant Oracle model over all Office-Home data with real labels, and use
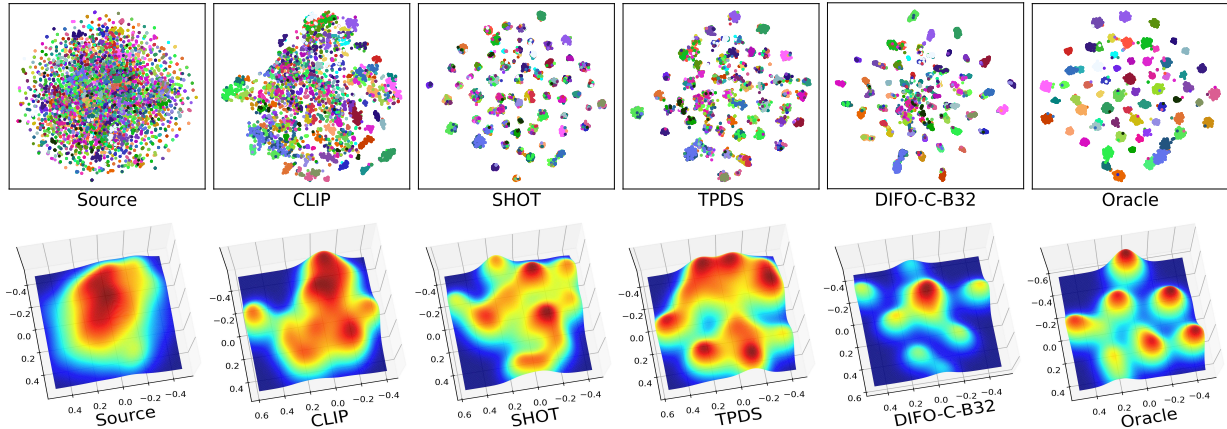
Figure 5. Feature distribution visualization comparison on transfer task Ar→Cl in Office-Home. Oracle is trained on target domain Cl using the ground-truth labels. Different colors stand for different categories. **Top**: t-SNE feature distribution over 65 categories. **Bottom**: The corresponding 3D density charts. For easy view, the first 10 categories were used in this plot.
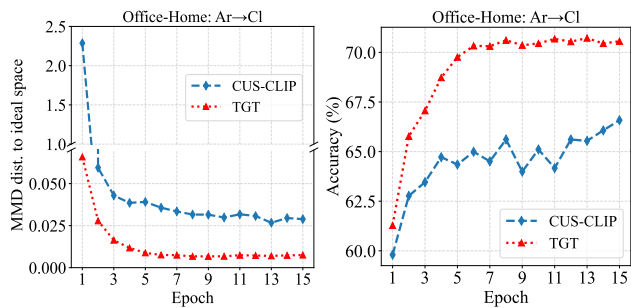


Figure 6. The evolving dynamics of MMD distance during adaption of Ar→Cl in Office-Home. **Left** and **Right** present the varying curves of MMD distance and accuracy, respectively

the logits to express the ideal task-specific space $\mathcal{O}$. After that, an analysis is conducted on the transfer task Ar→Cl. During this adaptation, there are $T$ (epoch number) intermediate target models and customized CLIP models. We feedforward the target data through each intermediate model and take the logits as a space. Thus, we obtain $T$ intermediate target feature spaces $\{\mathcal{U}_k\}_{k=1}^T$ and $T$ intermediate customized CLIP feature spaces $\{\mathcal{V}\}_{k=1}^T$. Within this context, these intermediate spaces can depict the task-specific distillation to $\mathcal{O}$. In practice, the CLIP image encoder's backbone is set to ViT-B/32.

In the left of Fig. 6, we give the MMD distance change curve of $\{\mathcal{U}_k\}_{k=1}^T$ (in red, termed TGT) and $\{\mathcal{V}\}_{k=1}^T$ (in blue, termed CUS-CLIP), taking $\mathcal{O}$ as the original space. It is seen that at early epochs (1~4), TGT and CUS-CLIP sharply decrease and then maintain a gradual decrease in the following epochs. Meanwhile, this change is consistent with the accuracy varying shown in the right of Fig. 6.

These results indicate that our DIFO encourages task-specific knowledge adaptation due to converging the ideal task-specific space. Besides, we observe two details. First, after epoch 1, CUS-CLIP's distance reduces by **2.2**, which is **58.6** time of TGT's decrease of **0.038**. This is because CLIP represents a heterogeneous space of vision-language, much different from the vision space $\mathcal{O}$. Furthermore, the large distance decrease confirms the effect of customization. Second, the synchronized distance reductions of CUS-CLIP and TGT indicate the interaction between the target model and CLIP is a crucial design for task-specific distillation.

## 5. Conclusion

We present an innovative approach, referred to as DIFO, designed to tackle the SFDA problem. To the best of our knowledge, this marks the initial endeavor to address SFDA by leveraging a pretrained ViL foundation model, departing from previous approaches that predominantly concentrated on self-mining auxiliary information. DIFO is featured with alternating between customization of the ViL model and the transfer of task-specific knowledge from the customized ViL model. We introduce two pivotal designs: a mutual information-based alignment for ViL customization and a most-likely category encouragement for more precise adaptation of task-specific knowledge. Our method's effectiveness is validated by state-of-the-art experimental results across four challenging datasets.

# Source-Free Domain Adaptation with Frozen Multimodal Foundation Model

## Supplementary Material

## 6. A Proof of Theorem 1.

**Restatement of Theorem 1** *Given two random variables $X$, $Y$. Their mutual information $\mathrm{I}(X,Y)$ and KL divergence $D_{\mathrm{KL}}(X||Y)$ satisfy the unequal relationship as follows.*

$$-\mathrm{I}(X,Y) \leq D_{\mathrm{KL}}(X||Y). \qquad (8)$$

*Proof.* Suppose the probability density function (PDF) of $X$ and $Y$ are $p(\boldsymbol{x})$ and $p(\boldsymbol{y})$, respectively; their join PDF is $p(\boldsymbol{x},\boldsymbol{y})$. We have

$$\mathrm{I}(X,Y) = \sum p(\boldsymbol{x},\boldsymbol{y}) \log \frac{p(\boldsymbol{x},\boldsymbol{y})}{p(\boldsymbol{x}) \cdot p(\boldsymbol{y})}$$
$$= D_{\mathrm{KL}}(p(\boldsymbol{x},\boldsymbol{y}) || p(\boldsymbol{x}) \cdot p(\boldsymbol{y})).$$

Well known, the KL divergence is non-negative [7]. Thus,

$$-\mathrm{I}(X,Y) \leq 0 \leq D_{\mathrm{KL}}(X||Y)$$

## 7. Evaluation Datasets

We evaluate four standard benchmarks below.
- **Office-31** [33] is a small-scaled dataset including three domains, i.e., Amazon (A), Webcam (W), and Dslr (D), all of which are taken of real-world objects in various office environments. The dataset has 4,652 images of 31 categories in total. Images in (A) are online e-commerce pictures. (W) and (D) consist of low-resolution and high-resolution pictures.
- **Office-Home** [45] is a medium-scale dataset that is mainly used for domain adaptation, all of which contains 15k images belonging to 65 categories from working or family environments. The dataset has four distinct domains, i.e., Artistic images (Ar), Clip Art (Cl), Product images (Pr), and Real-word images (Rw).
- **VisDA** [29] is a challenging large-scale dataset with 12 types of synthetic to real transfer recognition tasks. The source domain contains 152k synthetic images (Sy), whilst the target domain has 55k real object images (Re) from the famous Microsoft COCO dataset.
- **DomainNet-126** [30] is another large-scale dataset. As a subset of DomainNet containing 600k images of 345 classes from 6 domains of different image styles, this dataset has 145k images from 126 classes, sampled from 4 domains, Clipart (C), Painting (P), Real (R), Sketch (S), as [34] identify severe noisy labels in the dataset.

## 8. Implementation Details

**Souce model pre-training.** For all transfer tasks on the three datasets, we train the source model $\theta_s$ on the source domain in a supervised manner using the following objective of the classic cross-entropy loss with smooth label, like other methods [25, 40, 49].

$$L_{\mathrm{s}}(\mathcal{X}_s, \mathcal{Y}_s; \theta_s) = -\frac{1}{n_s} \sum_{i=1}^{n_s} \sum_{c=1}^{C} \tilde{l}_{i,c}^s \log p_{i,c}^s,$$

where $n_s$ is the number of the source data, $p_{i,c}^s$ is the $c$-th element of $\boldsymbol{p}_i^s = \theta_s(\boldsymbol{x}_i^s)$ that is the category probability vector of input instance $\boldsymbol{x}_i^s$ after $\theta_s$ mapping; $\tilde{l}_{i,c}^s$ is the $c$-th element of the smooth label [28] $\tilde{\boldsymbol{l}}_i^s = (1-\sigma)\boldsymbol{l}_i^s + \sigma/C$, in which $\boldsymbol{l}_i^s$ is a one-hot encoding of hard label $y_i^s$ and $\sigma = 0.1$. The source dataset is divided into the training set and testing set in a 0.9:0.1 ratio.

**Network setting.** The DIFO model contains two network branches. In the target model branch, the feature extractor consists of a deep architecture and a fully-connected layer followed by a batch-normalization layer. Same to the previous work [25, 27, 32, 48, 49], the deep architecture is transferred from the deep models pre-trained on ImageNet (i.e., ResNet-50 is used on **Office-31**, **Office-Home** and **DomainNet-126**, whilst ResNet-101 is adopted on **VisDA**). The ending classifier is a fully-connected layer with weight normalization. On the other hand, the ViL model branch chooses the most adopted CLIP as the implementation where the text encoder's transformer-based architecture follows modification proposed in [31] as the backbone. Regarding the image encoder, we adopt two versions corresponding to the two implementations of DIFO in this paper, including DIFO-C-B32 and DIFO-C-RN. Specifically, in DIFO-C-B32, image encoders follow ViT-B/32 architecture proposed in CLIP [31] while DIFO-C-RN uses ResNet [11] as the backbone. The same as the target model mentioned above, ResNet-101 is adopted on **VisDA** and ResNet-50 is used on the rest datasets.

**Parameter setting.** For the trade-off parameter $\alpha$ and $\beta$ in the objective $L_{\mathrm{PC}}$ (Eq. (6)) and $L_{\mathrm{MKA}}$ (Eq. (7)) is set to 1.0 and 0.4 on all datasets, respectively. The parameter of Exponential distribution $\lambda$ in Eq. (4) is specified to 10.0. The temperature parameters in Eq. (5) are $\tau = 0.1$. The number of the most-likely categories is set to $N = 2$.

**Training setting.** We adopt the batch size of 64, SGD optimizer with momentum 0.9 and 15 training epochs on all datasets. The prompt template for initiation is the mostly used *'a photo of a [CLASS].'* [31] where [CLASS] stands for the class name. All experiments are conducted with PyTorch on a single GPU of NVIDIA RTX.

Table 7. Full results (%) of Closed-set SFDA on **VisDA**. **SF** and **M** mean source-free and multimodal, respectively.

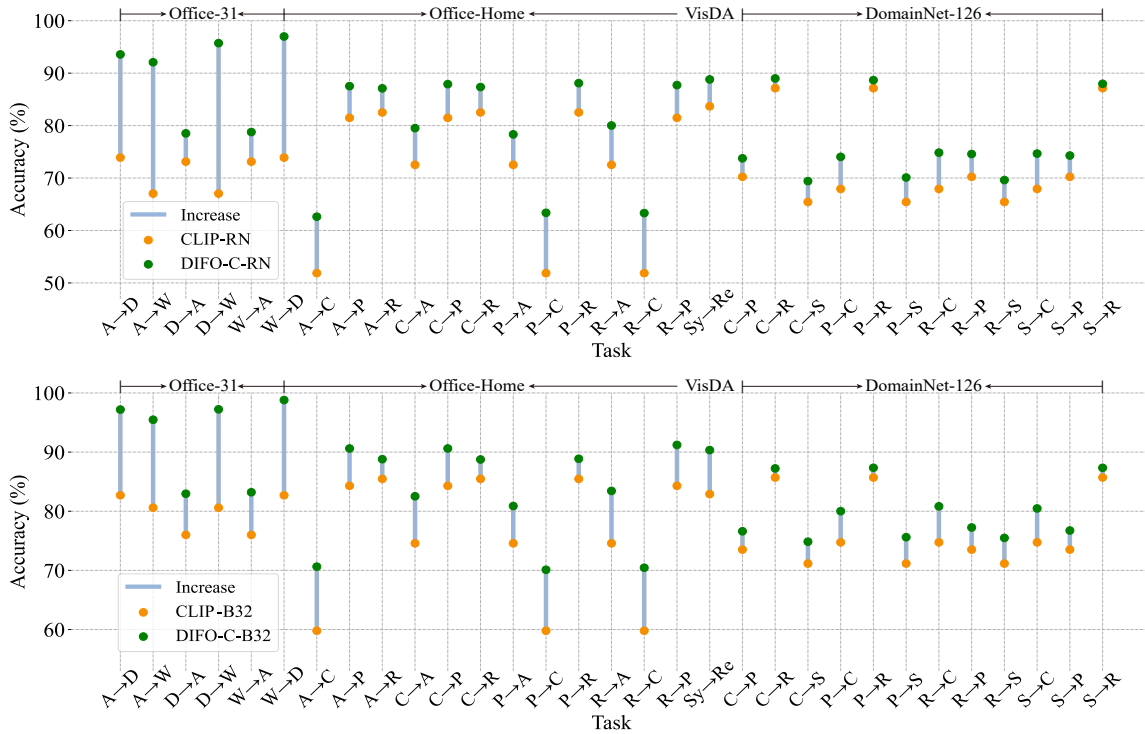| Method | Venue | SF | M | plane | bcycl | bus | car | horse | knife | mcycl | person | plant | sktbrd | train | truck | Perclass |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source | - | - | - | 60.7 | 21.7 | 50.8 | 68.5 | 71.8 | 5.4 | 86.4 | 20.2 | 67.1 | 43.3 | 83.3 | 10.6 | 49.2 |
| DAPL-RN [9] | TNNLS23 | ✗ | ✓ | **98.1** | 83.1 | 88.8 | 77.9 | 97.4 | 91.5 | 94.2 | 79.7 | 88.6 | 89.3 | 92.5 | 62.0 | 86.9 |
| PADCLIP-RN [18] | ICCV23 | ✗ | ✓ | 96.7 | 88.8 | 87.0 | 82.8 | 97.1 | 93.0 | 91.3 | 83.0 | 95.5 | 91.8 | 91.5 | 63.0 | 88.5 |
| ADCLIP-RN [36] | ICCVW23 | ✗ | ✓ | **98.1** | 83.6 | **91.2** | 76.6 | **98.1** | 93.4 | **96.0** | 81.4 | 86.4 | 91.5 | 92.1 | 64.2 | 87.7 |
| SHOT [25] | ICML20 | ✓ | ✗ | 95.0 | 87.4 | 80.9 | 57.6 | 93.9 | 94.1 | 79.4 | 80.4 | 90.9 | 89.8 | 85.8 | 57.5 | 82.7 |
| NRC [49] | NIPS21 | ✓ | ✗ | 96.8 | 91.3 | 82.4 | 62.4 | 96.2 | 95.9 | 86.1 | **90.7** | 94.8 | 94.1 | 90.4 | 59.7 | 85.9 |
| GKD [38] | IROS21 | ✓ | ✗ | 95.3 | 87.6 | 81.7 | 58.1 | 93.9 | 94.0 | 80.0 | 80.0 | 91.2 | 91.0 | 86.9 | 56.1 | 83.0 |
| AaD [50] | NIPS22 | ✓ | ✗ | 97.4 | 90.5 | 80.8 | 76.2 | 97.3 | 96.1 | 89.8 | 82.9 | 95.5 | 93.0 | 92.0 | 64.7 | 88.0 |
| AdaCon [2] | CVPR22 | ✓ | ✗ | 97.0 | 84.7 | 84.0 | 77.3 | 96.7 | 93.8 | 91.9 | 84.8 | 94.3 | 93.1 | 94.1 | 49.7 | 86.8 |
| CoWA [20] | ICML22 | ✓ | ✗ | 96.2 | 89.7 | 83.9 | 73.8 | 96.4 | **97.4** | 89.3 | 86.8 | 94.6 | 92.1 | 88.7 | 53.8 | 86.9 |
| SCLM [40] | NN22 | ✓ | ✗ | 97.1 | 90.7 | 85.6 | 62.0 | 97.3 | 94.6 | 81.8 | 84.3 | 93.6 | 92.8 | 88.0 | 55.9 | 85.3 |
| ELR [51] | ICLR23 | ✓ | ✗ | 97.1 | 89.7 | 82.7 | 62.0 | 96.2 | 97.0 | 87.6 | 81.2 | 93.7 | 94.1 | 90.2 | 58.6 | 85.8 |
| PLUE [26] | CVPR23 | ✓ | ✗ | 94.4 | **91.7** | 89.0 | 70.5 | 96.6 | 94.9 | 92.2 | 88.8 | 92.9 | 95.3 | 91.4 | 61.6 | 88.3 |
| TPDS [41] | IJCV23 | ✓ | ✗ | 97.6 | 91.5 | 89.7 | 83.4 | 97.5 | 96.3 | 92.2 | 82.4 | **96.0** | 94.1 | 90.9 | 40.4 | 87.6 |
| **DIFO**-C-RN | - | ✓ | ✓ | 97.7 | 87.6 | 90.5 | **83.6** | 96.7 | 95.8 | 94.8 | 74.1 | 92.4 | 93.8 | 92.9 | 65.5 | 88.8 |
| **DIFO**-C-B32 | - | ✓ | ✓ | 97.5 | 89.0 | 90.8 | 83.5 | 97.8 | 97.3 | 93.2 | 83.5 | 95.2 | **96.8** | **93.7** | **65.9** | **90.3** |



Figure 7. Transfer performance comparison of **DIFO** and CLIP on all tasks of the four evaluation datasets. **Top: DIFO**-C-RN **v.s.** CLIP-RN. **Bottom: DIFO**-C-B32 **v.s.** CLIP-B32.

# 9. Supplementation of Full Experiment Results

**Full results on VisDA.** As the supplement of results on VisDA, Tab. 7 presents the full classification details over the 12 categories. It is seen that DIFO-C-RN and DIFO-C-B32 obtain the best results in 7/12 categories compared with SFDA methods. Meanwhile, DIFO-C-RN and DIFO-C-B32 are on top of the second best UDA results in 8/12 categories. Also, we note that the UDA method of ADCLIP beats DIFO-C-RN and DIFO-C-B32 on four transfer tasks. It

is understandable that ADCLIP use the labelled source data, whilst our method cannot access the source data. Despite this, DIFO still presents advantages over these source data-required method (see the average accuracy).

**Full results of comparison to CLIP.** As the supplementation of these domain-grouped results reported in the paper, Fig. 7 gives a comprehensive visualization comparison with CLIP in the perspective of all 31 transfer tasks on the four evaluation datasets. It is seen that the results of

Table 8. Full results (%) of Partial-set SFDA and Open-set SFDA on **Office-Home**.

| Partial-set SFDA | Venue | Ar→Cl | Ar→Pr | Ar→Rw | Cl→Ar | Cl→Pr | Cl→Rw | Pr→Ar | Pr→Cl | Pr→Rw | Rw→Ar | Rw→Cl | Rw→Pr | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source | – | 45.2 | 70.4 | 81.0 | 56.2 | 60.8 | 66.2 | 60.9 | 40.1 | 76.2 | 70.8 | 48.5 | 77.3 | 62.8 |
| SHOT [25] | ICML20 | 64.8 | 85.2 | **92.7** | 76.3 | 77.6 | 88.8 | 79.7 | 64.3 | 89.5 | 80.6 | 66.4 | 85.8 | 79.3 |
| HCL [12] | NIPS21 | 65.6 | 85.2 | **92.7** | 77.3 | 76.2 | 87.2 | 78.2 | 66.0 | 89.1 | 81.5 | 68.4 | 87.3 | 79.6 |
| CoWA [20] | ICML22 | 69.6 | 93.2 | 92.3 | 78.9 | 81.3 | 92.1 | 79.8 | 71.7 | 90.0 | 83.8 | **72.2** | **93.7** | 83.2 |
| AaD [50] | NIPS22 | 67.0 | 83.5 | 93.1 | 80.5 | 76.0 | 87.6 | 78.1 | 65.6 | 90.2 | 83.5 | 64.3 | 87.3 | 79.7 |
| CRS [52] | CVPR23 | 68.6 | 85.1 | 90.9 | 80.1 | 79.4 | 86.3 | 79.2 | 66.1 | 90.5 | 82.2 | 69.5 | 89.3 | 80.6 |
| **DIFO**-C-B32 | – | **70.2** | **91.7** | 91.5 | **87.8** | **92.6** | **92.9** | **87.3** | **70.7** | **92.9** | **88.5** | 69.6 | 91.5 | **85.6** |
| Open-set SFDA | Venue | Ar→Cl | Ar→Pr | Ar→Rw | Cl→Ar | Cl→Pr | Cl→Rw | Pr→Ar | Pr→Cl | Pr→Rw | Rw→Ar | Rw→Cl | Rw→Pr | Avg. |
| Source | – | 36.3 | 54.8 | 69.1 | 33.8 | 44.4 | 49.2 | 36.8 | 29.2 | 56.8 | 51.4 | 35.1 | 62.3 | 46.6 |
| SHOT [25] | ICML20 | 64.5 | 80.4 | 84.7 | 63.1 | 75.4 | 81.2 | 65.3 | 59.3 | 83.3 | 69.6 | 64.6 | 82.3 | 72.8 |
| HCL [12] | NIPS21 | 64.0 | 78.6 | 82.4 | 64.5 | 73.1 | 80.1 | 64.8 | 59.8 | 75.3 | **78.1** | **69.3** | 81.5 | 72.6 |
| CoWA [20] | ICML22 | 63.3 | 79.2 | 85.4 | 67.6 | **83.6** | 82.0 | 66.9 | 56.9 | 81.1 | 68.5 | 57.9 | **85.9** | 73.2 |
| AaD [50] | NIPS22 | 63.7 | 77.3 | 80.4 | 66.0 | 72.6 | 77.6 | 69.1 | **62.5** | 79.8 | 71.8 | 62.3 | 78.6 | 71.8 |
| CRS [52] | CVPR23 | **65.2** | 76.6 | 80.2 | 66.2 | 75.3 | 77.8 | 70.4 | 61.8 | 79.3 | 71.1 | 61.1 | 78.3 | 73.2 |
| **DIFO**-C-B32 | – | 64.5 | **86.2** | **87.9** | **68.2** | 79.3 | **86.1** | **67.2** | 62.1 | **88.3** | 71.9 | 65.3 | 84.4 | **75.9** |

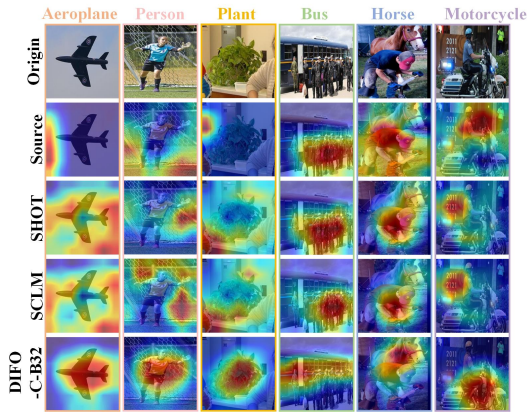

Figure 8. Grad-CAM visualization of **DIFO**-C-B32 and typical comparison methods on toy samples selected from VisDA.



Figure 9. The evolving dynamics of model learning attention based on **DIFO**-C-B32. The red bounding box indicates the failure case.

DIFO (marked by green circles) are above CLIP (marked by orange circles) on all tasks, whether we use DIFO-C-RN or DIFO-C-B32.

**Full results of Partial-set and Open-set SFDA.** As the supplementation of these average results in Tab. 5, Tab. 8 gives the full classification accuracy over 12 transfer tasks in the **Office-Home** dataset. As the top in Tab. 8, DIFO-C-B32 obtains best results on 9/12 tasks in the Partial-set SFDA and on the half tasks in the Open-set SFDA.

## 10. Expanded Model Analysis

**Grad-CAM visualization.** In Fig. 8, we present the Grad-CAM visualization [35] comparison with the source model and two typical SFDA methods, SHOT and SCLM, based on self-supervised learning without ViL model help. For the single object-contained images (see 1∼3 column), DIFO-C-B32's attention focuses on the target object, whilst other
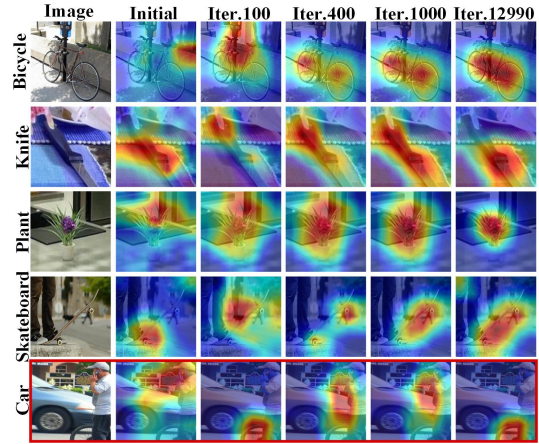
methods cover the entire image. Regarding the multi-object-contained images (see 4∼6 column), DIFO-C-B32's attention is more consistent with the target semantics given by the real labels than other methods focusing on the wrong object. These results explain the effectiveness of DIFO-C-B32 integrating the domain generality of the ViL model and the task specificity of the source model.

**Attention-based evolving dynamics.** To better understand the working of DIFO, this part visualizes the evolving dynamics of model learning attention during the training phase. For a clear view, we display the Grad-CAM visualization results at some typical iterations, as shown in Fig. 9. Among the rightly classified images (the top four rows), the attention smoothly concentrates to the discriminative visual patch. In contrast, the attention of the misclassified image (the last row) converges to the meaningless one.

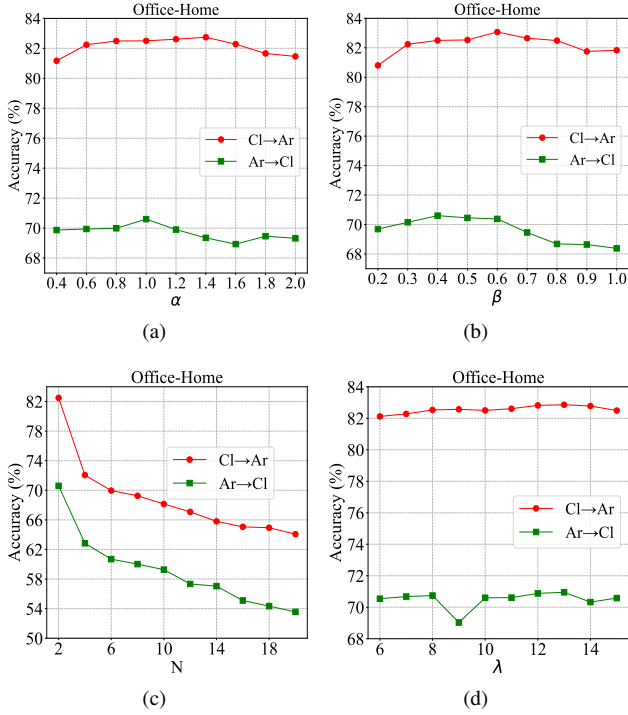**Sensitivity of hyper-parameter.** In the DIFO method, $\alpha$, $\beta$

Figure 10. Performance sensitivity of the hyper-parameters. From (a) to (d), the four sub-figures present the accuracy changing as the parameters $\alpha$, $\beta$, $N$ and $\lambda$ varying, respectively.
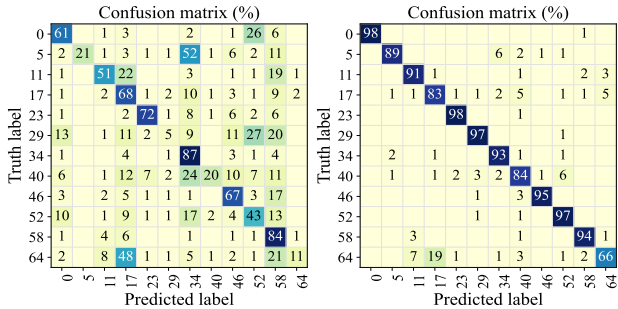


Figure 11. The confusion matrix for 12-way classification on **VisDA**. **Left:** Source model result, **Right:** DIFO-C-B32 result.

are trade-off parameters in objective $L_{\mathrm{PC}}$ (see Eq. (6)) and $L_{\mathrm{MKA}}$ (see Eq. (7)). $\lambda$ is the parameter of Exponential distribution in Eq. (4), whilst $N$ is the number of the most-likely categories. This part discusses their performance sensitivity based on the symmetric transfer tasks Cl→Ar and Ar→Cl in the **Office-Home** dataset. As depicted in Fig. 10 (a), (b) and (d), when these parameters changes, there are no evident drops in the accuracy variation curves. This indicates that DIFO is insensitive to parameters $\alpha$, $\beta$ and $\lambda$. As for $N$, the accuracy gradually decreases as $N$ increases. This phenomenon is consistent with our expectation that small $N$
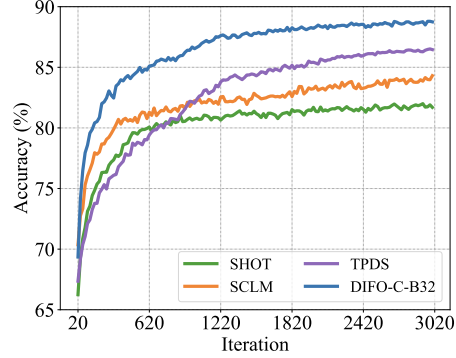


Figure 12. Classification accuracy varying curve comparison on **VisDA** during the adaptation phase.

is better and a large value will introduce the semantic noise.

**Confusion matrix.** To present a quantitative observation on the category, this part gives the confusion matrix based on the classification results on the **VisDA** dataset. For comparison, we show the confusion matrix of the source model at the left side of Fig. 11. In the no-adaptation case, the misclassified data scatter over the matrix. After adaptation, the misclassified data are evidently corrected by DIFO-C-B32 at the right side of Fig. 11. It is seen that DIFO-C-B32 improves performance on all categories, and on some categories achieving significant growth. For instance, in the second category, the performance promotes by **68**% (from **21**% to **89**%).

**Training stability.** Training stability is a vital characteristic of supervised learning methods. Based on the large-size dataset **VisDA**, we present the adaptation details of DIFO-C-B32 using the accuracy varying curves on the target domain. For comparison, the curves of typical self-supervised methods, SHOT, SCLM and TPDS, are also depicted. As shown in Fig. 12, the accuracy gradually increases to the maximum. This result confirms the training stability of DIFO-C-B32. Also, DIFO-C-B32 converges much faster than SHOT, SCLM and TPDS. It indicates that introducing task-specific knowledge from the ViL model is helpful in boosting the source model adaptation.

## References

[1] Alex Andonian, Shixing Chen, and Raffay Hamid. Robust cross-modal representation learning with progressive self-distillation. In *CVPR*, 2022. 2

[2] Dian Chen, Dequan Wang, Trevor Darrell, and Sayna Ebrahimi. Contrastive test-time adaptation. In *CVPR*, 2022. 5, 6, 2

[3] Weijie Chen, Luojun Lin, Shicai Yang, Di Xie, Shiliang Pu, and Yueting Zhuang. Self-supervised noisy label learning for source-free unsupervised domain adaptation. In *IROS*, 2022. 2

[4] Junhyeong Cho, Gilhyun Nam, Sungyeon Kim, Hunmin Yang, and Suha Kwak. Promptstyler: Prompt-driven style generation for source-free domain generalization. In *ICCV*, 2023. 2

[5] Ning Ding, Yixing Xu, Yehui Tang, Chao Xu, Yunhe Wang, and Dacheng Tao. Source-free domain adaptation via distribution estimation. In *CVPR*, 2022. 1, 2

[6] Yuntao Du, Haiyang Yang, Mingcai Chen, Juan Jiang, Hongtao Luo, and Chongjun Wang. Generation, augmentation, and alignment: A pseudo-source domain based method for source-free domain adaptation. *arXiv:2109.04015*, 2021. 2

[7] Shinto Eguchi and John Copas. Interpreting kullback–leibler divergence with the neyman–pearson lemma. *Journal of Multivariate Analysis*, 97(9):2034–2040, 2006. 1

[8] Yaroslav Ganin and Victor S. Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*, 2015. 1

[9] Chunjiang Ge, Rui Huang, Mixue Xie, Zihang Lai, Shiji Song, Shuang Li, and Gao Huang. Domain adaptation via prompt learning. *IEEE Transactions on Neural Networks and Learning Systems*, 2023. 2, 5, 6

[10] Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on vision transformer. *IEEE transactions on pattern analysis and machine intelligence*, 45 (1):87–110, 2022. 5

[11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5, 1

[12] Jiaxing Huang, Dayan Guan, Aoran Xiao, and Shijian Lu. Model adaptation: Historical contrastive learning for unsupervised domain adaptation without source data. In *NeurIPS*, 2021. 1, 5, 7, 3

[13] Xu Ji, Joao F Henriques, and Andrea Vedaldi. Invariant information clustering for unsupervised image classification and segmentation. In *CVPR*, 2019. 3

[14] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021. 2

[15] Guoliang Kang, Lu Jiang, Yi Yang, and Alexander G Hauptmann. Contrastive adaptation network for unsupervised domain adaptation. In *CVPR*, 2019. 1

[16] Jogendra Nath Kundu, Akshay R Kulkarni, Suvaansh Bhambri, Deepesh Mehta, Shreyas Anand Kulkarni, Varun Jampani, and Venkatesh Babu Radhakrishnan. Balancing discriminability and transferability for source-free domain adaptation. In *ICML*, 2022. 1

[17] Vinod K Kurmi, Venkatesh K Subramanian, and Vinay P Namboodiri. Domain impression: A source data free domain adaptation method. In *WACV*, 2021. 1

[18] Zhengfeng Lai, Noranart Vesdapunt, Ning Zhou, Jun Wu, Cong Phuoc Huynh, Xuelu Li, Kah Kuen Fu, and Chen-Nee Chuah. Padclip: Pseudo-labeling with adaptive debiasing in clip for unsupervised domain adaptation. In *ICCV*, 2023. 5, 6, 2

[19] Qicheng Lao, Xiang Jiang, and Mohammad Havaei. Hypothesis disparity regularized mutual information maximization. In *AAAI*, 2021. 1, 2

[20] Jonghyun Lee, Dahuin Jung, Junho Yim, and Sungroh Yoon. Confidence score for source-free unsupervised domain adaptation. In *ICML*, 2022. 5, 6, 7, 2, 3

[21] Jingjing Li, Zhekai Du, Lei Zhu, Zhengming Ding, Ke Lu, and Heng Tao Shen. Divergence-agnostic unsupervised domain adaptation by adversarial attacks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):8196–8211, 2021. 2

[22] Rui Li, Qianfen Jiao, Wenming Cao, Hau-San Wong, and Si Wu. Model adaptation: Unsupervised domain adaptation without ource data. In *CVPR*, 2020. 2

[23] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *CVPR*, 2023. 2

[24] J. Liang, R. He, Z. Sun, and T. Tan. Distant supervised centroid shift: A simple and efficient approach to visual domain adaptation. In *CVPR*, 2019. 1

[25] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In *ICML*, 2020. 2, 5, 6, 7, 1, 3

[26] Mattia Litrico, Alessio Del Bue, and Pietro Morerio. Guiding pseudo-labels with uncertainty estimation for source-free unsupervised domain adaptation. In *CVPR*, 2023. 5, 6, 2

[27] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I. Jordan. Conditional adversarial domain adaptation. In *NeurIPS*, 2018. 1

[28] R. Müller, S. Kornblith, and G. E Hinton. When does label smoothing help? In *NeurIPS*, 2019. 1

[29] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. Visda: The visual domain adaptation challenge. *arXiv:1710.06924*, 2017. 5, 1

[30] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *ICCV*, 2019. 5, 1

[31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 2, 5, 6, 1

[32] Subhankar Roy, Martin Trapp, Andrea Pilzer, Juho Kannala, Nicu Sebe, Elisa Ricci, and Arno Solin. Uncertainty-guided source-free domain adaptation. In *ECCV*, 2022. 1

[33] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *ECCV*, 2010. 5, 1

[34] Kuniaki Saito, Donghyun Kim, Stan Sclaroff, Trevor Darrell, and Kate Saenko. Semi-supervised domain adaptation via minimax entropy. In *ICCV*, 2019. 1

[35] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Gradcam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, 2017. 3

[36] Mainak Singha, Harsh Pal, Ankit Jha, and Biplab Banerjee. Ad-clip: Adapting domains in prompt space using clip. In *ICCV Workshop*, 2023. 5, 6, 2

[37] Song Tang, Mao Ye, Pei Xu, and Xudong Li. Adaptive pedestrian detection by predicting classifier. *Neural Computing and Applications*, 31:1189–1200, 2019. 2

[38] S Tang, Yuji Shi, Zhiyuan Ma, Jian Li, Jianzhi Lyu, Qingdu Li, and Jianwei Zhang. Model adaptation through hypothesis transfer with gradual knowledge distillation. In *IROS*, 2021. 5, 6, 2

[39] Song Tang, Yan Yang, Zhiyuan Ma, Norman Hendrich, Fanyu Zeng, Shuzhi Sam Ge, Changshui Zhang, and Jianwei Zhang. Nearest neighborhood-based deep clustering for source data-absent unsupervised domain adaptation. *arXiv:2107.12585*, 2021. 2

[40] Song Tang, Yan Zou, Zihao Song, Jianzhi Lyu, Lijuan Chen, Mao Ye, Shouming Zhong, and Jianwei Zhang. Semantic consistency learning on manifold for source data-free unsupervised domain adaptation. *Neural Networks*, 152, 2022. 1, 2, 5, 6

[41] Song Tang, An Chang, Fabian Zhang, Xiatian Zhu, Mao Ye, and Changshui Zhang. Source-free domain adaptation via target prediction distribution searching. *International Journal of Computer Vision*, pages 1–19, 2023. 4, 5, 6, 2

[42] Korawat Tanwisuth, Xinjie Fan, Huangjie Zheng, Shujian Zhang, Hao Zhang, Bo Chen, and Mingyuan Zhou. A prototype-oriented framework for unsupervised domain adaptation. In *NeurIPS*, 2021. 2

[43] Jiayi Tian, Jing Zhang, Wen Li, and Dong Xu. Vdm-da: Virtual domain modeling for source data-free domain adaptation. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(6):3749–3760, 2021. 1

[44] Jiayi Tian, Jing Zhang, Wen Li, and Dong Xu. Vdm-da: Virtual domain modeling for source data-free domain adaptation. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(6):3749–3760, 2021. 2

[45] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *CVPR*, 2017. 5, 1

[46] Fan Wang, Zhongyi Han, Yongshun Gong, and Yilong Yin. Exploring domain-invariant parameters for source free domain adaptation. In *CVPR*, 2022. 1

[47] Haifeng Xia, Handong Zhao, and Zhengming Ding. Adaptive adversarial network for source-free domain adaptation. In *CVPR*, 2021. 1

[48] Ruijia Xu, Guanbin Li, Jihan Yang, and Liang Lin. Larger norm more transferable: An adaptive feature norm approach for unsupervised domain adaptation. In *CVPR*, 2019. 1

[49] Shiqi Yang, Joost van de Weijer, Luis Herranz, Shangling Jui, et al. Exploiting the intrinsic neighborhood structure for source-free domain adaptation. In *NeurIPS*, 2021. 1, 4, 5, 6, 2

[50] Shiqi Yang, Yaxing Wang, Kai Wang, Shangling Jui, et al. Attracting and dispersing: A simple approach for source-free domain adaptation. In *NeurIPS*, 2022. 1, 5, 6, 7, 2, 3

[51] Li Yi, Gezheng Xu, Pengcheng Xu, Jiaqi Li, Ruizhi Pu, Charles Ling, A Ian McLeod, and Boyu Wang. When source-free domain adaptation meets learning with noisy labels. In *ICLR*, 2023. 5, 6, 2

[52] Yixin Zhang, Zilei Wang, and Weinan He. Class relationship embedded learning for source-free unsupervised domain adaptation. In *CVPR*, 2023. 5, 7, 3

[53] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022. 2

[54] ZongxianLee. A pytorch implementation of maximum mean discrepancies (MMD) loss. https://github.com/ZongxianLee/MMD_Loss.Pytorch, 2019. 7