# Domain-Specificity Inducing Transformers for Source-Free Domain Adaptation

Sunandini Sanyal* Ashish Ramayee Asokan* Suvaansh Bhambri* Akshay Kulkarni
Jogendra Nath Kundu R Venkatesh Babu
Vision and AI Lab, Indian Institute of Science, Bengaluru

## Abstract

*Conventional Domain Adaptation (DA) methods aim to learn domain-invariant feature representations to improve the target adaptation performance. However, we motivate that domain-specificity is equally important since in-domain trained models hold crucial domain-specific properties that are beneficial for adaptation. Hence, we propose to build a framework that supports disentanglement and learning of domain-specific factors and task-specific factors in a unified model. Motivated by the success of vision transformers in several multi-modal vision problems, we find that queries could be leveraged to extract the domain-specific factors. Hence, we propose a novel Domain-Specificity inducing Transformer (DSiT) framework [1] for disentangling and learning both domain-specific and task-specific factors. To achieve disentanglement, we propose to construct novel Domain-Representative Inputs (DRI) with domain-specific information to train a domain classifier with a novel domain token. We are the first to utilize vision transformers for domain adaptation in a privacy-oriented source-free setting, and our approach achieves state-of-the-art performance on single-source, multi-source, and multi-target benchmarks.*

## 1. Introduction

Machine learning models often fail to generalize to unseen domains due to the discrepancy between training (source) and test (target) data distributions (*i.e. domain shift*). This results in poor deployment performance and is also critical for applications like autonomous driving [10] or medical imaging [13]. Unsupervised domain adaptation (DA) techniques seek to address this challenging scenario by transferring task-specific knowledge from a labeled source domain to an unlabeled target domain. However, DA works [16] require concurrent access to source and target data. Such a constraint is highly impractical since data tends to be proprietary and cannot be easily shared. Hence, we focus on Source-Free DA [36] that operates under a practical

---

*Equal Contribution
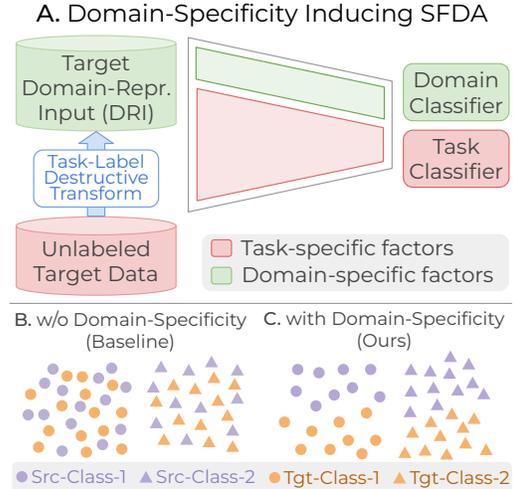[1] Project Page: http://val.cds.iisc.ac.in/DSiT-SFDA/



Figure 1. **A.** We induce domain-specificity by disentangling domain- and task-specific factors within the model. A task-label-destructive transform produces novel Domain-Representative Inputs (DRI) to learn domain-specific factors via domain classification. **B.** Conventional DA methods preserve domain-invariance, resulting in only task-oriented clusters in the feature space. **C.** Our proposed disentanglement ensures that different domains are well-clustered.

setting of sharing only the source model between a *vendor* and a *client*, without any data-sharing.

Conventional Domain Adaptation works [16] aim to learn task-discriminative features that are domain-invariant, *i.e.* indistinguishable w.r.t. domain shift. Intuitively, if the feature distributions of the source and target domain are similar, then a low error on the source domain translates to the target, shown theoretically by [3]. But domain-invariance does not always result in optimal target performance because a specific domain might be far away from the support of a domain-invariant model [15]. Further, supervised in-domain trained models (where train and test datasets come from the same domain) usually perform better as they hold useful domain-specific properties. Thus, we motivate the concept of *domain-specificity to improve the target adaptation performance*.

In source-free DA, while adapting a model to a new

target domain, the major problem is to preserve the task knowledge from the source domain. Prior works address this by aligning the target model feature space with that of the source [42]. However, we argue that it is equally important for a model to learn domain-specific information while preserving task-specific knowledge. Hence, in our work, we seek a solution to an important question, *"How do we develop a framework that enables us to improve domain-specificity while also retaining task-specificity?"*.

Thus, we seek a framework that supports the disentanglement of task-specific factors and domain-specific factors, thereby providing us better control over them. A well-disentangled framework would allow the learning of both domain-specific and task-specific factors in the model simultaneously, as shown in Fig. 1A. This not only yields better performance but also ensures that the different domains are well-clustered in the feature space of the model (Fig. 1C) compared to the baseline approach (Fig. 1B). However, the key question remains, *"How do we devise a method that encourages the disentanglement and learning of domain-specific and task-specific factors for better adaptation?"*

Recently, vision transformers have demonstrated remarkable performance in several vision tasks [31, 44]. They contain a multi-head self-attention mechanism that attends to all image patches and provides a global context. Motivated by this fundamental difference of a global context, we explore the possibility of a disentanglement framework with transformers since the domain information is inherently a global, higher-order statistic [8] that may not be adequately captured in CNNs. Inspired by the multi-modal works [23] that utilize queries to extract domain-specific information from a particular modality, we propose to enable the disentanglement of domain-specific and task-specific features through the query weights as part of our novel *Domain-Specificity inducing Transformer* (DSiT) framework.

Concretely, we induce domain-specificity by updating only the query weights via domain classifier training (Fig.1A). The remaining weights are updated via task classifier training. To further inculcate the disentanglement, we train the domain-classifier with novel *Domain-Representative Inputs* (DRI) where a task-label-destructive transform removes the task-specific information. We also propose a novel *domain-specificity disentanglement criterion* to evaluate the disentanglement of domain-specific and task-specific factors and demonstrate the disentanglement induced by our proposed framework, DSiT. We outline the contributions of our work as follows:

- We investigate and provide insights on how domain-specificity can be leveraged to improve DA. To this end, we propose a novel, unified Domain-specificity-inducing Transformer (DSiT) to disentangle and learn task-specific and domain-specific factors.

- We utilize query weights to enable the disentanglement in DSiT with a novel training algorithm that well supports our insights. We also introduce novel Domain-Representative Inputs (DRI) to further enhance the disentanglement.

- We define a novel domain-specificity disentanglement criterion to determine if the domain-specific and task-specific factors are well-disentangled.

- We achieve state-of-the-art performance on source-free benchmarks across single-source, multi-source, and multi-target DA while also introducing the first source-free DA benchmarks for transformers.

## 2. Related Works

**Domain Adaptation.** Prior DA works minimize the domain gap via adversarial distribution matching [20, 22] or minimize statistical distances [49, 55, 56, 74]. However, SFDA works [32, 33, 34, 42, 43] use information maximization and pseudo-labeling to match the target features with the source features. [39, 71] perform neighborhood clustering and regularization in the target domain. Our work also considers an SFDA setup [35, 36] that focuses on a practical vendor-client setting where vendor and client may use cooperative/same learning strategies, without data sharing.

**Transformer-based DA.** Vision transformers are self-attention-based architectures with state-of-the-art performance on several vision tasks like object recognition [5], semantic segmentation [21], *etc*. The application of transformers in DA scenarios is relatively less explored. CD-Trans [68] proposes cross-attention between source and target image pairs for domain alignment. SSRT [58] uses a self-training mechanism while [70] does adversarial training. TransDA [69] incorporates a self-attention layer on top of a ResNet backbone for SFDA, and hence is not a pure Vision Transformer (ViT) based solution. Unlike prior SFDA works, we propose a ViT-based solution for the first time, with a focus on enhancing domain-specificity.

**Domain-specificity for DA.** Prior works such as [37, 38] focus on style and content disentanglement and aim to achieve domain invariance by removing the style (domain-specific) information from images while preserving only the content information. However, we argue that domain invariance doesn't guarantee optimal performance. We draw motivation from domain generalization works [15] that propose constructing domain prototypes from unlabeled target samples to learn domain-specific features. Similarly, DiDA [4] proposes feature disentanglement into common and domain-specific features to improve adaptation performance. Furthermore, DMG [7] proposed to balance specificity and invariance via learning domain-specific masks and [6] proposed to learn domain-specific batch-

normalization parameters. Inspired by these, we propose to leverage domain-specificity in transformers to improve the target adaptation performance in a source-free DA setting.

## 3. Approach

**Problem setup.** For closed-set DA, we consider a labeled source dataset $\mathcal{D}_s = \{(x_s, y_s) : x_s \in \mathcal{X}, y_s \in \mathcal{C}_g\}$ where $\mathcal{X}$ denotes the input space and $\mathcal{C}_g$ denotes the goal task label set. We denote the unlabeled target dataset by $\mathcal{D}_t = \{x_t : x_t \in \mathcal{X}\}$. The task is to predict the label for each target sample $x_t$ from the label set $\mathcal{C}_g$. Following [68], we use ViT-B [12] as the backbone feature extractor $h : \mathcal{X} \to \{\mathcal{Z}_c, \mathcal{Z}_1, \mathcal{Z}_2, ...\mathcal{Z}_{N_P}\}$ where $\mathcal{Z}_c$ represents the class token feature-space and $\mathcal{Z}_1, \mathcal{Z}_2, ...\mathcal{Z}_{N_P}$ are patch token feature-spaces ($N_P$ is the number of patches). A goal task classifier trained on the class token is denoted as $f_g : \mathcal{Z}_c \to \mathcal{C}_g$. In the source-free vendor-client setup [36], we operate under the practical constraint where data sharing between vendor and client is prohibited. The vendor trains a model on the source domain and shares only the model with the client. The client uses the vendor-side model to adapt to the target domain. Unlike usual DA works [16] that advocate domain-invariant learning, we investigate how domain-specificity can be leveraged for DA. First, we discuss why domain-specificity is useful through the following insight.

**Insight 1. (Domain-specificity leads to improved DA)** *In the source-free DA setting, the aim is to achieve good target accuracy for the client-side model. Since an in-domain trained model better represents domain-specific factors as compared to a domain-invariant model, inculcating domain specificity on the client-side plays a major role in improving the target adaptation performance.*

**Remarks.** In a vendor-client setting, a client wishes to improve the goal task performance of the model on the target data. Conventional DA methods [16] often devise a strategy to learn domain-invariant features that contain only the task knowledge, which could be used to generalize across multiple domains. However, the optimal task-specific features from a domain-specific model might be far away from the task-specific features of a domain-invariant model [15]. Hence, such approaches are unsuitable for source-free DA, where the goal of improved target performance can be better achieved if we inculcate in-domain knowledge to build a domain-specific model. Therefore, we note that it is equally important for a client-side model to learn "domain-specific factors" containing the crucial in-domain characteristic knowledge along with the "task-specific factors" that hold information on the goal task.

As Insight 1 motivates domain-specificity, a natural question arises - "*How can we enable and control domain-specificity while also learning the task-specific factors?*"

**Insight 2. (Disentanglement of domain-specific and task-specific factors to control domain specificity)** *Disentanglement of domain-specific and task-specific factors provides a way to learn the two orthogonal factors together within the same model, enabling better control over them.*

**Remarks.** Since source and target domains can be far apart, *i.e.* high variance in domain-specific factors, it is desirable for a domain-specific model to learn a *disentangled* feature space that can adequately capture task-specific as well as domain-specific factors. Prior SFDA works either aim to align the target domain features towards the source domain [42] or update all parameters, risking the loss of source-side task-specific knowledge [71]. Conversely, our proposed disentanglement allows simultaneous and separate learning of the domain-specific and task-specific factors. Moreover, these domain-specific factors could be leveraged to orient the task-specific factors [15] for a new target domain.

Next, we seek an answer to the following question, "*How do we devise a unified architecture to enable disentanglement of domain-specific and task-specific factors?*"

### 3.1. Exploring the potential of transformers towards domain specificity

We explore the various possibilities of disentangling domain-specific and task-specific factors in the self-attention module of a transformer. Prior vision transformer works have demonstrated superior feature alignment capabilities within the attention module of a transformer, even across multi-modal scenarios such as text-to-vision [23, 61]. Moreover, transformers are inherently more robust to noise [68] and capture long-range dependencies via self-attention across different patches compared to the localized context that is captured in CNNs. Hence, we propose to adopt transformer architectures for inducing domain-specificity and are the first to use them in a source-free DA setting.

**a) Query for domain-specificity.** Vaswani et al. [63] introduced self-attention as a compatibility function between a Query and a corresponding Key. The final output is computed as a weighted sum of the Value where the weights are assigned as per the self-attention. Prior multi-modal works [23, 41] use modality-specific queries (e.g. from text or image modality) as input to the cross-attention. While we draw inspiration from these works, our novel approach explicitly enforces domain-specificity in the "self-attention query" (*i.e.* inside every transformer layer), different from the cross-attention approach of multi-modal works. This helps to better control the task-specific factors, where a domain-specific Query learns to acquire the domain-specific information from the input Value, making the overall output domain-specific in nature. While there could be more sophisticated ways to inculcate domain-specificity, such as hypernetworks [18] or auxiliary models [25], we propose this simple strategy that validates our insights.
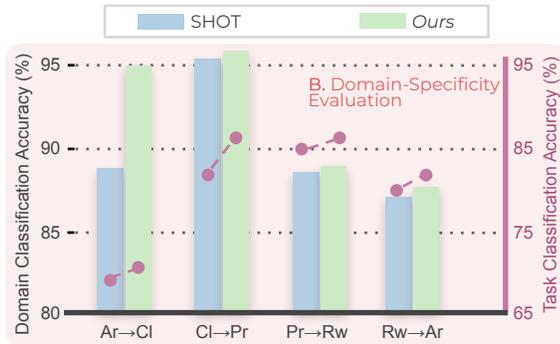
Figure 2. We evaluate the domain-specificity of our adapted model w.r.t. SHOT on Office-Home. Task acc. is goal task accuracy and Domain acc. is the accuracy of the binary source-target domain classifier (as explained in Sec. 3.1b). Our higher task accuracy (in pink) indicates better task-specificity and better disentanglement.

**b) Evaluating domain-specificity and task-specificity.** To analyze the effectiveness of our proposed strategy, we examine the domain-specificity of the ViT class token by training a linear domain classifier with the source and target domain samples. The accuracy is computed on four settings Ar→Cl, Cl→Pr, Pr→Rw, and Rw→Ar on the Office-Home dataset. From Fig. 2, we observe that our approach achieves a higher domain classification accuracy across all four settings. Note that we access both source and target domain data only for this analysis. Intuitively, we examine the domain-specificity inculcated in the class tokens as they are used by the task classifier. We also report the task performance to highlight the task-specific knowledge of the model in Fig. 2. We observe that our approach achieves higher goal task accuracy as well. This implies that our proposed approach captures the domain-specific and task-specific factors well, leading to better adaptation performance.

### 3.2. Training algorithm

We propose **D**omain-**S**pecificity inducing **T**ransformers for Source-Free DA (**DSiT**). In the following sections, we describe the proposed approach to train DSiT.

#### 3.2.1 Vendor-side source training

We train DSiT in two steps. First, we perform domain-specificity disentanglement via domain classifier training. The second step involves goal task classifier training. Next, we delve into the details of each step.

**a) Domain-specificity disentanglement** (Fig. 3B) We first induce domain-specificity by training a domain classifier $f_d$ with a novel domain token $z_d \in \mathcal{Z}_d$ from the backbone $h$. For this, the vendor first prepares augmented datasets $\mathcal{D}_s^{(i)} = \{(x_s^{[i]}, y_s, y_d) \ \forall \ (x_s, y_s) \in \mathcal{D}_s\} \ \forall \ i \in [N_a]$ by augmenting each source sample $x_s$ ($N_a$ is the number of augmentations). Here, the augmentation $\mathcal{A}_i : \mathcal{X} \to \mathcal{X}$ is applied to get $x_s^{[i]} = \mathcal{A}_i(x_s)$. Each input is assigned a domain label $y_d = i$ where $i$ denotes the augmentation label. We

use five label-preserving augmentations that simulate novel domains [35] (see Suppl. for more details).

We motivate the use of augmentations with two supports. First, as we operate in a source-free setting, access to multiple domains cannot be assumed, and these augmented domains can be used to inculcate domain-specificity. Further, this domain-specificity inculcation can be shared between vendor and client by sharing only the augmentation information. Second, as shown in Fig. 4, an augmented source domain may be closer to another augmented target domain (blue lines) w.r.t. the original source-target (red line). We also confirm this numerically (reported in Suppl). Thus, we have access to more diverse domains (with lower domain gaps) which need to be well-separated by the domain classifier $f_d$, allowing for better domain-specificity.

While we could train $f_d$ with only augmented samples to inculcate domain-specificity, we also need to disentangle the task-specific information as discussed under Insight 2. Thus, we next describe a novel input representation that better represents domain-specific factors by separating the task-specific factors, *i.e.* enabling better disentanglement.

**Domain-Representative Input (DRI)** is constructed at the input level to preserve only domain-specific information by applying a task-label-destructive transformation of patch-shuffling [46] as shown in Fig. 3A. Intuitively, domain information is a higher-order statistic [8] that is preserved after patch-shuffling while class information is lost. The shuffling of patches is extremely vital for amplifying the input domain-specific factors for the domain-specificity training step. Hence, DRIs are not just used as an augmentation for the goal task training, as shown in prior transformer works [52], but are a means to inculcate domain specificity. We empirically demonstrate in Sec. 4.2b that DRI guides the goal task performance better than only augmented inputs, because of improved disentanglement. We also provide additional baselines in Table 5 to validate the unsuitability of DRI as a general augmentation for the goal task training. As per our analysis in Sec. 3.1, we train only the transformer query weights, to exclusively hold domain-specific knowledge, with the domain classification loss,

$$\min_{\theta_Q, \theta_{f_d}} \mathbb{E}_{(x,y_d) \in \cup_i \mathcal{D}_s^{(i)}} [\mathcal{L}_{dom}] \text{ where } \mathcal{L}_{dom} = \mathcal{L}_{ce}(f_d(z_d), y_d)$$

(1)

where $z_d$ is the domain-token output from backbone $h$ and $\theta_Q = \cup_j W_{Q_j}$; $j$ is an index over the backbone layers. Note that DRI is used for $\mathcal{L}_{dom}$. For simplicity, we re-use $\mathcal{D}_s^{(i)}$ to include the task-label-destructive transform (Fig. 3A).

**b) Domain-specific goal task training** (Fig. 3C) After one round of domain classifier training, we train DSiT for the goal task. Here, we update key and value weights $W_K$ and $W_V$ along with all other parameters (except query weights $W_Q$). Hence, the frozen queries are disentangled
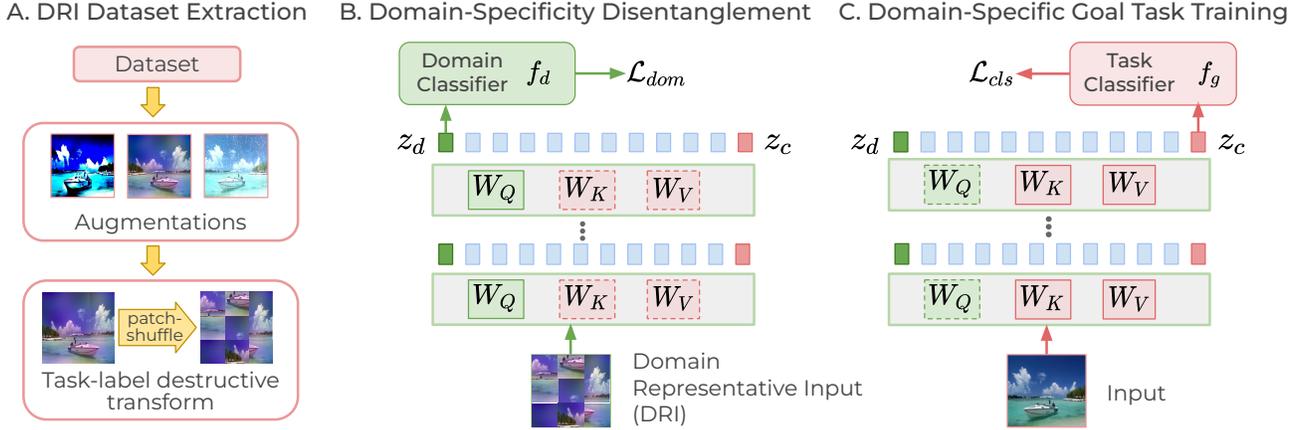
**A. DRI Dataset Extraction** **B. Domain-Specificity Disentanglement** **C. Domain-Specific Goal Task Training**

Figure 3. **DSiT Training: A. DRI dataset extraction:** DRI are obtained via patch-shuffling of augmented inputs. **B. Domain-Specificity Disentanglement:** Domain classifier $f_d$ is trained with domain classification loss $\mathcal{L}_{dom}$ using DRI, updating only query weights $W_Q$ via a domain token $z_d$, while other weights including $W_K$ and $W_V$ are frozen (dotted boundary). **C. Domain-Specific Goal Task Training:** The task classifier is trained with task classification loss $\mathcal{L}_{cls}$ updating all weights except $W_Q$.

and hold crucial domain-specific information. The vendor trains a source model consisting of the backbone $h$ and task classifier $f_g$ using the source dataset $\mathcal{D}_s$ with the task-classification loss as follows,

$$\min_{\theta_h \backslash \theta_Q, \theta_{f_g}} \mathbb{E}_{(x_s, y_s) \in \mathcal{D}_s} [\mathcal{L}_{cls}] \text{ where } \mathcal{L}_{cls} = \mathcal{L}_{ce}(f_g(z_c), y_c)$$
(2)

where $z_c$ is the class-token output from backbone $h$, and $\theta_h \backslash \theta_Q$ are the parameters of the backbone $h$ excluding all $W_Q$ weights while $\theta_{f_g}$ are the parameters of classifier $f_g$. The two steps of domain-specificity disentanglement (Eq. 1) and goal task training (Eq. 2) are performed in tandem, one after the other. See Suppl. for details.

### 3.2.2 Client-side target adaptation

The vendor shares the source-trained, disentangled DSiT model with the client for target adaptation. The client also follows the same process of domain-specificity disentanglement as the vendor. To disentangle the target-domain-specific factors and task-specific factors for improved target adaptation (as per Insight 1), the client applies the same augmentations to the target data to generate DRIs for domain classifier training with the domain classification loss $\mathcal{L}_{dom}$ (as described in Sec. 3.2.1a and Fig. 3B). For goal task training, the client uses the standard information maximization and diversity losses [42, 71, 72]. The overall client-side objective is as follows:

$$\min_{\theta_h \backslash \theta_Q, \theta_{f_g}} \mathbb{E}_{\mathcal{D}_t} [\mathcal{L}_{im} + \mathcal{L}_{div}] + \min_{\theta_Q, \theta_{f_d}} \mathbb{E}_{\cup_i \mathcal{D}_t^{(i)}} [\mathcal{L}_{dom}]$$
(3)

We detail $\mathcal{L}_{im}$ and $\mathcal{L}_{div}$ in the Suppl. as these are commonly used and not our main contributions. Note that only original target data $\mathcal{D}_t$ is used for $\mathcal{L}_{im}$ and $\mathcal{L}_{div}$. As in Eq. 1, DRI



Figure 4. Novel augmented domains can be simulated where the domain gap between an augmented source and an augmented target (blue lines) can be less than the original domain gap (red line). With a domain classifier for augmented domains, we induce better domain-specificity as multiple domains with lower domain gaps are well-separated.

are used for $\mathcal{L}_{dom}$ here. For simplicity, we re-use $\mathcal{D}_t^{(i)}$ to include the task-label-destructive transform (Fig. 3A).

### 3.3. Domain-specificity disentanglement criterion

To evaluate the disentanglement of domain-specific and task-specific factors, we propose a criterion considering the following three feature-space cosine similarity metrics. Let $\gamma_{cls}$ be the intra-class, inter-domain similarity, $\gamma_{dom}$ be the intra-domain, inter-class similarity, and $\gamma_{all}$ be the inter-class, inter-domain similarity. These are computed as the similarity between pairs of features averaged over the possible input pairs for each input. For example, possible input pairs for $\gamma_{cls}$ have to come from the same class but different domains. See Suppl. for more details. Based on these, we define a criterion to determine if the domain-specific and task-specific factors of the model are well-disentangled.

Table 1. Single-Source Domain Adaptation (SSDA) on Office-Home benchmark. SF indicates *source-free* adaptation. ResNet-based methods (top) and Transformer-based methods (bottom). * indicates results taken from CDTrans [68]. **Bold** numbers indicate the best results among SFDA methods. <span style="color:green">Green</span> results indicate our method's improvements over the corresponding methods.

| Method | SF | Office-Home | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 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 |
| ResNet-50 [19] | ✗ | 34.9 | 50.0 | 58.0 | 37.4 | 41.9 | 46.2 | 38.5 | 31.2 | 60.4 | 53.9 | 41.2 | 59.9 | 46.1 |
| SENTRY [51] | ✗ | 61.8 | 77.4 | 80.1 | 66.3 | 71.6 | 74.7 | 66.8 | 63.0 | 80.9 | 74.0 | 66.3 | 84.1 | 72.2 |
| SCDA [40] | ✗ | 60.7 | 76.4 | 82.8 | 69.8 | 77.5 | 78.4 | 68.9 | 59.0 | 82.7 | 74.9 | 61.8 | 84.5 | 73.1 |
| A$^2$Net [67] | ✓ | 58.4 | 79.0 | 82.4 | 67.5 | 79.3 | 78.9 | 68.0 | 56.2 | 82.9 | 74.1 | 60.5 | 85.0 | 72.8 |
| GSFDA [72] | ✓ | 57.9 | 78.6 | 81.0 | 66.7 | 77.2 | 77.2 | 65.6 | 56.0 | 82.2 | 72.0 | 57.8 | 83.4 | 71.3 |
| CPGA [53] | ✓ | 59.3 | 78.1 | 79.8 | 65.4 | 75.5 | 76.4 | 65.7 | 58.0 | 81.0 | 72.0 | 64.4 | 83.3 | 71.6 |
| NRC [71] | ✓ | 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 |
| SHOT [42] | ✓ | 57.1 | 78.1 | 81.5 | 68.0 | 78.2 | 78.1 | 67.4 | 54.9 | 82.2 | 73.3 | 58.8 | 84.3 | 71.8 |
| SHOT++ [43] | ✓ | 57.9 | 79.7 | 82.5 | 68.5 | 79.6 | 79.3 | 68.5 | 57.0 | 83.0 | 73.7 | 60.7 | 84.9 | 73.0 |
| TVT [70] | ✗ | 74.8 | 86.8 | 89.4 | 82.7 | 87.9 | 88.2 | 79.8 | 71.9 | 90.1 | 85.4 | 74.6 | 90.5 | 83.5 |
| SSRT-B [58] | ✗ | 75.1 | 88.9 | 91.0 | 85.1 | 88.2 | 89.9 | 85.0 | 74.2 | 91.2 | 85.7 | 78.5 | 91.7 | 85.4 |
| CDTrans [68] | ✗ | 68.8 | 85.0 | 86.9 | 81.5 | 87.1 | 87.3 | 79.6 | 63.3 | 88.2 | 82.0 | 66.0 | 90.6 | 80.5 |
| SHOT-B* [43] | ✓ | 67.1 | 83.5 | 85.5 | 76.6 | 83.4 | 83.7 | 76.3 | 65.3 | 85.3 | 80.4 | 66.7 | 83.4 | 78.1 <span style="color:green">(+2.4)</span> |
| DIPE [66] | ✓ | 66.0 | 80.6 | 85.6 | 77.1 | 83.5 | 83.4 | 75.3 | 63.3 | 85.1 | 81.6 | 67.7 | 89.6 | 78.2 <span style="color:green">(+2.3)</span> |
| Mixup [34] | ✓ | 65.3 | 82.1 | 86.5 | 77.3 | 81.7 | 82.4 | 77.1 | 65.7 | 84.6 | 81.2 | **70.1** | 88.3 | 78.5 <span style="color:green">(+2.0)</span> |
| DSiT-B (*Ours*) | ✓ | **69.2** | **83.5** | **87.3** | **80.7** | **86.1** | **86.2** | **77.9** | **67.9** | **86.6** | **82.4** | 68.3 | **89.8** | **80.5** |

**Definition 1. (Domain-Specificity Disentanglement Criterion)** *The domain-specific and task-specific factors of a model are well-disentangled if*

$$\gamma_{cls}, \gamma_{dom} > \gamma_{all} \text{ and } |\gamma_{cls} - \gamma_{dom}| < \tau \quad (4)$$

**Remarks.** Here, $\tau$ is a threshold. In other words, $\gamma_{cls}$ and $\gamma_{dom}$ should be higher than $\gamma_{all}$ and the absolute difference between $\gamma_{cls}$ and $\gamma_{dom}$ should be low. Based on the definitions, $\gamma_{cls}$ indicates task-specificity since the used samples belong to the same class but different domains. Similarly, $\gamma_{dom}$ represents the domain-specificity. $\gamma_{all}$ is the similarity between samples from different classes and different domains, which should be lower than $\gamma_{dom}$ and $\gamma_{cls}$ with any domain-specific and task-specific knowledge, respectively. Along with the first condition, if the difference between task-specificity and domain-specificity is low, *i.e.* both are equally present, the two will be well-disentangled in the model. In Sec. 4.2c, we empirically verify that our proposed approach satisfies Definition 1.

## 4. Experiments

We evaluate the effectiveness of our proposed approach on the standard DA benchmarks across several settings.

**Datasets.** We demonstrate the efficacy of our work on the following three standard object recognition DA benchmarks. **Office-31** [55] is a benchmark consisting of three domains under office environments - Amazon (**A**), DSLR (**D**), and Webcam (**W**) with 31 classes each. **Office-Home** [65] consists of images of everyday objects divided into four domains − Artistic (**Ar**), ClipArt (**Cl**), Product (**Pr**), and

Real-world (**Rw**), each with 65 classes. **VisDA** [50] is a large-scale dataset for adapting algorithms from synthetic domains to real domains. There are 152,397 synthetic images in the source domain and 55,388 real-world images in the target domain. **DomainNet** [49] is the most challenging benchmark with six domains, each with 345 classes: ClipArt (**C**), Real (**R**), Infograph (**I**), Painting (**P**), Sketch (**S**), and Quickdraw (**Q**).

**Implementation details.** We follow the experimental settings of CDTrans [68] and use DeiT-Base [60] as the backbone ViT for fair comparisons. The DeiT-Base architecture consists of 12 layers, where each layer consists of 12 self-attention heads collectively termed as multi-head self-attention. We use an input size of $224 \times 224$ and a patch size of $16 \times 16$ for all our experiments, resulting in $14 \times 14$ patches (196 total) per input. DeiT-B contains an additional distillation token, however, the rest of the architecture is the same as a ViT-B backbone. For optimizing the training objectives, we use Stochastic Gradient Descent (SGD) with momentum 0.9, and a weight decay ratio of $1 \times 10^{-4}$. Refer to the Suppl. for the complete implementation details.

### 4.1. Comparisons with prior works

We compare with SOTA methods on single-source, and multi-source DA in Table 1, 2, 3 and 4 and multi-target DA (see Suppl. for details).

**a) Single-source domain adaptation (SSDA).** Table 1 2, and 3 present comparisons between our proposed approach DSiT and prior SSDA works. For the Office-Home benchmark, our approach achieves state-of-the-art performance

Table 2. Single-Source Domain Adaptation (SSDA) on Office-31 and VisDA benchmarks. SF indicates *source-free* adaptation. ResNet-based methods (top) and Transformer-based methods (bottom). * indicates results taken from CDTrans [68]. **Bold** numbers indicate the best results among SFDA methods. <span style="color:green">Green</span> results indicate our method's improvements over the corresponding methods.

| Method | SF | Office-31 | | | | | | | VisDA |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | A→W | D→W | W→D | A→D | D→A | W→A | Avg. | S→R |
| ResNet-50 [19] | ✗ | 68.9 | 68.4 | 62.5 | 96.7 | 60.7 | 99.3 | 76.1 | 52.4 |
| CAN [30] | ✗ | 94.5 | 99.1 | 99.8 | 95.0 | 78.0 | 77.0 | 90.6 | 87.2 |
| FixBi [47] | ✗ | 96.1 | 99.3 | 100.0 | 95.0 | 78.7 | 79.4 | 91.4 | 87.2 |
| CDAN+RADA [28] | ✗ | 96.2 | 99.3 | 100.0 | 96.1 | 77.5 | 77.4 | 91.1 | 76.3 |
| CPGA [53] | ✓ | 94.1 | 98.4 | 99.8 | 94.4 | 76.0 | 76.6 | 89.9 | 84.1 |
| HCL [24] | ✓ | 91.3 | 98.2 | 100.0 | 90.8 | 72.7 | 72.7 | 87.6 | 83.5 |
| VDM-DA [59] | ✓ | 94.1 | 98.0 | 100.0 | 93.2 | 75.8 | 77.1 | 89.7 | 85.1 |
| A$^2$Net [67] | ✓ | 94.0 | 99.2 | 100.0 | 94.5 | 76.7 | 76.1 | 90.1 | 84.3 |
| NRC [71] | ✓ | 90.8 | 99.0 | 100.0 | 96.0 | 75.3 | 75.0 | 89.4 | 85.9 |
| SHOT [42] | ✓ | 90.1 | 98.4 | 99.9 | 94.0 | 74.7 | 74.3 | 88.6 | 82.9 |
| SHOT++ [43] | ✓ | 90.4 | 98.7 | 99.9 | 94.3 | 76.2 | 75.8 | 89.2 | 87.3 |
| TVT [70] | ✗ | 96.4 | 99.4 | 100.0 | 96.4 | 84.9 | 86.1 | 93.8 | 83.9 |
| CGDM-B* [14] | ✗ | 95.3 | 97.6 | 99.8 | 94.6 | 78.8 | 81.2 | 91.2 | 82.3 |
| CD-Trans [68] | ✗ | 96.7 | 99.0 | 100.0 | 97.0 | 81.1 | 81.9 | 92.6 | 88.4 |
| SSRT-B [58] | ✗ | 97.7 | 99.2 | 100.0 | 98.6 | 83.5 | 82.2 | 93.5 | 88.7 |
| DIPE [66] | ✓ | 95.5 | 98.5 | 100.0 | 94.8 | 77.5 | 77.1 | 90.5 <span style="color:green">(+2.5)</span> | 82.8 <span style="color:green">(+4.8)</span> |
| Mixup [34] | ✓ | 96.1 | 98.6 | 100.0 | 95.4 | 80.2 | 80.1 | 91.7 <span style="color:green">(+1.3)</span> | 86.3 <span style="color:green">(+1.3)</span> |
| SHOT-B* | ✓ | 94.3 | 99.0 | 100.0 | 95.3 | 79.4 | 80.2 | 91.4 <span style="color:green">(+1.6)</span> | 85.9 <span style="color:green">(+1.7)</span> |
| DSiT-B (*Ours*) | ✓ | **97.2** | **99.1** | **100.0** | **98.0** | **81.7** | **81.8** | **93.0** | **87.6** |

Table 3. Single-Source Domain Adaptation (SSDA) results on the DomainNet dataset. * indicates results taken from [58]. NSF indicates Non-Source-Free, SF indicates Source-Free and SO indicates Source-Only.

| SCDA [75] (NSF) | clp | inf | pnt | qdr | rel | skt | Avg. |
| --- | --- | --- | --- | --- | --- | --- | --- |
| clp | - | 20.4 | 43.3 | 15.2 | 59.3 | 46.5 | 36.9 |
| inf | 32.7 | - | 34.5 | 6.3 | 47.6 | 29.2 | 30.1 |
| pnt | 46.4 | 19.9 | - | 8.1 | 58.8 | 42.9 | 35.2 |
| qdr | 31.1 | 6.6 | 18.0 | - | 28.8 | 22.0 | 21.3 |
| rel | 55.5 | 23.7 | 52.9 | 9.5 | - | 45.2 | 37.4 |
| skt | 55.8 | 20.1 | 46.5 | 15.0 | 56.7 | - | 38.8 |
| Avg. | 44.3 | 18.1 | 39.0 | 10.8 | 50.2 | 37.2 | 33.3 |

| DeiT-B [60] (SO) | clp | inf | pnt | qdr | rel | skt | Avg. |
| --- | --- | --- | --- | --- | --- | --- | --- |
| clp | - | 24.3 | 49.6 | 15.8 | 65.3 | 52.1 | 41.4 |
| inf | 45.9 | - | 45.9 | 6.7 | 61.4 | 39.5 | 39.9 |
| pnt | 53.2 | 23.8 | - | 6.5 | 66.4 | 44.7 | 38.9 |
| qdr | 31.9 | 6.8 | 15.4 | - | 23.4 | 20.6 | 19.6 |
| rel | 59.0 | 25.8 | 56.3 | 9.16 | - | 44.8 | 39.0 |
| skt | 60.6 | 20.6 | 48.4 | 16.5 | 61.2 | - | 41.5 |
| Avg. | 50.1 | 20.3 | 43.1 | 10.9 | 55.5 | 40.3 | 36.7 |

| SHOT-B [42] (SF) | clp | inf | pnt | qdr | rel | skt | Avg. |
| --- | --- | --- | --- | --- | --- | --- | --- |
| clp | - | 27.0 | 49.7 | 16.5 | 65.4 | 53.2 | 46.1 |
| inf | 46.4 | - | 45.9 | 7.4 | 60.6 | 40.1 | 40.1 |
| pnt | 54.6 | 25.7 | - | 8.1 | 66.3 | 49.0 | 40.7 |
| qdr | 33.3 | 6.8 | 15.5 | - | 23.8 | 24.0 | 20.7 |
| rel | 59.3 | 28.1 | 57.4 | 9.0 | - | 47.3 | 40.2 |
| skt | 64.0 | 26.5 | 55.0 | 18.2 | 63.8 | - | 45.5 |
| Avg. | 51.5 | 26.6 | 44.7 | 11.8 | 56.0 | 42.7 | 38.9 |

| CDTrans* [68] (NSF) | clp | inf | pnt | qdr | rel | skt | Avg. |
| --- | --- | --- | --- | --- | --- | --- | --- |
| clp | - | 27.9 | 57.6 | 27.9 | 73.0 | 58.8 | 49.0 |
| inf | 58.6 | - | 53.4 | 9.6 | 71.1 | 47.6 | 48.1 |
| pnt | 60.7 | 24.0 | - | 13.0 | 69.8 | 49.6 | 43.4 |
| qdr | 2.9 | 0.4 | 0.3 | - | 0.7 | 4.7 | 1.8 |
| rel | 49.3 | 18.7 | 47.8 | 9.4 | - | 33.5 | 31.7 |
| skt | 66.8 | 23.7 | 54.6 | 27.5 | 68.0 | - | 48.1 |
| Avg. | 47.7 | 18.9 | 42.7 | 17.5 | 56.5 | 38.8 | 37.0 |

| SSRT-B* [58] (NSF) | clp | inf | pnt | qdr | rel | skt | Avg. |
| --- | --- | --- | --- | --- | --- | --- | --- |
| clp | - | 33.8 | 60.2 | 19.4 | 75.8 | 59.8 | 49.8 |
| inf | 55.5 | - | 54.0 | 9.0 | 68.2 | 44.7 | 46.3 |
| pnt | 61.7 | 28.5 | - | 8.4 | 71.4 | 55.2 | 45.0 |
| qdr | 42.5 | 8.8 | 24.2 | - | 37.6 | 33.6 | 29.3 |
| rel | 69.9 | 37.1 | 66.0 | 10.1 | - | 58.9 | 48.4 |
| skt | 70.6 | 32.8 | 62.2 | 21.7 | 73.2 | - | 52.1 |
| Avg. | 60.0 | 28.2 | 53.3 | 13.7 | 65.3 | 50.4 | 45.2 |

| DSiT (*Ours*) (SF) | clp | inf | pnt | qdr | rel | skt | Avg. |
| --- | --- | --- | --- | --- | --- | --- | --- |
| clp | - | 27.2 | 51.8 | 23.1 | 70.2 | 54.7 | 45.4 |
| inf | 52.3 | - | 48.8 | 12.8 | 68.3 | 44.2 | 45.3 |
| pnt | 59.2 | 26.1 | - | 14.5 | 71.5 | 51.4 | 44.5 |
| qdr | 38.1 | 8.3 | 21.2 | - | 37.2 | 27.6 | 26.5 |
| rel | 60.4 | 28.0 | 57.8 | 13.1 | - | 49.7 | 41.8 |
| skt | 66.3 | 27.5 | 56.0 | 24.4 | 70.2 | - | 48.9 |
| Avg. | 55.3 | 23.4 | 47.1 | 17.6 | 63.5 | 45.5 | 42.1 |

among source-free works (Table 1). DSiT outperforms the source-free prior works SHOT-B* by 2.4%, Mixup by 2%, and DIPE by 2.3% on Office-Home and yields comparable performance to the non-source-free transformer-based prior work CDTrans. Similarly, on the Office-31 benchmark (Table 2), our approach improves over the source-free works SHOT-B* by 1.6%, Mixup by 1.3% and DIPE by 2.5% and gives competitive results compared to non-source-free works. On the more challenging VisDA dataset (Table 2), our approach outperforms source-free SHOT-B* baseline by 1.7%. DSiT also outperforms all the existing source-

free methods by a significant margin on the most challenging benchmark DomainNet (Table 3), and also achieves a 5.1% improvement over CDTrans, despite the latter being non-source-free. Refer to Suppl. for the complete table.

**b) Multi-source domain adaptation (MSDA).** In Table 4, we compare our approach with source-free and non-source-free prior works for multi-source DA. We improve over the source-free baseline SHOT-B* by 1% on Office-Home. Despite the source-free constraint that we address with our method, we outperform even non-source-free works.

Overall, these benchmark results highlight the efficacy of

Table 4. Multi-Source DA (MSDA) on Office-Home benchmark. SF indicates *source-free* adaptation. ResNet-based methods (top) and Transformer-based methods (bottom).

| Method | SF | Office-Home | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | | →Ar | →Cl | →Pr | →Rw | Avg. |
| Source-combine | ✗ | 58.0 | 57.3 | 74.2 | 77.9 | 66.9 |
| WDA [1] | ✗ | 71.9 | 61.4 | 84.1 | 82.3 | 74.9 |
| SImpAl [64] | ✗ | 70.8 | 56.3 | 80.2 | 81.5 | 72.2 |
| CMSDA [57] | ✗ | 71.5 | 67.7 | 84.1 | 82.9 | 76.6 |
| DECIS [2] | ✓ | 74.5 | 59.4 | 84.4 | 83.6 | 75.5 |
| SHOT++ [43] | ✓ | 73.1 | 61.3 | 84.3 | 84.0 | 75.7 |
| CAiDA [11] | ✓ | 75.2 | 60.5 | 84.7 | 84.2 | 76.2 |
| SHOT [42] | ✓ | 72.2 | 59.3 | 82.8 | 82.9 | 74.3 |
| SHOT-B* | ✓ | 83.9 | 71.8 | 89.7 | 89.4 | 83.7 |
| DSiT-B (*Ours*) | ✓ | **84.4** | **73.8** | **90.7** | **89.7** | **84.7** (+1.0) |

Table 5. Ablation study for DRI as a general augmentation on Office-Home SSDA benchmark with the DeiT-B backbone.

| Method | Ar→Cl | Cl→Pr | Pr→Rw | Rw→Ar | Avg. |
| --- | --- | --- | --- | --- | --- |
| SHOT-B [42] | 67.1 | 83.4 | 85.3 | 80.4 | 79.1 |
| SHOT-B + DRI | 65.8 | 79.9 | 84.2 | 79.6 | 77.4 |
| CDTrans [68] | 68.8 | 87.1 | 88.2 | 82.0 | 81.5 |
| CDTrans + DRI | 59.7 | 81.6 | 84.2 | 81.5 | 76.8 |

Table 6. Ablation study for various stages of training on Office-Home SSDA benchmark with the DeiT-B backbone (*average over 4 settings*). SO: Source-Only, DST: Domain-Specificity Training

| Training Phase | Model | # | DST | DRI | Avg. |
| --- | --- | --- | --- | --- | --- |
| Vendor-side | SO | 1 | ✗ | ✗ | 74.8 |
| | DSiT | 2 | ✓ | ✗ | 75.6 (+0.8) |
| | | 3 | ✓ | ✓ | 76.6 (+1.8) |
| Client-side | SHOT | 4 | ✗ | ✗ | 79.0 |
| | DSiT | 5 | ✓ | ✗ | 80.0 (+1.0) |
| | | 6 | ✓ | ✓ | **81.1** (+2.1) |

Table 7. Empirical evaluation of similarity metrics shows that domain-specificity $\gamma_{dom}$ and task-specificity $\gamma_{cls}$ are closer for our DSiT along with better DA performance, indicating better disentanglement than the baseline.

| Office-Home | $\gamma_{cls}$ | $\gamma_{dom}$ | $\gamma_{all}$ | SSDA |
| --- | --- | --- | --- | --- |
| SHOT-B [42] | 0.84 | 0.74 | 0.73 | 78.1 |
| DSiT-B (*Ours*) | 0.81 | 0.78 | 0.71 | **80.5** |

the proposed disentanglement of domain-specific and task-specific factors (Insight 2) across three diverse settings of single-source, multi-source, and multi-target DA.

## 4.2. Analysis

In Table 6, we provide an ablation study for various steps in the training of our approach.

**a) Effect of inculcating domain-specificity.** In Table 6, we compare the effect of domain-specificity training (without DRI) on both vendor-side and client-side training. On vendor-side (#1 vs. #2), we observe improvements of 0.8% while the improvements on client-side (#4 vs. #5) are 1.0%. This is in line with our Insight 1 that domain-specificity supports target adaptation performance (even without DRI).

**b) Effect of DRI.** As discussed in Sec. 3.2, domain-representative inputs (DRI) reduce the impact of task-specific information while preserving and aiding in the learning of domain-specific information. This further amplifies the impact of domain-specificity on the goal task, which results in improvements of 1% on the vendor-side (#2 vs. #3) and 1.1% on the client-side (#5 vs. #6).

**c) Empirical analysis of Definition 1.** As per the discussion under Definition 1, we report the three similarity metrics for the baseline SHOT and our DSiT in Table 7. The condition of $\gamma_{dom}, \gamma_{cls} > \gamma_{all}$ is satisfied by both models (since $\gamma_{all}$ denotes similarity between any class and any domain samples which would be lower in most cases). $\gamma_{dom}$ and $\gamma_{cls}$ for DSiT are closer in magnitude than SHOT, indicating its better disentanglement. This is further strength-

ened by the improved DA performance of our proposed method DSiT compared to the source-free baseline SHOT.

**d) Performance of DRI as a general augmentation.** To preserve and enhance domain-specificity, we utilize DRIs only for the domain-classifier training (Eq. 1), which helps to improve the target adaptation performance (as per Insight 1). However, DRIs are unsuitable for task-classifier training as the task label information is destroyed by patch-shuffling. In Table 5, we use DRI-augmented inputs for task classifier training with SHOT-B [42] and CDTrans [68] baselines and observe significant drops of 1.7% and 4.7% on Office-Home (4 settings) with respect to the corresponding baselines, respectively. This drop is expected and validates our insights that DRI is not suitable as a general augmentation for DA, but is extremely crucial for domain-specificity training (+1% improvement using DRI in Table 6).

## 5. Conclusion

In this work, we study the concept of domain-specificity in source-free DA. We provide insights to analyze how domain-specificity could be leveraged to improve target adaptation performance. We, therefore, propose a novel Domain-specificity-inducing Transformer (DSiT) where we leverage the queries of a vision transformer to induce domain-specificity and train the unified model to enable a disentanglement of task- and domain-specificity. Based on our insights, we also introduce a novel Domain-Specificity Disentanglement criterion to determine if the task-specific and domain-specific factors are well disentangled. The proposed approach outperforms several state-of-the-art benchmarks for single-source, multi-source, and multi-target DA.

# Supplementary Material

In this supplementary material, we provide more details on the training algorithm, experimental settings, additional comparisons, and analysis experiments. We have released our code on our project page: `https://val.cds.iisc.ac.in/DSiT-SFDA/`. The remainder of the supplementary material is structured as shown below:

- Sec. A: Approach (Algorithm 1)
- Sec. B: Implementation Details
  - Sec. B.1: Domain Augmentations (Fig. 5A)
  - Sec. B.2: DRI Dataset Extraction (Fig. 5B)
  - Sec. B.3 Domain-specificity criterion
  - Sec. B.4: Experimental Settings
- Sec. C Additional experimental results
  - Sec. C.1 Extended comparisons (Table 9, 13)
  - Sec. C.2 Vendor-side DSiT results (Table 10)
  - Sec. C.3 Model adaptation setting (Table 11)
  - Sec. C.4 Perf. on different backbones (Table 12)
  - Sec. C.5 Analysis for augmentations (Table 15)
  - Sec. C.6: Sensitivity analysis of DSiT (Table 16)
  - Sec. C.7 Training time comparisons (Table 14)
  - Sec. C.8: Statistical significance (Table 17)
  - Sec. C.9: Target adaptation losses (Table 18)
  - Sec. C.10: Effect of DRI grid-size (Fig. 6)

## A. Approach

Table 8 shows a complete list of the notations used in the paper. We summarize our full approach in Algorithm 1 and describe the details of the approach in this section.

**(a) Target adaptation losses.** For the client-side target adaptation, we use the Information Maximization loss formulation from SHOT [42], which consists of two terms: entropy loss $\mathcal{L}_{im}$ and diversity loss $\mathcal{L}_{div}$. The entropy loss $\mathcal{L}_{im}$ ensures that the confidence of the model towards a label is high. The diversity loss $\mathcal{L}_{div}$ ensures that the model's predictions are well-balanced across all classes and prevents the model from producing degenerate solutions. We define the two terms as follows:

$$\mathcal{L}_{im} = - \mathop{\mathbb{E}}_{x_t \in \mathcal{X}} \sum_{k=1}^{K} \delta_k(f_g(z_c)) \log \delta_k(f_g(z_c)) \quad (5)$$

Table 8. List of all the notations used throughout the paper.

| | Symbol | Description |
|---|---|---|
| **Models** | $h$ | Backbone feature extractor |
| | $f_g$ | Goal task classifier |
| | $f_d$ | Domain classifier |
| **Transformers** | $z_c$ | Class token of last layer |
| | $z_d$ | Domain token of last layer |
| | $N_P$ | Number of patch tokens |
| | $W_Q$ | Query weights |
| | $W_K$ | Key weights |
| | $W_V$ | Value weights |
| | $\theta_Q$ | Query weights of all layers |
| | $\theta_K$ | Key weights of all layers |
| | $\theta_V$ | Value weights of all layers |
| **Datasets** | $\mathcal{D}_s$ | Labeled source dataset |
| | $\mathcal{D}_t$ | Unlabeled target dataset |
| | $\mathcal{A}_i$ | $i^{\text{th}}$ augmentation function |
| | $\mathcal{D}_s^{[i]}$ | $i^{\text{th}}$ augmented source dataset |
| | $\mathcal{D}_t^{[i]}$ | $i^{\text{th}}$ augmented target dataset |
| | $(x_s, y_s)$ | Labeled source sample |
| | $(x_s^{[i]}, y_s, y_d)$ | Augmented source sample |
| | $x_t$ | Unlabeled target sample |
| | $(x_t^{[i]}, y_d)$ | Target augmented sample |
| **Spaces** | $\mathcal{X}$ | Input space |
| | $\mathcal{C}_g$ | Label set for goal task |
| | $\mathcal{Z}_c$ | Class token feature space |
| | $\mathcal{Z}_d$ | Domain token feature space |
| | $\mathcal{Z}_1, \ldots, \mathcal{Z}_{N_P}$ | Patch token |
| **Losses** | $\mathcal{L}_{dom}$ | Domain classification loss |
| | $\mathcal{L}_{cls}$ | Task classification loss |
| | $\mathcal{L}_{im}$ | Entropy loss |
| | $\mathcal{L}_{div}$ | Diversity loss |
| **Criterion** | $\gamma_{dom}$ | Domain specificity |
| | $\gamma_{cls}$ | Task specificity |
| | $\gamma_{all}$ | Inter-class-inter-domain similarity |
| | $\tau$ | Threshold |

$$\mathcal{L}_{div} = \sum_{k=1}^{K} \hat{p_k} \log \hat{p_k} = KL(\hat{p}, \frac{1}{K} 1_K) - \log K \quad (6)$$

where $\delta_k(a) = \frac{\exp(a_k)}{\sum_i \exp(a_i)}$ represents the $k^{\text{th}}$ element in the softmax output of $a \in \mathbb{R}^K$, and $z_c$ is the class-token from $h$ for an input $x_t$. We optimize all parameters of the transformer backbone $h$, except the query weights $\theta_Q$ as follows,

$$\min_{\theta_h \backslash \theta_Q, f_g} \mathop{\mathbb{E}}_{\mathcal{D}_t} [\mathcal{L}_{im} + \mathcal{L}_{div}] \quad (7)$$

We also utilize the clustering method of SHOT [42] for self-supervised pseudo-labeling. First, we obtain the class centroids in the target domain via weighted k-means clustering,

$$c_k = \frac{\sum_{x_t \in \mathcal{X}} \delta_k(f_g(z_c)) z_c}{\sum_{x_t \in \mathcal{X}} \delta_k(f_g(z_c))} \quad (8)$$

The centroid characterizes the labels for the samples. In order to obtain a pseudo-label, we choose the closest centroid based on the cosine distance as follows,

$$\hat{y}_c = \arg\min_k D_c(z_c, c_k) \qquad (9)$$

where $D_c$ denotes the cosine-distance in the class-token feature space $\mathcal{Z}_c$ between the centroid and the input sample features $z_c$. As the model keeps training, the centroids are updated after every few iterations, and pseudo-labels are assigned according to the new centroids.

**(b) Preliminaries on Transformers.** Recently, Vision Transformers (ViT) have been shown to improve significantly on several vision tasks [12]. Self-attention is one of the most important components in the transformer architecture. A ViT takes an image as input $x \in \mathcal{X} = \mathbb{R}^{H \times W \times C}$ in the form of patches of size $(P, P)$, where $H, W$ is the image size and $C$ is the number of channels. The total number of patches is denoted as $N_P = H \times W / P^2$. For self-attention, each patch is projected into $Q, K, V$ with a set of weights $W_Q, W_K, W_V$ respectively. The self-attention [63] is computed as follows,

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V \qquad (10)$$

where $d_k$ is the dimension of the keys/queries.

## B. Implementation Details

In this section, we describe our analysis and benchmarking experiments, which include the augmentation strategies, DRI dataset creation, backbone, and optimization details

### B.1. Domain Augmentations

To induce domain-specificity, we use five label-preserving augmentations to simulate virtual domains (Fig. 5A):
**a) FDA augmentation:** We use FDA [73] to stylize an image with a fixed style-transfer set of images [26]. This is done by superimposing the amplitude spectrum of the style images onto the input image.
**b) Weather augmentations:** We employ frost and snow augmentations [29] to augment the input images.
**c) AdaIN augmentation:** In this augmentation [26], we alter the feature statistics through an instance normalization layer [62] that stylizes the images using the same reference style image set as in FDA.
**d) Cartoon augmentation:** We employ cartoonization-based augmentations [29] to convert inputs to cartoon-like images with reduced texture.
**e) Style augmentation:** We use stylization from Jackson et al. [27]. No controllable parameters are available and style is chosen without a reference style image.

---

**Algorithm 1** DSiT Training Algorithm

---

**Vendor-side training**

1: **Input:** Let $\mathcal{D}_s$ be source data , $\mathcal{D}_s^{[i]}$ be augmented DRI dataset for each augmentation $\mathcal{A}_i$, ImageNet pretrained DeiT-B backbone $h$ from [68], randomly initialized goal classifier $f_g$ and randomly initialized domain classifier $f_d$.

2: **for** $iter < MaxIter$ **do**:

   *Goal task training*

3:   **for** $iter < MaxTaskIters$ **do**:
4:     Sample batch from $\mathcal{D}_s$
5:     Compute $\mathcal{L}_{cls}$ using Eq. 2 (main paper)
6:     **update** $\theta_h \setminus \theta_Q, \theta_{f_g}$ by minimizing $\mathcal{L}_{cls}$
7:   **end for**

   *Domain classifier training*

8:   **for** $iter < MaxDomainIters$ **do**:
9:     Sample batch of DRI from $\mathcal{D}_s^{[i]}$
10:     Compute $\mathcal{L}_{dom}$ using Eq. 1 (main paper)
11:     **update** $\theta_Q, \theta_{f_d}$ by minimizing $\mathcal{L}_{dom}$
12:   **end for**
                    ▷ The two steps are carried out alternatively
13: **end for**

**Client-side training**

14: **Input:** Target data $\mathcal{D}_t$, Target augmented DRI data $\mathcal{D}_t^{[i]}$, source-side pretrained backbone $h$, goal classifier $f_g$ and domain classifier $f_d$.

15: **for** $iter < MaxIter$ **do**:

   *Goal Task Training*

16:   **for** $iter < MaxTaskIters$ **do**:
17:     Sample batch from $\mathcal{D}_t$
18:     Compute $\mathcal{L}_{im}$ and $\mathcal{L}_{div}$ using Eq. 5, 6 (suppl.)
19:     **update** $\theta_h \setminus \theta_Q, \theta_{f_g}$ by minimizing $\mathcal{L}_{im} + \mathcal{L}_{div}$
20:   **end for**

   *Domain classifier training*

21:   **for** $iter < MaxDomainIters$ **do**:
22:     Sample batch of DRI from $\mathcal{D}_t^{[i]}$
23:     Compute $\mathcal{L}_{dom}$ using Eq. 1 (main paper)
24:     **update** $\theta_Q, \theta_{f_d}$ by minimizing $\mathcal{L}_{dom}$
25:   **end for**
                    ▷ The two steps are carried out alternatively.
26: **end for**

---

### B.2. DRI Dataset Extraction

The Domain-Representative Inputs (DRI) are created using augmentations as shown in Fig. 5. An input image is first augmented to simulate a virtual domain. Note that only one augmentation is used at a time. After this, the image
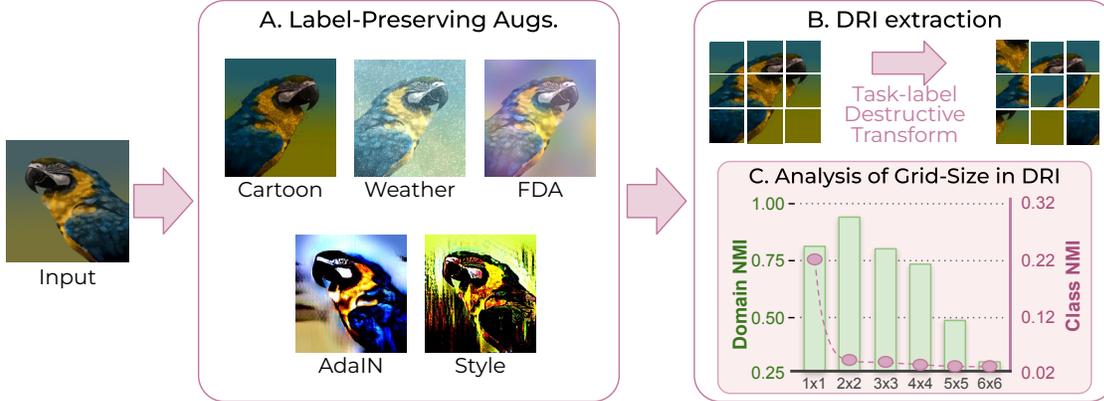
Figure 5. **A.** Label-preserving augmentations are first applied to the input to simulate novel domains. **B.** Then, the task-destructive transform of patch-shuffling is used to obtain the DRI image. **C.** We analyze the Domain-NMI and Class-NMI for different grid-sizes used in patch-shuffling. An example of $3 \times 3$ shuffling is shown in B.

is shuffled across patches to obtain a DRI image. The extent of patch shuffling is done such that the domain information is still intact, however the task-label information is lost. Following prior works [46], we use normalized mutual information (NMI) to assess the consistency between the feature clusters formed by a self-supervised learning algorithm on the transformed images and the class/domain labels (see Fig. 5C). To obtain NMI, the training images are first subjected to the class-destructive transformation to produce DRI images, and these images are then subjected to self-supervised learning to produce class and domain invariant features. In order to assign a domain or class label to each cluster, we finally apply clustering to the learned features. For the self-supervised learning, we employ SimCLR [9] on the DRI images and apply Gaussian mixture-based clustering to the learned features to obtain either domain or class labels for domain-NMI and class-NMI respectively.

From Figure 5, we see that the Domain NMI rises within a certain range as the number of grid partitions increases, whereas the Class-NMI sharply declines. These results demonstrate that domain-specific features can be learned by using an appropriate grid partition size. Hence, for all our experiments, we have used a grid shuffling size of $4 \times 4$ for representing DRI inputs.

## B.3. Domain-specificity disentanglement criterion

As discussed in section 3.3 (main) paper we define the a domain-specificity disentanglement criterion based on three parameters: $\gamma_{cls}$: intra-class, inter-domain similarity, $\gamma_{dom}$: intra-domain inter-class and $\gamma_{all}$ denotes the inter-class, inter-domain similarity. We define the criterion of domain-specificity disentanglement as follows:

$$\gamma_{dom} = \mathop{\mathbb{E}}_{\mathcal{D}_s \cup \mathcal{D}_t} \mathcal{D}_c(z_{c_1}, z_{c_2}), \text{where } y_{c_1} \neq y_{c_2}, y_{d_1} = y_{d_2} \quad (11)$$

$$\gamma_{cls} = \mathop{\mathbb{E}}_{\mathcal{D}_s \cup \mathcal{D}_t} \mathcal{D}_c(z_{c_1}, z_{c_2}), \text{where } y_{c_1} = y_{c_2}, y_{d_1} \neq y_{d_2} \quad (12)$$

$$\gamma_{dom} = \mathop{\mathbb{E}}_{\mathcal{D}_s \cup \mathcal{D}_t} \mathcal{D}_c(z_{c_1}, z_{c_2}), \text{where } y_{c_1} \neq y_{c_2}, y_{d_1} \neq y_{d_2} \quad (13)$$

where $\mathcal{D}_c(z_{c_1}, z_{c_2})$ denotes cosine similarity between the class-token features of two inputs $x_1, x_2$ with corresponding class labels $y_{c_1}, y_{c_2}$ and domain labels $y_{d_1}, y_{d_2}$.

**How to choose the threshold $\tau$?** We empirically evaluated the metrics $\gamma_{cls}$, $\gamma_{dom}$ and $\gamma_{all}$ in Table 6 (main paper) and found that task-specificity $\gamma_{cls}$ and domain-specificity $\gamma_{dom}$ are closer for DSiT (*Ours*) than the SHOT-B baseline. Based on our observations, we choose a threshold of $0.05$.

## B.4. Experimental settings

**Backbone details.** For our experiments, we use DeiT-Base [60] which has 86M parameters, pretrained on ImageNet-1k dataset. DeiT-Base architecture consists of 12 layers, where each layer consists of multi-head self-attention with 12 heads. The input to the transformer is an RGB image that is divided into $16 \times 16$-sized patches. Therefore, $P = 16$ and $N_P = 14$ for all our experiments. DeiT-B contains an additional distillation token, however, the rest of the architecture is the same as a ViT-B backbone.

**Optimization details.** For optimizing the training objectives, we use Stochastic Gradient Descent (SGD) with a momentum of $0.9$, and a weight decay ratio of $1 \times 10^{-4}$. The learning rate is set to $5 \times 10^{-3}$ for fine-tuning the domain classifier on the target domain. For the Goal Task training, we use a learning rate of $8 \times 10^{-3}$ for OfficeHome and VisDA, $8 \times 10^{-2}$ for Office-31, and $2 \times 10^{-3}$ for Domain-Net. The Goal Task Training and Domain-specific disentanglement Training for the vendor-side source domain are done for 20 epochs, of which 10 are used for warm-up with a warm-up factor of $0.01$. In the client-side target adapta-

Table 9. Single-Source Domain Adaptation (SSDA) results on the DomainNet dataset. * indicates results taken from [58].

| ResNet-101 [19] | clp | inf | pnt | qdr | rel | skt | Avg. |
|---|---|---|---|---|---|---|---|
| clp | - | 19.3 | 37.5 | 11.1 | 52.2 | 41.0 | 32.2 |
| inf | 30.2 | - | 31.2 | 3.6 | 44.0 | 27.9 | 27.4 |
| pnt | 39.6 | 18.7 | - | 4.9 | 54.5 | 36.3 | 30.8 |
| qdr | 7.0 | 0.9 | 1.4 | - | 4.1 | 8.3 | 4.3 |
| rel | 48.4 | 22.2 | 49.4 | 6.4 | - | 38.8 | 33.0 |
| skt | 46.9 | 15.4 | 37.0 | 10.9 | 47.0 | - | 31.4 |
| Avg. | 34.4 | 15.3 | 31.3 | 7.4 | 40.4 | 30.5 | **26.6** |

| CDAN [45] | clp | inf | pnt | qdr | rel | skt | Avg. |
|---|---|---|---|---|---|---|---|
| clp | - | 20.4 | 36.6 | 9.0 | 50.7 | 42.3 | 31.8 |
| inf | 27.5 | - | 25.7 | 1.8 | 34.7 | 20.1 | 22.0 |
| pnt | 42.6 | 20.0 | - | 2.5 | 55.6 | 38.5 | 31.8 |
| qdr | 21.0 | 4.5 | 8.1 | - | 14.3 | 15.7 | 12.7 |
| rel | 51.9 | 23.3 | 50.4 | 5.4 | - | 41.4 | 34.5 |
| skt | 50.8 | 20.3 | 43.0 | 2.9 | 50.8 | - | 33.6 |
| Avg. | 38.8 | 17.7 | 32.8 | 4.3 | 41.2 | 31.6 | **27.7** |

| MIMFTL [17] | clp | inf | pnt | qdr | rel | skt | Avg. |
|---|---|---|---|---|---|---|---|
| clp | - | 15.1 | 35.6 | 10.7 | 51.5 | 43.1 | 31.2 |
| inf | 32.1 | - | 31.0 | 2.9 | 48.5 | 31.0 | 29.1 |
| pnt | 40.1 | 14.7 | - | 4.2 | 55.4 | 36.8 | 30.2 |
| qdr | 18.8 | 3.1 | 5.0 | - | 16.0 | 13.8 | 11.3 |
| rel | 48.5 | 19.0 | 47.6 | 5.8 | - | 39.4 | 32.1 |
| skt | 51.7 | 16.5 | 40.3 | 12.3 | 53.5 | - | 34.9 |
| Avg. | 38.2 | 13.7 | 31.9 | 7.2 | 45.0 | 32.8 | **28.1** |

| MDD+SCDA [75] | clp | inf | pnt | qdr | rel | skt | Avg. |
|---|---|---|---|---|---|---|---|
| clp | - | 20.4 | 43.3 | 15.2 | 59.3 | 46.5 | 36.9 |
| inf | 32.7 | - | 34.5 | 6.3 | 47.6 | 29.2 | 30.1 |
| pnt | 46.4 | 19.9 | - | 8.1 | 58.8 | 42.9 | 35.2 |
| qdr | 31.1 | 6.6 | 18.0 | - | 28.8 | 22.0 | 21.3 |
| rel | 55.5 | 23.7 | 52.9 | 9.5 | - | 45.2 | 37.4 |
| skt | 55.8 | 20.1 | 46.5 | 15.0 | 56.7 | - | 38.8 |
| Avg. | 44.3 | 18.1 | 39.0 | 10.8 | 50.2 | 37.2 | **33.3** |

| DeiT-B [60] | clp | inf | pnt | qdr | rel | skt | Avg. |
|---|---|---|---|---|---|---|---|
| clp | - | 24.3 | 49.6 | 15.8 | 65.3 | 52.1 | 41.4 |
| inf | 45.9 | - | 45.9 | 6.7 | 61.4 | 39.5 | 39.9 |
| pnt | 53.2 | 23.8 | - | 6.5 | 66.4 | 44.7 | 38.9 |
| qdr | 31.9 | 6.8 | 15.4 | - | 23.4 | 20.6 | 19.6 |
| rel | 59.0 | 25.8 | 56.3 | 8.1 | - | 44.8 | 39.0 |
| skt | 60.6 | 20.6 | 48.4 | 16.5 | 61.2 | - | 41.5 |
| Avg. | 50.1 | 20.3 | 43.1 | 10.9 | 55.5 | 40.3 | **36.7** |

| SHOT-B [42] | clp | inf | pnt | qdr | rel | skt | Avg. |
|---|---|---|---|---|---|---|---|
| clp | - | 27.0 | 49.7 | 16.5 | 65.4 | 53.2 | 46.1 |
| inf | 46.4 | - | 45.9 | 7.4 | 60.6 | 40.1 | 40.1 |
| pnt | 54.6 | 25.7 | - | 8.1 | 66.3 | 49.0 | 40.7 |
| qdr | 33.3 | 6.8 | 15.5 | - | 23.8 | 24.0 | 20.7 |
| rel | 59.3 | 28.1 | 57.4 | 9.0 | - | 47.3 | 40.2 |
| skt | 64.0 | 26.5 | 55.0 | 18.2 | 63.8 | - | 45.5 |
| Avg. | 51.5 | 26.6 | 44.7 | 11.8 | 56.0 | 42.7 | **38.9** |

| CDTrans* [68] | clp | inf | pnt | qdr | rel | skt | Avg. |
|---|---|---|---|---|---|---|---|
| clp | - | 27.9 | 57.6 | 27.9 | 73.0 | 58.8 | 49.0 |
| inf | 58.6 | - | 53.4 | 9.6 | 71.1 | 47.6 | 48.1 |
| pnt | 60.7 | 24.0 | - | 13.0 | 69.8 | 49.6 | 43.4 |
| qdr | 2.9 | 0.4 | 0.3 | - | 0.7 | 4.7 | 1.8 |
| rel | 49.3 | 18.7 | 47.8 | 9.4 | - | 33.5 | 31.7 |
| skt | 66.8 | 23.7 | 54.6 | 27.5 | 68.0 | - | 48.1 |
| Avg. | 47.7 | 18.9 | 42.7 | 17.5 | 56.5 | 38.8 | **37.0** |

| SSRT-B* [58] | clp | inf | pnt | qdr | rel | skt | Avg. |
|---|---|---|---|---|---|---|---|
| clp | - | 33.8 | 60.2 | 19.4 | 75.8 | 59.8 | 49.8 |
| inf | 55.5 | - | 54.0 | 9.0 | 68.2 | 44.7 | 46.3 |
| pnt | 61.7 | 28.5 | - | 8.4 | 71.4 | 55.2 | 45.0 |
| qdr | 42.5 | 8.8 | 24.2 | - | 37.6 | 33.6 | 29.3 |
| rel | 69.9 | 37.1 | 66.0 | 10.1 | - | 58.9 | 48.4 |
| skt | 70.6 | 32.8 | 62.2 | 21.7 | 73.2 | - | 52.1 |
| Avg. | 60.0 | 28.2 | 53.3 | 13.7 | 65.3 | 50.4 | **45.2** |

| DSiT (Ours) | clp | inf | pnt | qdr | rel | skt | Avg. |
|---|---|---|---|---|---|---|---|
| clp | - | 27.2 | 51.8 | 23.1 | 70.2 | 54.7 | 45.4 |
| inf | 52.3 | - | 48.8 | 12.8 | 68.3 | 44.2 | 45.3 |
| pnt | 59.2 | 26.1 | - | 14.5 | 71.5 | 51.4 | 44.5 |
| qdr | 38.1 | 8.3 | 21.2 | - | 37.2 | 27.6 | 26.5 |
| rel | 60.4 | 28.0 | 57.8 | 13.1 | - | 49.7 | 41.8 |
| skt | 66.3 | 27.5 | 56.0 | 24.4 | 70.2 | - | 48.9 |
| Avg. | 55.3 | 23.4 | 47.1 | 17.6 | 63.5 | 45.5 | **42.1** |

Table 10. Vendor-side Performance of Single-Source Domain Adaptation (SSDA) on Office-Home, DomainNet and VisDA.

| Training stage | Method | Office-Home (4 settings) | VisDA | DomainNet |
|---|---|---|---|---|
| Vendor-side | SHOT | 74.8 | 67.4 | 36.7 |
| | DSiT | 76.6 (+1.8) | 68.7 (+1.3) | 37.8 (+1.1) |
| Client-side | SHOT | 79.0 | 85.9 | 38.3 |
| | DSiT | 81.1 (+2.1) | 87.5 (+1.6) | 42.1 (+3.8) |

Table 11. Single-Source Domain Adaptation (SSDA) on Office-Home (4 settings) for different vendor-side and client-side adaptation strategies.

| # | Vendor-side | Client-side | Ar→Cl | Cl→Pr | Pr→Rw | Rw→Ar | Avg |
|---|---|---|---|---|---|---|---|
| 1. | SHOT | SHOT | 67.1 | 83.4 | 85.3 | 80.4 | 79.1 |
| 2. | SHOT | DSiT | 68.4 | 85.7 | 86.8 | 81.8 | 80.7 (+1.6) |
| 3. | DSiT | SHOT | 67.1 | 83.0 | 84.9 | 81.0 | 79.0 |
| 4. | DSiT | DSiT | 69.2 | 86.8 | 86.6 | 82.4 | 81.3 (+2.3) |

tion, the goal task training (*task classifier training*) is carried out for 2 epochs, followed by the domain specificity disentanglement (*domain classifier training*) until a domain classification accuracy of 80% is achieved. These two steps are carried out alternatively for an effective 40 epochs of task-classifier training, the same as CDTrans [68]. We use an NVIDIA RTX A5000 GPU with 64GB RAM and 24GB GPU memory to train our models. Our code takes a total training time of approximately 5 hours for Target adaptation training for the Office-Home dataset.

## C. Additional Experimental Results

### C.1. Extended Comparisons

**Comparisons on single-source domain adaptation (SSDA):** In Tables 9, we show additional comparisons for our method with existing SSDA works on the DomainNet benchmark. We achieve significant improvements over existing works, especially on CDTrans [68] despite it being a non-source-free method. It is worth noting that CDTrans uses the entire domain during the training and evaluation steps, while we train on the *train* split and evaluate on the *test* split, same as SSRT [58].

**Multi-target domain adaptation (MTDA):** In Table 13, we provide a quantitative comparison with the prior arts for multi-target domain adaptation on Office-Home. The performance improvement is quite prominent (+2.0%) over the source-free prior art (SHOT-B), and the proposed approach also yields comparable performance to non-source-free prior arts, D-CGCT and CDAN+DCL [54], which mainly focus on domain invariant features.

### C.2. Vendor-side DSiT Performance

Our DSiT approach incorporates a novel Domain-Specificity Training (DST) that improves the vendor-side performance over the standard source-only baseline (shown in Table 10). Further, we observe significant gains from vendor to client-side (4.5% and 4.3% for Office Home and DomainNet, Table 10). This shows that our vendor-side DST positively aids client-side DST.

Table 12. Single-Source Domain Adaptation (SSDA) on Office-Home on ViT-S and DeiT-S Backbone. SF indicates *source-free* adaptation. "-S" denotes Small ViT Backbone.

| Method | SF | Office-Home | | | | |
|---|---|---|---|---|---|---|
| | | Ar→Cl | Cl→Pr | Pr→Rw | Rw→Ar | Avg. |
| CDTrans-S [68] | ✗ | 60.7 | 75.6 | 84.4 | 77.0 | 74.4 |
| SHOT-S | ✓ | **56.3** | 73.7 | 81.3 | 76.7 | 71.9 |
| **DSiT-S (*Ours*)** | ✓ | 55.3 | **77.4** | **83.0** | **76.9** | **73.1 (+1.2)** |
| SHOT-B | ✓ | 69.07 | 85.31 | 88.13 | 83.89 | 81.6 |
| **DSiT-B (*Ours*)** | ✓ | **71.84** | **87.18** | 88.11 | 83.4 | **82.6 (+1.0)** |

Table 13. Multi-Target Domain Adaptation (MTDA) on Office-Home. SF indicates *source-free* adaptation. ResNet-based methods (top) and Transformer-based methods (bottom).

| Method | SF | Office-Home | | | | |
|---|---|---|---|---|---|---|
| | | Ar→ | Cl→ | Pr→ | Rw→ | Avg. |
| MT-MTDA [48] | ✗ | 64.6 | 66.4 | 59.2 | 67.1 | 64.3 |
| CDAN+DCL [45] | ✗ | 63.0 | 66.3 | 60.0 | 67.0 | 64.1 |
| D-CGCT [54] | ✗ | 70.5 | 71.6 | 66.0 | 71.2 | 69.8 |
| D-CGCT-B [54] | ✗ | 77.0 | 78.5 | 77.9 | 80.9 | 78.6 |
| SHOT-B* | ✓ | 75.4 | 79.3 | 73.6 | 77.1 | 76.4 |
| **DSiT-B (*Ours*)** | ✓ | **77.3** | **83.4** | 75.6 | 76.8 | **78.3 (+1.9)** |

## C.3. Performance in a Model Adaptation Setting

Our DSiT approach works well even for a model adaptation setting, where we perform DST only on client-side without any specialized training on the vendor-side (#2, Table 11). However, we get the best results when DST is done on both vendor and client-side (#4, Table 11). We also observe that SHOT target adaptation (TA) with our vendor-side DSiT model (#2) gives the same performance as the baseline (#1). This indicates that our proposed TA (#4) is able to better leverage our vendor-side model to yield improved adaptation performance.

Table 14. Training time comparison of our approach DSiT vs SHOT on Office-Home (Rw→Ar)

| Method | Training time (in min) | | | | | Inf. time (ms) | Acc. |
|---|---|---|---|---|---|---|---|
| | Src. train | Src. DST | Tgt. adapt | Tgt. DST | Total time | | |
| SHOT-B | 12 | - | 17 | - | 29 | 3.6 | 80.4 |
| *Ours* | 12 | 109 | - | 258 | 270 | 3.6 | **82.4** |

## C.4. Performance on different backbones

We report results in Table 12 for DeiT-S backbone (with 22M parameters) pre-trained on ImageNet and observe that our approach improves over the baseline SHOT-S baseline by 1.2%. Note that "-S" denotes Small. We also report the results over ViT-B backbone which is trained on ImageNet-21K dataset. Over ViT-B, our approach shows an improvement of 1.0% over the SHOT baseline.

Table 15. Analysis for $\mathcal{A}$-distance of three augmentations on 4 settings of Office-Home (SSDA).

| Aug. | Ar→Cl | Cl→Pr | Pr→Rw | Rw→Ar |
|---|---|---|---|---|
| FDA | 0.857 | 0.730 | 0.829 | 0.286 |
| **Original** | **1.049** | **0.852** | **0.834** | **0.504** |
| AdaIN | 1.072 | 0.648 | 0.842 | 0.136 |

Table 16. Sensitivity Analysis on Single-Source Domain Adaptation (SSDA) on Office-Home. (4 settings)

| Epochs | Ar → Cl | Cl → Pr | Pr → Rw | Rw → Ar | Avg. |
|---|---|---|---|---|---|
| 1 | 64.1 | 79.8 | 84.7 | 79.6 | 77.0 |
| 2 | 69.2 | 86.1 | 86.6 | 82.4 | 81.1 |
| 3 | 69.6 | 86.7 | 87.3 | 82.4 | 81.5 |
| 5 | 69.5 | 86.3 | 87.1 | 82.5 | 81.3 |

Table 17. Significance experiments of DSiT-B (Ours) on Single-Source DA (SSDA) on Office-Home (4 settings).

| Ar → Cl | Cl → Pr | Pr → Rw | Rw → Ar | Avg. |
|---|---|---|---|---|
| 69.2 ±0.1 | 86.1 ±0.3 | 86.6 ±0.3 | 82.4 ±0.6 | 81.1 ±0.1 |

## C.5. Experimental analysis for augmentations

In Table 15, we show the $\mathcal{A}$-distance (domain-gap) between augmented source and target domains on Office-Home using the class token of a source-trained DeiT-B. The domain gap for FDA is lower than the original source-target while it is higher for AdaIN. This validates Fig. 4 (main paper) which illustrated that augmented domains may be closer or farther than original domains.

Table 18. Ablation study for the three components of the client-side adaptation on 4 settings of Office-Home. *PL* indicates pseudo-labeling

| Method | Ar → Cl | Cl → Pr | Pr → Rw | Rw → Ar | Avg. |
|---|---|---|---|---|---|
| Source Baseline | 62.5 | 79.4 | 84.3 | 79.2 | 76.4 |
| $\mathcal{L}_{im}$ | 62.1 | 79.7 | 80.1 | 73.9 | 74.0 |
| $\mathcal{L}_{im} + \mathcal{L}_{div}$ | 68.2 | 86.0 | 86.6 | 81.3 | 80.5 |
| $\mathcal{L}_{im} + \mathcal{L}_{div} + PL$ | 69.2 | 86.1 | 86.6 | 82.4 | 81.1 |

## C.6. Sensitivity Analysis of Alternate Training

In our approach, we perform alternate rounds of training of the domain and the task classifier. We usually train the task classifier for a few epochs, followed by the domain classifier training. In this analysis, we vary the number of epochs of task classifier training from 1 to 5 epochs in an alternate round and observe its impact on task accuracy. Table 16 shows that the task accuracy increases at 2 epochs and is maximum at 3 epochs of training. In all our experiments, we report results with 2 epochs of task classifier training.

## C.7. Training time comparisons

We provide detailed training and inference time comparisons of our method with SHOT-B in Table 14. The DST
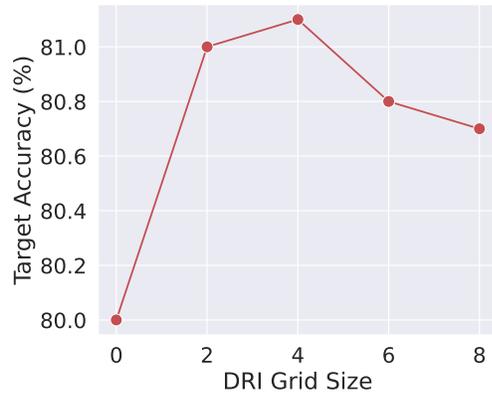
Figure 6. Sensitivity analysis on DRI grid-size for Single-Source DA on four settings of Office-Home.

training time is higher due to augmented images being computed at training time, which can be easily reduced by loading pre-computed augmented images. We point out that the inference time remains the same, highlighting the fact that the same DeiT-Base architecture is used for both methods.

### C.8. Statistical Significance

We report mean and standard deviation over 3 runs for three Office-Home settings in Table 17. We observe that the standard deviation (0.1 to 0.3) is very low w.r.t. to our gains ($\sim$2%) over SHOT-B.

### C.9. Effect of target adaptation losses

We perform an ablation study on 4 settings of the Office-Home dataset to analyze the influence of each component of the target adaptation objective described in Section A, and present the results in Table 18. Target adaptation with the entropy loss $\mathcal{L}_{im}$ alone shows sub-optimal results, even when compared to the source-trained baselines, which is also observed by [42]. Adding the diversity loss $\mathcal{L}_{div}$ shows comparatively better performance, indicating that balancing the classifier's predictions across all classes is essential. Lastly, the self-supervised pseudo-labeling *PL* also improves the performance further, demonstrating its importance towards the client-side adaptation.

### C.10. Effect of DRI grid-size

Here, we study the effect of varying the DRI grid size for the domain classifier training to determine its influence on the target accuracy (Figure 6). The target accuracy gradually increases upon increasing the grid size from 1 to 4, with the best results being observed at grid size 4. Beyond this point, the performance begins to drop, which can be attributed to the excessive destruction of information caused by over-partitioning of the images. To achieve a balance between the effect of the task-destructive transformation and its impact on the target accuracy, we use a grid size of $4 \times 4$ for all our experiments.

## References

[1] Surbhi Aggarwal, Jogendra Nath Kundu, R. Venkatesh Babu, and Anirban Chakraborty. WAMDA: Weighted alignment of sources for multi-source domain adaptation. In *BMVC*, 2020. 8

[2] Sk. Miraj Ahmed, Dripta S. Raychaudhuri, S. Paul, Samet Oymak, and Amit K. Roy-Chowdhury. Unsupervised multi-source domain adaptation without access to source data. In *CVPR*, 2021. 8

[3] Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wortman Vaughan. A theory of learning from different domains. *Machine learning*, 79(1-2):151–175, 2010. 1

[4] Jinming Cao, Oren Katzir, Peng Jiang, Dani Lischinski, Daniel Cohen-Or, Changhe Tu, and Yangyan Li. DiDA: Iterative boosting of disentangled synthesis and domain adaptation. In *IEEE ITME*, 2021. 2

[5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 2

[6] Woong-Gi Chang, Tackgeun You, Seonguk Seo, Suha Kwak, and Bohyung Han. Domain-specific batch normalization for unsupervised domain adaptation. In *CVPR*, 2019. 2

[7] Prithvijit Chattopadhyay, Yogesh Balaji, and Judy Hoffman. Learning to balance specificity and invariance for in and out of domain generalization. In *ECCV*, 2020. 2

[8] Chao Chen, Zhihang Fu, Zhihong Chen, Sheng Jin, Zhaowei Cheng, Xinyu Jin, and Xian-Sheng Hua. HoMM: Higher-order moment matching for unsupervised domain adaptation. In *AAAI*, 2020. 2, 4

[9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020. 11

[10] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. 1

[11] Jiahua Dong, Zhen Fang, Anjin Liu, Gan Sun, and Tongliang Liu. Confident anchor-induced multi-source free domain adaptation. In *NeurIPS*, 2021. 8

[12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 3, 10

[13] Qi Dou, Cheng Ouyang, Cheng Chen, Hao Chen, and Pheng-Ann Heng. Unsupervised cross-modality domain adaptation of convnets for biomedical image segmentations with adversarial loss. In *IJCAI*, 2018. 1

[14] Zhekai Du, Jingjing Li, Hongzu Su, Lei Zhu, and Ke Lu. Cross-domain gradient discrepancy minimization for unsupervised domain adaptation. In *CVPR*, 2021. 7

[15] Abhimanyu Dubey, Vignesh Ramanathan, Alex Pentland, and Dhruv Mahajan. Adaptive methods for real-world domain generalization. In *CVPR*, 2021. 1, 2, 3

[16] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario

Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The Journal of Machine Learning Research*, 17(1):2096–2030, 2016. 1, 3

[17] Jian Gao, Yang Hua, Guosheng Hu, Chi Wang, and Neil M Robertson. Reducing distributional uncertainty by mutual information maximisation and transferable feature learning. In *ECCV*, 2020. 12

[18] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. In *ICLR*, 2017. 3

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 6, 7, 12

[20] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. CyCADA: Cycle-consistent adversarial domain adaptation. In *ICML*, 2018. 2

[21] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. DAFormer: Improving network architectures and training strategies for domain-adaptive semantic segmentation. In *CVPR*, 2022. 2

[22] Lanqing Hu, Meina Kan, Shiguang Shan, and Xilin Chen. Duplex generative adversarial network for unsupervised domain adaptation. In *CVPR*, 2018. 2

[23] Ronghang Hu and Amanpreet Singh. UniT: Multimodal multitask learning with a unified transformer. In *ICCV*, 2021. 2, 3

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

[25] Sheng-Wei Huang, Che-Tsung Lin, Shu-Ping Chen, Yen-Yi Wu, Po-Hao Hsu, and Shang-Hong Lai. Auggan: Cross domain adaptation with gan-based data augmentation. In *ECCV*, 2018. 3

[26] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, 2017. 10

[27] Philip TG Jackson, Amir Atapour Abarghouei, Stephen Bonner, Toby P Breckon, and Boguslaw Obara. Style augmentation: data augmentation via style randomization. In *CVPR workshops*, 2019. 10

[28] Xin Jin, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. Re-energizing domain discriminator with sample relabeling for adversarial domain adaptation. In *ICCV*, 2021. 7

[29] Jung. imgaug. In *https://github.com/aleju/imgaug*, 2020. 10

[30] Guoliang Kang, Lu Jiang, Yi Yang, and Alexander G Hauptmann. Contrastive adaptation network for unsupervised domain adaptation. In *CVPR*, 2019. 7

[31] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM Computing Surveys (CSUR)*, 2021. 2

[32] Jogendra Nath Kundu, Suvaansh Bhambri, Akshay Kulkarni, Hiran Sarkar, Varun Jampani, and R Venkatesh Babu. Concurrent subsidiary supervision for unsupervised source-free domain adaptation. In *ECCV*, 2022. 2

[33] Jogendra Nath Kundu, Suvaansh Bhambri, Akshay R Kulkarni, Hiran Sarkar, Varun Jampani, and R Venkatesh Babu. Subsidiary prototype alignment for universal domain adaptation. In *NeurIPS*, 2022. 2

[34] Jogendra Nath Kundu, Akshay Kulkarni, Suvaansh Bham-

bri, Deepesh Mehta, Shreyas Kulkarni, Varun Jampani, and R. Venkatesh Babu. Balancing discriminability and transferability for source-free domain adaptation. In *ICML*, 2022. 2, 6, 7

[35] Jogendra Nath Kundu, Akshay Kulkarni, Amit Singh, Varun Jampani, and R. Venkatesh Babu. Generalize then adapt: Source-free domain adaptive semantic segmentation. In *ICCV*, 2021. 2, 4

[36] Jogendra Nath Kundu, Naveen Venkat, M V Rahul, and R. Venkatesh Babu. Universal source-free domain adaptation. In *CVPR*, 2020. 1, 2, 3

[37] Seunghun Lee, Sunghyun Cho, and Sunghoon Im. DRANet: Disentangling representation and adaptation networks for unsupervised cross-domain adaptation. In *CVPR*, 2021. 2

[38] Suhyeon Lee, Junhyuk Hyun, Hongje Seong, and Euntai Kim. Unsupervised domain adaptation for semantic segmentation by content transfer. In *AAAI*, 2021. 2

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

[40] Shuang Li, Mixue Xie, Fangrui Lv, Chi Harold Liu, Jian Liang, Chen Qin, and Wei Li. Semantic concentration for domain adaptation. In *ICCV*, 2021. 6

[41] Xiangyu Li, Yonghong Hou, Pichao Wang, Zhimin Gao, Mingliang Xu, and Wanqing Li. Trear: Transformer-based rgb-d egocentric action recognition. *IEEE Transactions on Cognitive and Developmental Systems*, 14(1):246–252, 2021. 3

[42] 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, 3, 5, 6, 7, 8, 9, 12, 14

[43] Jian Liang, Dapeng Hu, Yunbo Wang, Ran He, and Jiashi Feng. Source data-absent unsupervised domain adaptation through hypothesis transfer and labeling transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. 2, 6, 7, 8

[44] Yang Liu, Yao Zhang, Yixin Wang, Feng Hou, Jin Yuan, Jiang Tian, Yang Zhang, Zhongchao Shi, Jianping Fan, and Zhiqiang He. A survey of visual transformers. *arXiv preprint arXiv:2111.06091*, 2021. 2

[45] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. Conditional adversarial domain adaptation. In *NeurIPS*, 2017. 12, 13

[46] Yu Mitsuzumi, Go Irie, Daiki Ikami, and Takashi Shibata. Generalized domain adaptation. In *CVPR*, 2021. 4, 11

[47] Jaemin Na, Heechul Jung, Hyung Jin Chang, and Wonjun Hwang. FixBi: Bridging domain spaces for unsupervised domain adaptation. In *CVPR*, 2021. 7

[48] Le Thanh Nguyen-Meidine, Atif Belal, Madhu Kiran, Jose Dolz, Louis-Antoine Blais-Morin, and Eric Granger. Unsupervised multi-target domain adaptation through knowledge distillation. In *WACV*, 2021. 13

[49] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *ICCV*, 2019. 2, 6

[50] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. VisDA: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924*,

2017. 6

[51] Viraj Prabhu, Shivam Khare, Deeksha Kartik, and Judy Hoffman. SENTRY: Selective entropy optimization via committee consistency for unsupervised domain adaptation. In *ICCV*, 2021. 6

[52] Yao Qin, Chiyuan Zhang, Ting Chen, Balaji Lakshminarayanan, Alex Beutel, and Xuezhi Wang. Understanding and improving robustness of vision transformers through patch-based negative augmentation. In *NeurIPS Workshop*, 2021. 4

[53] Zhen Qiu, Yifan Zhang, Hongbin Lin, Shuaicheng Niu, Yanxia Liu, Qing Du, and Mingkui Tan. Source-free domain adaptation via avatar prototype generation and adaptation. In *IJCAI*, 2021. 6, 7

[54] Subhankar Roy, Evgeny Krivosheev, Zhun Zhong, Nicu Sebe, and Elisa Ricci. Curriculum graph co-teaching for multi-target domain adaptation. In *CVPR*, 2021. 12, 13

[55] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *ECCV*, 2010. 2, 6

[56] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*, 2018. 2

[57] Marin Scalbert, Maria Vakalopoulou, and Florent Couzini'e-Devy. Multi-source domain adaptation via supervised contrastive learning and confident consistency regularization. In *BMVC*, 2021. 8

[58] Tao Sun, Cheng Lu, Tianshuo Zhang, and Haibin Ling. Safe self-refinement for transformer-based domain adaptation. In *CVPR*, 2022. 2, 6, 7, 12

[59] 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*, 2021. 7

[60] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021. 6, 7, 11, 12

[61] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *ACL*, 2019. 3

[62] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *CVPR*, 2017. 10

[63] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 3, 10

[64] Naveen Venkat, Jogendra Nath Kundu, Durgesh Kumar Singh, Ambareesh Revanur, and R. Venkatesh Babu. Your classifier can secretly suffice multi-source domain adaptation. In *NeurIPS*, 2020. 8

[65] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *CVPR*, 2017. 6

[66] Fan Wang, Zhongyi Han, Yongshun Gong, and Yilong Yin. Exploring domain-invariant parameters for source free do-
main adaptation. In *CVPR*, 2022. 6, 7

[67] Haifeng Xia, Handong Zhao, and Zhengming Ding. Adaptive adversarial network for source-free domain adaptation. In *ICCV*, 2021. 6, 7

[68] Tongkun Xu, Weihua Chen, Pichao Wang, Fan Wang, Hao Li, and Rong Jin. CDTrans: Cross-domain transformer for unsupervised domain adaptation. In *ICLR*, 2022. 2, 3, 6, 7, 8, 10, 12, 13

[69] Guanglei Yang, Hao Tang, Zhun Zhong, Mingli Ding, Ling Shao, Nicu Sebe, and Elisa Ricci. Transformer-based source-free domain adaptation. In *APIN*, 2022. 2

[70] Jinyu Yang, Jingjing Liu, Ning Xu, and Junzhou Huang. TVT: Transferable vision transformer for unsupervised domain adaptation. In *WACV*, 2023. 2, 6, 7

[71] Shiqi Yang, Yaxing Wang, Joost van de Weijer, Luis Herranz, and Shangling Jui. Exploiting the intrinsic neighborhood structure for source-free domain adaptation. In *NeurIPS*, 2021. 2, 3, 5, 6, 7

[72] Shiqi Yang, Yaxing Wang, Joost van de Weijer, Luis Herranz, and Shangling Jui. Generalized source-free domain adaptation. In *ICCV*, 2021. 5, 6

[73] Yanchao Yang and Stefano Soatto. FDA: Fourier domain adaptation for semantic segmentation. In *CVPR*, 2020. 10

[74] Kun Zhang, Bernhard Schölkopf, Krikamol Muandet, and Zhikun Wang. Domain adaptation under target and conditional shift. In *ICML*, 2013. 2

[75] Yuchen Zhang, Tianle Liu, Mingsheng Long, and Michael Jordan. Bridging theory and algorithm for domain adaptation. In *ICML*, 2019. 7, 12