

Concurrent Subsidiary Supervision for Unsupervised Source-Free Domain Adaptation

Jogendra Nath Kundu^{1*}, Suvaansh Bhambri^{1*}, Akshay Kulkarni^{1*},
Hiran Sarkar¹, Varun Jampani², and R. Venkatesh Babu¹

¹ Indian Institute of Science

² Google Research

Abstract. The prime challenge in unsupervised domain adaptation (DA) is to mitigate the domain shift between the source and target domains. Prior DA works show that pretext tasks could be used to mitigate this domain shift by learning domain invariant representations. However, in practice, we find that most existing pretext tasks are ineffective against other established techniques. Thus, we theoretically analyze how and when a subsidiary pretext task could be leveraged to assist the goal task of a given DA problem and develop objective subsidiary task suitability criteria. Based on this criteria, we devise a novel process of sticker intervention and cast sticker classification as a supervised subsidiary DA problem concurrent to the goal task unsupervised DA. Our approach not only improves goal task adaptation performance, but also facilitates privacy-oriented source-free DA i.e. without concurrent source-target access. Experiments on the standard Office-31, Office-Home, DomainNet, and VisDA benchmarks demonstrate our superiority for both single-source and multi-source source-free DA. Our approach also complements existing non-source-free works, achieving leading performance.

1 Introduction

The prevalent trend in supervised deep learning systems is to assume that training and testing data follow the same distribution. However, such models often fail [7] when deployed in a new environment (target domain) due to the discrepancy in the training (source domain) and target distributions. A standard approach to deal with this problem of *domain shift* is Unsupervised Domain Adaptation (DA) [13,41], which aims to minimize the domain discrepancy [4] between source and target. The prime challenge in DA is to facilitate the effective utilization of the unlabeled samples while adapting to the target domain.

Drawing motivation from self-supervised pretext task literature [47,16], recent DA works [6,43] have adopted subsidiary tasks as side-objectives to improve the adaptation performance. The intuition is that subsidiary task objectives enforce learning of domain-generic representations, leading to improved domain alignment [66] and consequently, better feature clustering for unlabeled target

* Equal contribution | Webpage: <https://sites.google.com/view/sticker-sfda>

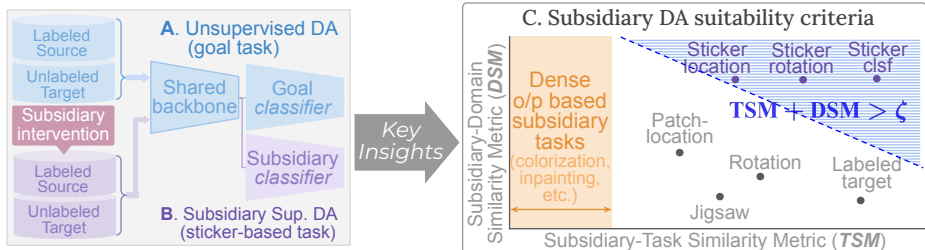


Fig. 1. We tackle **A.** unsupervised goal task DA by introducing **B.** a concurrent subsidiary supervised DA. **C.** Our theoretical insights reveal that subsidiary tasks having both higher TSM (X-axis) and DSM (Y-axis) are most suitable for concurrent goal-subsidary adaptation (*i.e.* the shaded blue area). The proposed sticker-based tasks better suit concurrent goal-subsidary DA among other self-supervised pretext tasks.

[43]. We aim to design a similar framework but, contrary to prior works, we adopt a novel perspective of subsidiary supervised DA for the subsidiary task concurrent to unsupervised goal task DA. Specifically, the framework involves a shared backbone with a goal classifier and a subsidiary classifier (Fig. 1A, B).

To better understand how subsidiary supervised DA objectives support goal task DA, we intend to theoretically analyze the proposed framework. While several subsidiary tasks are available in the literature, there has been little attention on identifying the desirable properties of a subsidiary task that would better aid the unsupervised DA. A recent self-supervised work [71] studied the effectiveness of pretraining with existing subsidiary tasks [47,16] on different downstream supervised settings such as fine-grained or medical image classification [50,72]. We argue that our intended theoretical analysis is necessary to understand the same for DA settings as DA presents a different set of challenges compared to downstream supervised learning paradigms.

Thus, we attempt to answer two interconnected questions,

1. *How does subsidiary supervised DA help goal task unsupervised DA?*
2. *What kind of subsidiary tasks better suit concurrent goal-subsidary DA?*

For the first question, we uncover theoretical insights based on generalization bounds in DA [4,84]. These bounds define distribution shift or domain discrepancy between source and target as the worst discrepancy for a given hypothesis space. We analyze the effect of adding the subsidiary supervised DA problem on the hypothesis space of the shared backbone. Based on this, we find that a higher domain similarity between goal and subsidiary task samples leads to a lower domain discrepancy. This leads to better adaptation for concurrent goal-subsidary DA w.r.t. naive goal DA. Further, we observe that a higher goal-subsidary task similarity aids effective learning of both tasks with the shared backbone, which is crucial for subsidiary DA to positively impact the goal DA.

For the second question, we first devise a subsidiary-domain similarity metric (DSM) and a subsidiary-task similarity metric (TSM) to measure the domain similarity and task similarity between any subsidiary task with a given goal task. Based on our theoretical insights, we propose a subsidiary task suitability criteria using both DSM and TSM to identify *DA-assistive* subsidiary tasks. With this

criteria, we evaluate the commonly used subsidiary tasks from the pretext task literature like rotation prediction [43], patch location [66], and jigsaw permutation prediction [6] in Fig. 1C. We observe that these existing tasks have significantly low DSM. On the other hand, dense output based tasks like colorization [33] or inpainting [52] severely lack in TSM as goal task is classification-based. Understanding these limitations, we devise a sticker-intervention that facilitates domain preservation (high DSM) and propose a range of sticker-based subsidiary tasks (Fig. 2). For general shape-based goal tasks, it turns out that sticker classification task has the best TSM among other sticker-based tasks. This yields higher adaptation performance thereby validating the proposed criteria.

To evaluate our theoretical insights and the proposed concurrent subsidiary DA, we particularly focus on source-free DA regime [34,29]. In this, the source and target data are not concurrently accessible while model



Fig. 2. Sticker intervention involves mixup of input with a masked sticker. We devise the following sticker-based tasks; **A.** locating the quadrant of the sticker, **B.** predicting sticker rotation, **C.** classifying sticker category.

sharing is permitted. While this challenging setting holds immense practical value by working within the data privacy regulations, we choose source-free DA as it can prominently highlight our advantages. The well-developed discrepancy minimization techniques, tailored to general DA scenarios, guide the adaptation more significantly than our proposed approach but cannot be used for source-free DA. Further, existing source-free works [37] rely heavily on pseudo-label based self-training on target data. Our proposed subsidiary supervised adaptation implicitly regularizes target-side self-training, leading to improved adaptation.

To summarize, our main contributions are:

- We introduce concurrent subsidiary supervised DA, for a subsidiary task, that not only improves unsupervised goal task DA but also facilitates source-free adaptation. We provide theoretical insights to analyze the impact of subsidiary DA on the domain discrepancy, and hence, the goal task DA.
- Based on our insights, we devise a subsidiary DA suitability criteria to identify *DA-assistive* subsidiary tasks that better aid the unsupervised goal task DA. We also propose novel sticker intervention based subsidiary tasks that demonstrate the efficacy of the criteria.
- Our proposed approach achieves state-of-the-art performance on source-free single-source DA (SSDA) as well as source-free multi-source DA (MSDA) for image classification. The proposed approach also complements existing non-source-free works, achieving leading performance.

2 Related Work

Pretext tasks in self-supervised learning. Pretext tasks are used to learn deep feature representations from unlabeled data, in a self-supervised manner,

for downstream tasks. There are several pretext tasks such as image inpainting [52], colorization [82,33,83], spatial context prediction [9], contrastive predictive coding [48], image rotation [16], and jigsaw puzzle solving [47]. Pretext tasks are commonly used for pre-training on unlabeled data followed by finetuning on labeled data. Conversely, we perform supervised DA for the pretext-like task along with the unsupervised goal task DA, resulting in a representation that aligns the domains while maintaining the goal task performance.

Source-free DA. Recently, several methods have investigated source-free DA. USFDA [31] and FS [32] investigate universal DA [80] and open-set DA [60], in a source-free setting by synthesizing training samples to make the decision boundaries compact. SHOT [37,38], NRC [78] maximize mutual information and propose pseudo-labeling, using global structure to match target features to that of a fixed source classifier. To provide adaptation supervision, 3C-GAN [34] generates labeled target-style images from a GAN. Finally, SFDA [39], UR [64], and GtA [30] are semantic segmentation specific source-free DA techniques.

Pretext task based DA. Several DA works have demonstrated the efficacy of learning meaningful representations using pretext tasks. Early works [14,15] used reconstruction as a pretext task to extract domain-invariant features. [5] captured both domain-specific and shared features by separating the feature space into domain-private and domain-shared spaces. [6] used jigsaw puzzles as a side-objective to tackle domain generalization. [66] proposed that adaptation can be accomplished by learning many self-supervision tasks at the same time. [26] suggested a cross-domain SSL strategy for adaptation with minimal source labels based on instance discrimination [74]. [22] recommended employing SSL pretext tasks like rotation prediction and patch placement prediction. [59] solved the challenge of universal domain adaptation by unsupervised clustering. [57] employed easy labels for synthetic images, such as the surface normal, depth, and instance contour, to train a network. [11] employed SSL pretext tasks like rotation prediction as part of their domain generalization technique.

3 Approach

In this section, we introduce required preliminaries (Sec. 3.1), followed by theoretical insights (Sec. 3.2) that motivate our training algorithm design (Sec. 3.4).

3.1 Preliminaries

3.1.1 Goal task unsupervised DA. For closed set DA problem, 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} is the input space and \mathcal{C}_g denotes the label set for the goal task. x_s is drawn from the marginal distribution p_s . Let $\mathcal{D}_t = \{x_t : x_t \in \mathcal{X}\}$ be an unlabeled target dataset with $x_t \sim p_t$. The goal is to assign labels for each target image x_t . The usual approach [13,65,40] is to use a backbone feature extractor $h : \mathcal{X} \rightarrow \mathcal{Z}$ followed by a goal classifier $f_g : \mathcal{Z} \rightarrow \mathcal{C}_g$ (see Fig. 3). The expected source risk with h and an optimal labeling function $f_S : \mathcal{X} \rightarrow \mathcal{C}_g$, is $\epsilon_s(h) = \mathbb{E}_{x \sim p_s} [\mathbb{1}(f_g \circ h(x) \neq f_S(x))]$, where $\mathbb{1}(\cdot)$ is an indicator function. Similarly, $\epsilon_t(h)$ is the target risk with optimal labeling function $f_T : \mathcal{X} \rightarrow \mathcal{C}_g$. We restate the theoretical upper bound on target

risk from [84]. For backbone hypothesis $h \in \mathcal{H}$ with \mathcal{H} being the hypothesis space and a domain classifier $f_d : \mathcal{Z} \rightarrow \{0, 1\}$ (0 for source, 1 for target),

$$\epsilon_t(h) \leq \epsilon_s(h) + d_{\mathcal{H}}(p_s, p_t) + \lambda_g; \quad \lambda_g = \min \left\{ \mathbb{E}_{p_s} [\mathbb{1}(f_s(x) \neq f_T(x))], \mathbb{E}_{p_t} [\mathbb{1}(f_S(x) \neq f_T(x))] \right\}$$

where, $d_{\mathcal{H}}(p_s, p_t) = \sup_{h \in \mathcal{H}} \left| \mathbb{E}_{x \sim p_s} [\mathbb{1}(f_d \circ h(x) = 1)] - \mathbb{E}_{x \sim p_t} [\mathbb{1}(f_d \circ h(x) = 1)] \right|$ (1)

Here, $d_{\mathcal{H}}$ is the \mathcal{H} -divergence [4] that indicates the distribution shift or worst-case domain discrepancy between the two domains. λ_g is a constant that represents the optimal cross-domain error of the labeling functions. Thus, the target risk $\epsilon_t(h)$ is upper bounded by these two terms along with the source risk $\epsilon_s(h)$.

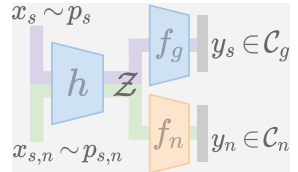


Fig. 3. Our method uses a shared backbone h with goal classifier f_g and subsidiary classifier f_n .

3.1.2 Subsidiary supervised DA. Next, we introduce a subsidiary supervised DA problem concurrent to the goal task unsupervised DA. To this end, we aim to devise a subsidiary classification task with a new label set \mathcal{C}_n . The label-set specific attributes are inflicted on $x \in \mathcal{X}$ via an intervention, to form supervised pairs. These pairs form labeled source, $(x_{s,n}, y_n) \in \mathcal{D}_{s,n}$ and labeled target, $(x_{t,n}, y_n) \in \mathcal{D}_{t,n}$ datasets. Here, the inputs $x_{s,n}$ and $x_{t,n}$ are drawn from marginal distributions $p_{s,n}$ and $p_{t,n}$ respectively. We also define the optimal labeling functions for source and target subsidiary task as $f_{s,n} : \mathcal{X} \rightarrow \mathcal{C}_n$ and $f_{T,n} : \mathcal{X} \rightarrow \mathcal{C}_n$. Next, the prediction mapping involves the shared goal-task backbone h followed by a subsidiary classifier $f_n : \mathcal{Z} \rightarrow \mathcal{C}_n$ (see Fig. 3). Here, the source-subsidiary task error is $\epsilon_{s,n}(h) = \mathbb{E}_{x \sim p_{s,n}} [\mathbb{1}(f_n \circ h(x) \neq f_{s,n}(x))]$. Similarly, $\epsilon_{t,n}(h)$ for target and λ_n defined as in Eq. 1. Thus, generalization bounds for subsidiary DA with the same \mathcal{H} is stated as,

$$\epsilon_{t,n}(h) \leq \epsilon_{s,n}(h) + d_{\mathcal{H}}(p_{s,n}, p_{t,n}) + \lambda_n \quad (2)$$

3.1.3 Metrics. We introduce two metrics that form the basis of our insights.

a) Subsidiary-Domain Similarity Metric (DSM), $\gamma_{DSM}(\cdot, \cdot)$. DSM measures the similarity between two domains as the inverse of the standard \mathcal{A} -distance [4]. \mathcal{A} -distance can be thought of as a proxy [13] for \mathcal{H} -divergence.

b) Subsidiary-Task Similarity Metric (TSM), $\gamma_{TSM}(\cdot, \cdot)$. TSM measures the task similarity of a subsidiary task w.r.t. the goal task. TSM is computed using the standard linear evaluation protocol [62] borrowed from transfer learning and self-supervised literature. It is the performance of a subsidiary-task linear classifier attached to a goal-task pretrained backbone feature extractor $h_{s,g}$. Intuitively, it indicates the extent of compatibility between the two tasks.

For a dataset pair of source-goal and source-subsidiary, *i.e.* $(\mathcal{D}_s, \mathcal{D}_{s,n})$;

$$\gamma_{DSM}(\mathcal{D}_s, \mathcal{D}_{s,n}) = 1 - \frac{1}{2} d_{\mathcal{A}}(\mathcal{D}_s, \mathcal{D}_{s,n}); \quad \gamma_{TSM}(\mathcal{D}_s, \mathcal{D}_{s,n}) = 1 - \min_{f_n} \hat{\epsilon}_{s,n}(h_{s,g}) \quad (3)$$

Here, $d_{\mathcal{A}}(\cdot, \cdot)$ denotes \mathcal{A} -distance and $\hat{\epsilon}_{s,n}(\cdot)$ denotes empirical error for subsidiary task on source data. Note that $0 \leq \hat{\epsilon}_{s,n}(h_{s,g}) \leq 1$ while $0 \leq d_{\mathcal{A}}(\mathcal{D}_1, \mathcal{D}_2) \leq 2$.

3.2 Theoretical insights

We analyze the impact of solving subsidiary supervised DA on the goal task unsupervised DA. We first consider the combined bounds (combining Eq. 1, 2),

$$\epsilon_t(h) + \epsilon_{t,n}(h) \leq \epsilon_s(h) + \epsilon_{s,n}(h) + d_{\mathcal{H}}(p_s, p_t) + d_{\mathcal{H}}(p_{s,n}, p_{t,n}) + \lambda_g + \lambda_n \quad (4)$$

Among the six terms on the right side, the two λ terms are constants as they do not involve the hypothesis h or hypothesis space \mathcal{H} . We analyze the source error duet, $\epsilon_s(h) + \epsilon_{s,n}(h)$, and the domain discrepancy duet $d_{\mathcal{H}}(p_s, p_t) + d_{\mathcal{H}}(p_{s,n}, p_{t,n})$.

3.2.1 Analyzing the domain discrepancy duet. First, we analyze w.r.t. the domain discrepancy duet by considering the following three configurations.

a) While performing only unsupervised goal task DA, the backbone optimization would operate on a limited hypothesis space $\mathcal{H}_g^{(uns)} \subset \mathcal{H}$ where $\mathcal{H}_g^{(uns)} = \{h \in \mathcal{H} : |\epsilon_t(h) - \epsilon_s(h)| \leq \zeta_g^{(uns)}\}$. Here, $\zeta_g^{(uns)}$ is a threshold on the source-target error gap.

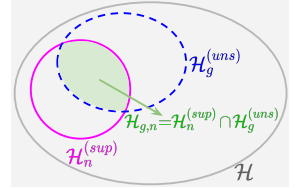


Fig. 4. Three configurations of hypothesis spaces.

b) While performing supervised adaptation only for subsidiary domain adaptation, the optimization would operate on a limited hypothesis space $\mathcal{H}_n^{(sup)} \subset \mathcal{H}$ i.e., $\mathcal{H}_n^{(sup)} = \{h \in \mathcal{H} : |\epsilon_{t,n}(h) - \epsilon_{s,n}(h)| \leq \zeta_n^{(sup)}\}$.

Here, $\zeta_n^{(sup)}$ is a threshold on the subsidiary-task source-target error gap.

c) While concurrently performing **a)** unsupervised goal task DA and **b)** subsidiary supervised DA (i.e. the proposed approach), the optimization would operate on a limited hypothesis space $\mathcal{H}_{g,n} \subset \mathcal{H}$. Specifically, $\mathcal{H}_{g,n} = \mathcal{H}_n^{(sup)} \cap \mathcal{H}_g^{(uns)}$. This is because the backbone is shared between the two DA tasks and hence, would be limited to the intersection space.

Different configurations lead to different \mathcal{H} -spaces and consequently, different \mathcal{H} -divergences. Comparing the \mathcal{H} -divergences leads us to the following insight.

Insight 1. (\mathcal{H} -divergence in concurrent goal DA and subsidiary DA)

The backbone hypothesis space for concurrent unsupervised goal DA and subsidiary supervised DA, i.e. $\mathcal{H}_{g,n} = \mathcal{H}_n^{(sup)} \cap \mathcal{H}_g^{(uns)}$ will yield a lower \mathcal{H} -divergence than $\mathcal{H}_g^{(uns)}$ (hypothesis space for only unsupervised goal task DA), i.e.

$$d_{\mathcal{H}_{g,n}}(p_s, p_t) \leq d_{\mathcal{H}_g^{(uns)}}(p_s, p_t) \quad \text{and} \quad d_{\mathcal{H}_{g,n}}(p_{s,n}, p_{t,n}) \leq d_{\mathcal{H}_g^{(uns)}}(p_{s,n}, p_{t,n}) \quad (5)$$

Remarks. In Eq. 1, $d_{\mathcal{H}}(p_s, p_t)$ is the supremum over the hypothesis space \mathcal{H} i.e. a worst-case measure. Since $\mathcal{H}_{g,n} \subset \mathcal{H}_g^{(uns)}$, $\mathcal{H}_{g,n}$ would have a lower \mathcal{H} -divergence as the worst-case hypothesis of $\mathcal{H}_g^{(uns)}$ may be absent in the subset $\mathcal{H}_{g,n}$. This applies to both pairs, (p_s, p_t) and $(p_{s,n}, p_{t,n})$. While a lower \mathcal{H} -divergence duet leads to improved goal DA, the equality may hold when the worst hypothesis of $\mathcal{H}_g^{(uns)}$ remains in $\mathcal{H}_{g,n}$. In such a case, concurrent DA would perform the same as naive goal DA. To this end, we put forward the following insight.

Insight 2. (When is concurrent DA strictly better than naive goal DA?) A subsidiary task supports the strict inequality $d_{\mathcal{H}_{g,n}}(p_s, p_t) < d_{\mathcal{H}_g^{(uns)}}(p_s, p_t)$ if with at least $(1 - \delta)$ probability, the subsidiary-domain similarity $\gamma_{DSM}(\mathcal{D}_s, \mathcal{D}_{s,n})$ exceeds a threshold ζ_d by no less than ξ ; $\mathbb{P}[\gamma_{DSM}(\mathcal{D}_s, \mathcal{D}_{s,n}) \geq \zeta_d - \xi] \geq 1 - \delta$.

Remarks. In other words, the strict inequalities in Eq. 5 would hold if the DSM $\gamma_{DSM}(\cdot, \cdot)$ exceeds a threshold ζ_d . The supports for this insight are twofold. First, a subsidiary task may heavily alter domain information [44], e.g. jigsaw shuffling [6]. Then, the backbone will be updated using out-of-domain samples which is undesirable as such samples are unlikely for inference. This will be avoided if Insight 2 is satisfied. Second, if DSM is high, we can approximate $p_s \approx p_{s,n}$ and $p_t \approx p_{t,n}$. Thus, more samples from subsidiary task data will be available for training the backbone to be domain-invariant (as subsidiary task uses samples from both the domains) i.e. reducing $d_{\mathcal{H}}$ against the same in naive goal DA.

3.2.2 Analyzing the source error duet. Now we analyze w.r.t. the source error duet of Eq. 4. While the \mathcal{H} -divergence is lower for concurrent goal task DA and subsidiary supervised DA, a logical concern is that simultaneous minimization of errors, i.e. $\epsilon_s(h) + \epsilon_{s,n}(h)$, for both tasks may be difficult with the shared backbone h . Further, it may happen that simultaneous training for both tasks in target domain may hamper the goal task performance as it is unsupervised. In such cases, the subsidiary task would be ill-equipped to assist the goal task adaptation. To avoid these, we propose another empirical criterion as follows.

Insight 3. (Goal and subsidiary task similarity for concurrent DA) Higher goal-subsidiary task similarity (TSM) aids effective minimization of both task errors with the shared backbone, which is crucial for subsidiary supervised DA to positively affect the goal task DA. The criterion is $\gamma_{TSM}(\mathcal{D}_s, \mathcal{D}_{s,n}) > \zeta_n$.

Remarks. Here, ζ_n is a threshold. The TSM γ_{TSM} indicates the compatibility of goal task features to support the subsidiary task. Intuitively, a higher TSM implies more overlap in the discriminative features of the two tasks, which would allow better simultaneous minimization of both task errors.

Based on Insight 1, concurrent subsidiary supervised DA and goal task DA yields a lower domain discrepancy. Further, based on Insight 2, a subsidiary task can be selected such that effective minimization of both source errors is possible simultaneously. Thus, using Eq. 1, we can infer that $\sup_{h \in \mathcal{H}_{g,n}} \epsilon_t(h) \leq \sup_{h \in \mathcal{H}_g^{(uns)}} \epsilon_t(h)$ i.e. a lower target error upper bound for our approach w.r.t. naive goal task DA. Now, we summarize the criteria (Insight 2, 3).

Definition 1. (Subsidiary DA suitability criteria) A subsidiary task is termed DA-assistive i.e. suitable for subsidiary supervised DA if the sum of DSM γ_{DSM} and TSM γ_{TSM} exceeds a threshold ζ ,

$$\gamma_{DSM}(\mathcal{D}_s, \mathcal{D}_{s,n}) + \gamma_{TSM}(\mathcal{D}_s, \mathcal{D}_{s,n}) > \zeta \quad (6)$$

Remarks. In other words, a subsidiary task which is domain-preserving and has high task similarity w.r.t. the goal task is DA-assistive i.e. suitable for subsidiary supervised DA to aid the goal task DA. We employ this criteria empirically for a diverse set of subsidiary tasks (shown in Fig. 1C). Next, we describe the

motivation for our proposed sticker intervention and corresponding subsidiary tasks as well as training algorithms tailored for source-free DA.

3.3 Sticker intervention based subsidiary task design

While one may consider pretext tasks from the self-supervised learning literature as candidates for subsidiary DA, almost all such tasks fail to satisfy subsidiary DA suitability criteria in Eq. 6. For instance, dense output based tasks such as colorization [82,33], inpainting [52], *etc.* exhibit markedly low task similarity (TSM) against the non-dense goal tasks. Further, the input intervention for certain pretext tasks such as jigsaw [6], patch-location[66], rotation [43,22], significantly alter the domain information leading to low domain similarity (DSM).

Insight 4. (Sticker-intervention based tasks well suit subsidiary DA)

Sticker intervention is the process of pasting a sticker x_n (i.e., a symbol with random texture and scale) on a given image sample $x_s \in \mathcal{D}_s$ to obtain a stickered sample, i.e. $x_{s,n} = \mathcal{T}(x_s, x_n) \in \mathcal{D}_{s,n}$. Following this, the subsidiary task could be defined as the classification of some sticker attribute (e.g. shape, location, or orientation). Such a formalization provides effective control to maximize $\gamma_{DSM}(\mathcal{D}_s, \mathcal{D}_{s,n})$ and $\gamma_{TSM}(\mathcal{D}_s, \mathcal{D}_{s,n})$, in line with our suitability criteria.

Remarks. The sticker intervention (Fig. 5) facilitates domain preservation while simultaneously supporting a range of subsidiary tasks. Since the proposed sticker intervention alters only a local area of the sample, the original content is not suppressed which in turn preserves the domain information, implying high DSM.



Fig. 5. Sticker intervention.

Following this, one can ablate over a range of sticker-based tasks in order to select a suitable subsidiary task based on the given goal task. Below, we discuss some possible subsidiary tasks under the sticker intervention.

a) Sticker location (Fig. 2A). We draw motivation from patch-location [66], where the task is to classify the quadrant to which a patch-input belongs. With sticker intervened images, the task is to classify the quadrant with the sticker. Our use of whole images as input is more domain-preserving than patch-input.

b) Sticker rotation (Fig. 2B). Motivated by the image rotation task [43], we propose sticker rotation task where the rotation of the sticker has to be classified (0° , 90° , 180° and 270° rotations possible). Note that our sticker rotation does not affect the domain information while rotating the entire image does.

c) Sticker classification (Fig. 2C). While the discriminative features in the previous two tasks were location and rotation, we propose sticker classification task with primary discriminative features as shape. In other words, the task is to classify the sticker shape (*i.e.* the symbol) given a stickered sample.

3.4 Training algorithm design under source-free constraints

For the standard DA setting with concurrent access to source and target data [13,65], the subsidiary supervised DA can be implemented simply by optimizing the subsidiary classification loss simultaneously for source and target. This would yield a lower domain discrepancy as discussed in Sec. 3.2. However, in the

more practical source-free setting [34,31] where concurrent source-target access is prohibited, this simple approach would not be possible. We believe the improvements will be prominent in source-free DA based on the following insight:

Insight 5. (Subsidiary DA better suits challenging source-free DA). Existing source-free DA works heavily rely on pseudo-label or clustering based self-training on unlabeled target with no obvious alternative. The proposed subsidiary supervised adaptation helps to implicitly regularize the target-side self-training, leading to improved adaptation performance. The subsidiary DA not only aids goal DA as a result of high DSM but also preserves the goal task inductive bias as a result of high TSM, while adhering to the source-free constraints.

Remarks. The source-free setting presents new challenges which highlight the advantages of our proposed method more prominently. This is because, the performance in non-source-free DA is strongly influenced by well-developed discrepancy minimization techniques. However, these techniques cannot be leveraged in a source-free setting due to their requirement of concurrent source-target data access. Thus, we primarily operate in the source-free regime to evaluate our theoretical insights and the proposed concurrent subsidiary supervised DA problem.

We perform the training in three steps. First two steps involve pre-training of goal task and subsidiary task respectively with source data. The final step involves adapting both tasks to target domain. For clarity, we first summarize available and intervened datasets required for training and their notations.

Datasets. The goal task source data is denoted by $(x_s, y_s) \in \mathcal{D}_s$ while the corresponding unlabeled target is denoted by $x_t \in \mathcal{D}_t$. The intervened stickered-source data, coupled with both goal and sticker task labels, is denoted by $(x_{s,n}, y_s, y_n) \in \mathcal{D}_{s,n}$. The corresponding stickered-target data, with only subsidiary sticker task labels, is denoted by $(x_{t,n}, y_n) \in \mathcal{D}_{t,n}$. We introduce a pseudo-OOS (out-of-source) dataset, $\mathcal{D}_s^{(od)}$ further in this section.

3.4.1. Goal task source pre-training (Fig. 6A). We train the backbone h and goal classifier f_g with source data \mathcal{D}_s and stickered-source data $\mathcal{D}_{s,n}$:

$$\min_{\theta_h, \theta_{f_g}} \mathbb{E}_{(x,y) \in \mathcal{D}_s \cup \mathcal{D}_{s,n}} [\mathcal{L}_{s,g}]; \quad \mathcal{L}_{s,g} = \mathcal{L}_{ce}(f_g \circ h(x), y) \quad (7)$$

Here, θ_h and θ_{f_g} are the parameters of h and f_g , \mathcal{L}_{ce} is the cross-entropy loss, y is the goal task label, and expectation is implemented by sampling mini-batches.

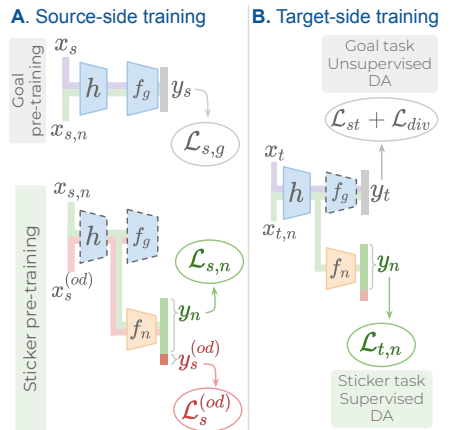


Fig. 6. **A.** Source-side training involves goal pre-training (Sec. 3.4.1) and sticker pre-training (Sec. 3.4.2). **B.** Target-side training involves concurrent goal-task unsupervised DA and sticker-task supervised DA (Sec. 3.4.3).

3.4.2. Sticker task source pre-training (Fig. 6A). We pretrain the sticker classifier f_n while inculcating the ability to reject samples out of the source distribution. Specifically, f_n predicts a $(|\mathcal{C}_n| + 1)$ -sized vector and is trained to classify *out-of-source* (OOS) samples to the $(|\mathcal{C}_n| + 1)^{\text{th}}$ class.

Insight 6. *The OOS node in the sticker classifier implicitly behaves as a domain discriminator from adversarial alignment methods. Minimizing the OOS probability only for the target data aligns the target with the source.*

Remarks. In source training, the OOS objective forces the sticker classifier to discriminate between source and OOS samples. This is done with the intuition that OOS samples simulate the role of target samples in adversarial alignment methods. This domain discriminatory knowledge will support future source-free target alignment. Concretely, the shared backbone can be adapted to the target, by minimizing OOS probability for target samples, as source knowledge is preserved by freezing f_g . Thus, we require OOS data to prepare f_n for adaptation.

Obtaining the OOS dataset. The naive approach is to use a dataset unrelated to the goal task label set. Conversely, we devise a pseudo-OOS dataset using only the already available source samples. Mitsuzumi *et al.* [44] show that, beyond a certain grid size, shuffling the grid patches makes the domain unrecognizable. Inspired by this, we generate the pseudo-OOS dataset by shuffling the grid patches of source images. We also perform the sticker intervention on shuffled images, at random, to further instill the differences between source and pseudo-OOS samples (see Suppl). Formally, $(x_s^{(od)}, y_s^{(od)}) \in \mathcal{D}_s^{(od)}$ where $y_s^{(od)}$ denotes the OOS category *i.e.* the $(|\mathcal{C}_n| + 1)^{\text{th}}$ category of f_n .

We train only the sticker classifier f_n , keeping backbone h and goal classifier f_g frozen, using cross-entropy loss \mathcal{L}_{ce} . With $\mathcal{L}_{s,n} = \mathcal{L}_{ce}(f_n \circ h(x_{s,n}), y_n)$, the overall objective for stickered source data $\mathcal{D}_{s,n}$ and pseudo-OOS data $\mathcal{D}_s^{(od)}$ is,

$$\min_{\theta_{f_n}} \mathbb{E}_{\mathcal{D}_{s,n}} [\mathcal{L}_{s,n}] + \mathbb{E}_{\mathcal{D}_s^{(od)}} [\mathcal{L}_s^{(od)}]; \quad \text{where } \mathcal{L}_s^{(od)} = \mathcal{L}_{ce}(f_n \circ h(x_s^{(od)}), y_s^{(od)}) \quad (8)$$

3.4.3. Source-free target adaptation (Fig. 6B). For unsupervised goal task adaptation, we use the general self training loss \mathcal{L}_{st} and diversity loss \mathcal{L}_{div} [37]. See Suppl. for more details. The goal task objective is given in Eq. 9 (left),

$$\min_{\theta_h} \mathbb{E}_{\mathcal{D}_t \cup \mathcal{D}_{t,n}} [\mathcal{L}_{st} + \mathcal{L}_{div}] \quad \text{and} \quad \min_{(\theta_h, \theta_{f_n})} \mathbb{E}_{\mathcal{D}_{t,n}} [\mathcal{L}_{t,n}]; \quad \mathcal{L}_{t,n} = \mathcal{L}_{ce}(f_n \circ h(x_{t,n}), y_n) \quad (9)$$

The goal classifier f_g is frozen to preserve its inductive bias and only the backbone h is updated for both original and stickered samples in Eq. 9 (left).

For subsidiary supervised sticker adaptation, we use a simple cross-entropy loss with sticker labels. We implicitly minimize OOS probability by maximizing label class probability. We observe that this works well and explicit minimization of OOS probability is not required. As per Insight 6, *out-of-target* (OOT) samples are not required. Further, using OOT samples to update the backbone could be undesirable as discussed under Insight 2. The objective is given in Eq. 9 (right). Both backbone h and sticker classifier f_n are updated as the task is supervised.

Table 1. Single-Source Domain Adaptation (SSDA) on Office-Home benchmarks. SF indicates *source-free* adaptation.

Method	SF	Office-Home													Avg
		Ar→Cl	Ar→Pr	Ar→Rw	Cl→Ar	Cl→Pr	Cl→Rw	Pr→Ar	Pr→Cl	Pr→Rw	Rw→Ar	Rw→Cl	Rw→Pr		
FixBi [45]	✗	58.1	77.3	80.4	67.7	79.5	78.1	65.8	57.9	81.7	76.4	62.9	86.7	72.7	
SENTRY[55]	✗	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 [35]	✗	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	
SHOT [37]	✓	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	
A ² Net [75]	✓	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 [79]	✓	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 [56]	✓	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 [78]	✓	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++[38]	✓	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	
<i>Ours</i>	✓	61.0	80.4	82.5	69.1	79.9	79.5	69.1	57.8	82.7	74.5	65.1	86.4	74.0	

4 Experiments

We provide the implementation details of our experiments and thoroughly evaluate our approach w.r.t. state-of-the-art prior works across multiple settings. Unless mentioned, *Ours* implies sticker classification as the subsidiary task.

4.1 Experimental setup

Datasets. We evaluate the effectiveness of our approach on four standard DA benchmarks. **Office-31** [58] benchmark consist of three domains under office environments: Amazon (**A**), DSLR (**D**), and Webcam (**W**), each with 31 object categories. **Office-Home** [70] is a more challenging dataset. It comprises of images of commonplace objects divided into four domains: Artistic (**Ar**), Clipart (**Cl**), Product (**Pr**), and Real-World (**Rw**), each with 65 classes. **VisDA** [54] is a large-scale dataset for synthetic-to-real domain adaptation. The source domain has 152,397 synthetic images, while the target domain has 55,388 real-world images. **DomainNet** [53] is the most challenging due to its highly diverse domains and huge class imbalance. It has 6 domains: Clipart (**C**), Real (**R**), Infograph (**I**), Painting (**P**), Sketch (**S**) and Quickdraw (**Q**) with 345 classes each.

Implementation details. We use a ResNet-50 [18] backbone for Office-Home, Office-31 and DomainNet, and ResNet-101 for VisDA, for a fair comparison with prior works. We employ the same network design as SHOT [37], *i.e.* replacing the classifier with a fully connected layer with batch norm [21] and another fully connected layer with weight normalization [61]. For the subsidiary classifier, we use the same architecture after ResLayer-3. The number of sticker classes is 10. See Suppl. for more details related to sticker intervention and sticker-based tasks. Following [2,37], we use label smoothing for source training using Adam [28] with learning rate 1e-3, momentum 0.9, and batch size 64. We use separate Adam optimizers for each loss term to avoid loss balancing hyperparameters.

4.2 Discussion

We provide an extensive ablation study of both the *source-side* and *target-side* training. Further, we show that our approach is compatible with existing non-source-free DA works and achieves faster and improved convergence.

Table 2. Multi-Source Domain Adaptation (MSDA) on DomainNet and Office-Home. We outperform *source-free* (denoted by SF) prior arts despite not using domain labels.

Method	SF	w/o Domain Labels	DomainNet								Office-Home				
			→C	→I	→P	→Q	→R	→S	Avg	→Ar	→Cl	→Pr	→Rw	Avg	
WAMDA [1]	✗	✗	59.3	21.8	52.1	9.5	65.0	47.7	42.6	71.9	61.4	84.1	82.3	74.9	
SImpAl ₅₀ [69]	✗	✗	66.4	26.5	56.6	18.9	68.0	55.5	48.6	70.8	56.3	80.2	81.5	72.2	
CMSDA [63]	✗	✗	70.9	26.5	57.5	21.3	68.1	59.4	50.4	71.5	67.7	84.1	82.9	76.6	
DRT [36]	✗	✗	71.0	31.6	61.0	12.3	71.4	60.7	51.3	-	-	-	-	-	
STEM [46]	✗	✗	72.0	28.2	61.5	25.7	72.6	60.2	53.4	-	-	-	-	-	
Source-combine	✗	✓	57.0	23.4	54.1	14.6	67.2	50.3	44.4	58.0	57.3	74.2	77.9	66.9	
SHOT [37]-Ens	✓	✗	58.6	25.2	55.3	15.3	70.5	52.4	46.2	72.2	59.3	82.8	82.9	74.3	
DECISION [2]	✓	✗	61.5	21.6	54.6	18.9	67.5	51.0	45.9	74.5	59.4	84.4	83.6	75.5	
SHOT++ [38]	✓	✗	-	-	-	-	-	-	-	73.1	61.3	84.3	84.0	75.7	
CAiDA [10]	✓	✗	-	-	-	-	-	-	-	75.2	60.5	84.7	84.2	76.2	
NRC [78]	✓	✓	65.8	24.1	56.0	16.0	69.2	53.4	47.4	70.6	60.0	84.6	83.5	74.7	
<i>Ours</i>	✓	✓	70.3	25.7	57.3	17.1	69.9	57.1	49.6	75.1	64.1	86.6	84.4	77.6	

4.2.1 Comparison with prior arts

a) Single Source Domain Adaptation (SSDA). We compare our proposed approach with prior source-free SSDA works in Table 1 and 4. Our approach outperforms source-free NRC [78] and SHOT++ [38] by 1.5% and 1.7% respectively on Office-31 (Table 4), and gives comparable performance to non-source-free works. On the larger and more challenging VisDA dataset, our approach surpasses NRC by 1.6% and SHOT++ by 1% (Table 4). On Office-Home (Table 1), our model achieves *state-of-the-art* results exceeding the source-free SHOT++ and the non-source-free method SCDA [35] by 1% and 0.9% respectively.

b) Multi Source Domain Adaptation (MSDA). In Table 2, we compare with the source-only baseline (*source-combine*) and source-free works. Even without domain labels, our approach achieves *state-of-the-art* results, even w.r.t. non-source-free works on Office-Home (+1%). On DomainNet, we outperform source-free works (+2.2%) with comparable results to non-source-free works.

c) Evaluating the subsidiary DA suitability criteria. We empirically evaluate the DSM and TSM for our sticker-based tasks as well as existing tasks borrowed from self-supervised literature in Fig. 7A, 7B. Compared to patch location [66] and image rotation [43], sticker location and sticker rotation tasks exhibit higher DSM and thus, are more suitable with better adaptation performance (also see Table 3). However, the sticker classification task is the most suitable due to its higher TSM as shape is the primary discriminative features, same as in the goal task. We observe a positive correlation between DA performance and both DSM and TSM, which empirically verifies our suitability criteria. In Table 3, we additionally compare dense output based tasks like colorization and inpainting, which give marginal gains compared to other subsidiary tasks.

Table 3. Subsidiary task comparisons on **Office-Home** for source-free DA. Here, baseline is same as #3 in Table 6.

Method	SSDA	MSDA
Baseline (B)	66.2	74.3
B + inpainting	66.3	74.5
B + colorization	66.8	74.7
B + jigsaw	67.0	74.8
B + patch-loc	67.6	75.0
B + rotation	67.9	75.4
B + sticker-loc	68.8	75.5
B + sticker-rot	69.0	75.7
B + sticker-clsf	69.7	76.2

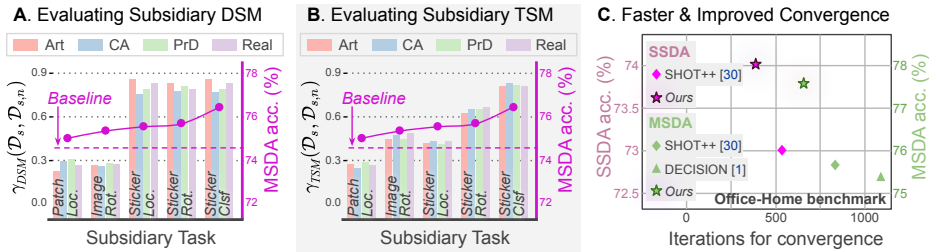


Fig. 7. We observe higher **A.** domain similarity (DSM) and **B.** task similarity (TSM) for our sticker-based tasks compared to existing subsidiary tasks like patch-location and image-rotation. This correlates with the better MSDA performance of sticker-based tasks on Office-Home and validates our criteria (Definition 1). **C.** Faster and improved convergence w.r.t. prior source-free works on both SSSDA and MSDA for Office-Home.

Table 4. Single-Source DA (SSDA) on Office-31 and VisDA. SF indicates *source-free*.

Method	SF	Office-31							VisDA	
		A→D	A→W	D→W	W→D	D→A	W→A	Avg	S → R	
CAN [25]	✗	95.0	94.5	99.1	99.8	78.0	77.0	90.6	87.2	
FixBi [45]	✗	95.0	96.1	99.3	100.0	78.7	79.4	91.4	87.2	
CDAN+RADA [23]	✗	96.1	96.2	99.3	100.0	77.5	77.4	91.1	76.3	
RFA [3]	✗	93.0	92.8	99.1	100.0	78.0	77.7	90.2	79.4	
SHOT [37]	✓	94.0	90.1	98.4	99.9	74.7	74.3	88.6	82.9	
CPGA [56]	✓	94.4	94.1	98.4	99.8	76.0	76.6	89.9	84.1	
HCL [19]	✓	90.8	91.3	98.2	100.0	72.7	72.7	87.6	83.5	
VDM-DA [68]	✓	93.2	94.1	98.0	100.0	75.8	77.1	89.7	85.1	
A ² Net [75]	✓	94.5	94.0	99.2	100.0	76.7	76.1	90.1	84.3	
NRC [78]	✓	96.0	90.8	99.0	100.0	75.3	75.0	89.4	85.9	
SHOT++ [38]	✓	94.3	90.4	98.7	99.9	76.2	75.8	89.2	87.3	
<i>Ours</i>	✓	96.1	94.5	99.2	100.0	77.1	78.5	90.9	88.2	

d) Faster and improved convergence. Fig. 7C illustrates the improved and faster convergence of our approach compared to source-free prior arts for both SSSDA and MSDA. The hypothesis space for concurrent subsidiary supervised DA and unsupervised goal task DA, $\mathcal{H}_{g,n}$, is a subset of the hypothesis space for only unsupervised goal task DA, $\mathcal{H}_g^{(uns)}$. Thus, we observe faster convergence for our approach. Further, as per Insight 1, the lower domain discrepancy leads to a lower target error *i.e.* improved convergence.

4.2.2 Ablation Study. Below, we discuss a thorough ablation study.

a) Effect of subsidiary supervised DA and OOS node. In Table 6, we compare the baseline *i.e.* only unsupervised goal task DA (#3) with the addition of only OOS classifier (#4). Here, a binary classifier is used for OOS detection. We observe gains of 0.8% and 0.6% for SSSDA and MSDA respectively. This indicates that only OOS helps, but subsidiary classifier is essential for further improvements. Next, we compare the baseline (#3) with concurrent goal-subsubsidiary DA without using OOS (#5). We observe an improvement of 3.5% and 1.9% for SSSDA and MSDA. Adding the OOS objective to the subsidiary supervised DA (#6 vs. #4) improves the source-target alignment as explained in Insight 6, resulting in improvements of 3.1% and 1.4% for SSSDA and MSDA.

Table 5. Evaluating compatibility of subsidiary DA with non-source-free DA works on Office-Home. SSDA and MSDA indicate single-source and multi-source DA.

Method	Office-Home	
	SSDA	MSDA
CDAN [41]	65.8	69.4
+ <i>Subsidiary-DA</i>	67.1	71.2
SRDC [67]	71.3	73.1
+ <i>Subsidiary-DA</i>	71.9	75.2
FixBi [45]	72.7	-
+ <i>Subsidiary-DA</i>	73.7	-
CMSDA [63]	-	76.6
+ <i>Subsidiary-DA</i>	-	78.1

Table 6. Ablation analysis. Here, *sticker-w-OOS-clsf* denotes learning with all the proposed components unlike in *only-OOS-clsf* (all losses except $\mathcal{L}_{s,n}, \mathcal{L}_{t,n}$) and *only-sticker-clsf* (all losses except $\mathcal{L}_s^{(od)}$). SF denotes source-free constraint.

# Variation	SF	Office-Home	
		SSDA	MSDA
1. Source-only baseline	-	60.2	66.9
2. + <i>sticker-w-OOS-clsf</i>	-	61.9	71.4
3. Adaptation baseline (B)	✓	66.2	74.3
4. B + <i>only-OOS-clsf</i>	✓	67.0	74.9
5. B + <i>only-sticker-clsf</i>	✓	69.7	76.2
6. B + <i>sticker-w-OOS-clsf</i>	✓	73.1	77.6
7. B + <i>sticker-w-OOS-clsf</i>	✗	74.5	78.3

b) Subsidiary-goal task similarity. As per Insight 3, higher goal-subsidiary task similarity is important for effective learning of both tasks. Thus, in Table 6, we compare the source-only baseline (#1) with only subsidiary supervised DA without goal task target adaptation (#2). We observe gains of 1.7% and 1.3% for SSDA and MSDA respectively. This illustrates the positive correlation between sticker classification and goal task even when target goal losses are not used.

4.2.3 Compatibility with non-source-free DA. In Table 5, we evaluate the compatibility of concurrent subsidiary supervised DA with existing non-source-free SSDA techniques [13,41,67]. MSDA results are obtained by combining the multiple sources for each target. Compared to the original reported results, all four perform better with our proposed subsidiary DA. Note that our non-source-free variant outperforms these results (#7 in Table 6).

5 Conclusion

In this work, we introduced concurrent subsidiary supervised DA for a pretext-like task to aid the unsupervised goal task DA. We provide theoretical insights to analyze the effect of subsidiary supervised DA on the domain discrepancy and consequently on the goal task adaptation. Based on the insights, we introduce a subsidiary DA suitability criteria to determine DA-assistive subsidiary tasks that improve the goal task DA performance. We also propose a novel sticker intervention based pretext task that follows our criteria. The proposed approach outperforms prior state-of-the-art source-free SSDA and MSDA works on four standard benchmarks, establishing the usefulness of our approach.

Acknowledgments. This work was supported by MeitY (Ministry of Electronics and Information Technology) project (No. 4(16)2019-ITEA), Govt. of India and a research grant by Google.

Supplementary Material: Concurrent Subsidiary Supervision for Unsupervised Source-Free Domain Adaptation

Supplementary Video

We provide a high-level summary video at <https://youtu.be/ENJMz-Eg87k>. We visually demonstrate the key insights of our work as well as illustrate the different subsidiary tasks and training algorithm used. We encourage the reader to go through the video for a better understanding of the key ideas.

Supplementary Document

In this document, we provide extensive implementation details, additional performance analysis and ablation studies. Towards reproducible research, we release our complete codebase and trained network weights at <https://github.com/val-iisc/StickerDA>. This supplementary is organized as follows:

- Section **A**: Notations (Table 7)
- Section **B**: Approach (Algo. 1)
 - Target adaptation (Sec. **B.1**)
 - Subsidiary DA suitability criteria (Sec. **B.2**)
- Section **C**: Implementation details
 - Sticker intervention (Sec. **C.1**, Fig. 8, 9)
 - Experimental settings (Sec. **C.2**)
- Section **D**: Analysis
 - Extended comparisons (Sec. **D.1**, Table 12, 8, 9)
 - Hyperparam. sensitivity (Sec. **D.2**, Table 10, Fig. 11, 10)
 - Domain discrepancy analysis (Sec. **D.3**, Fig. 10)
 - Domain alignment analysis (Sec. **D.4**, Fig. 10)
 - Efficiency analysis (Sec. **D.5**, Table 11)
 - Combining subsidiary tasks (Sec. **D.6**, Table 13)
 - Differences and relationships with prior-arts (Sec. **D.7**, Table 14, 15)

A Notations

We summarize the notations used in the paper in Table 7. The notations are listed under 5 groups: Models, Preliminaries, Datasets, Samples, and Spaces.

B Approach

We summarize our approach in Algo. 1 and provide details of the target adaptation objectives that were omitted from the main paper due to space constraints.

Table 7. Notation Table

	Symbol	Description
Models	h	Shared backbone feature extractor
	f_g	Goal task classifier
	f_n	Subsidiary task classifier
Preliminaries	p_s	Source marginal distribution
	p_t	Target marginal distribution
	ϵ_s	Source goal task error
	ϵ_t	Target goal task error
	$\epsilon_{s,n}$	Source subsidiary task error
	$\epsilon_{s,n}$	Target subsidiary task error
	$d_{\mathcal{H}}$	\mathcal{H} -divergence
	\mathcal{H}	Backbone hypothesis space
	$\mathcal{H}_g^{(uns)}$	\mathcal{H} -space for unsup. goal task
$\mathcal{H}_n^{(sup)}$	\mathcal{H} -space for sup. subsidiary task	
Datasets	\mathcal{D}_s	Labeled source dataset
	\mathcal{D}_t	Unlabeled target dataset
	$\mathcal{D}_{s,n}$	Subsidiary source dataset
	$\mathcal{D}_{t,n}$	Subsidiary target dataset
	$\mathcal{D}_s^{(od)}$	Pseudo-OOS dataset
Samples	(x_s, y_s)	Labeled source sample
	$(x_{s,n}, y_s, y_n)$	Labeled subsidiary source sample
	$(x_s^{(od)}, y_s^{(od)})$	Labeled pseudo-OOS sample
	x_t	Unlabeled target sample
	$(x_{t,n}, y_n)$	Subsidiary target sample
Spaces	\mathcal{X}	Input space
	\mathcal{Z}	Backbone feature space
	\mathcal{C}_g	Label set for goal task
	\mathcal{C}_n	Label set for subsidiary task

B.1 Target adaptation

Self-training loss. We apply self-supervision in the target domain to cluster target samples based on their neighborhood [78]. Each target sample in the feature space is aligned with its neighbor. As a result, the model learns a discriminative metric that translates a point to a semantically similar match. This is accomplished by reducing the entropy over point similarity. The model learns tightly clustered features as it moves neighboring points closer together, resulting in discriminative decision boundaries.

For each mini-batch of target features, we calculate the similarity to all target samples. Let $F_t^{(mb)} \in \mathbb{R}^{|\mathcal{D}_t| \times d}$ denote the memory bank which stores all target features and d denotes the dimensions for output features $f_g \circ h(x_t)$. Here, $|\mathcal{D}_t|$

denotes the number of samples in the target dataset. All stored features are L2-normalized. Specifically,

$$F_t^{(mb)} = [F_1, F_2, \dots, F_{|\mathcal{D}_t|}] \quad (10)$$

where F_j denotes the j^{th} item in $F_t^{(mb)}$. Let $f_i = h(x_i)$ denote the features of the current i^{th} mini-batch, and B_t denote the set of indices of the mini-batch samples in $F_t^{(mb)}$. The probability that f_i is a neighbor of the feature F_j is,

$$p_{i,j} = \frac{\exp(F_j^T f_i / \mathcal{T})}{\sum_{j=1, j \neq i}^{|\mathcal{D}_t|} \exp(F_j^T f_i / \mathcal{T})} \quad (11)$$

where the temperature parameter \mathcal{T} controls the number of neighbors. Then, the entropy *i.e.* the loss is defined as,

$$\mathcal{L}_{st} = -\frac{1}{|B_t|} \sum_{i \in B_t} \sum_{j=1, j \neq i}^{\mathcal{D}_t} p_{i,j} \log(p_{i,j}) \quad (12)$$

Diversity loss. We encourage the prediction to be balanced to avoid degenerate solutions, where the model predicts all data to a particular class (and does not predict other classes for any target sample). We employ the prediction diversity loss, which has been frequently used in clustering [17] and domain adaptation [37]. The diversity objective is,

$$\mathcal{L}_{div}(f_g \circ h(x)) = D_{KL}(\hat{p}, \frac{1}{|\mathcal{C}_g|} \mathbb{1}_{|\mathcal{C}_g|}) - \log |\mathcal{C}_g| \quad (13)$$

where $\mathbb{1}_{|\mathcal{C}_g|}$ represents a $|\mathcal{C}_g|$ -dimensional vector of ones, $\hat{p} = \mathbb{E}_{x_t \in \mathcal{D}_t} [\sigma(f_g \circ h(x_t))]$ is average output embedding for entire target dataset, and σ denotes softmax.

B.2 Subsidiary DA suitability criteria

B.2.1 Subsidiary-Domain Similarity Metric (DSM). As discussed in Sec. 3.1.3 of the main paper, we define *subsidiary-domain similarity metric*, γ_{DSM} as the inverse of the \mathcal{H} -divergence between the two domains. We follow [13] and use the \mathcal{A} -distance [4] between the goal task dataset \mathcal{D}_s and the subsidiary task dataset $\mathcal{D}_{s,n}$ as a proxy for \mathcal{H} -divergence. We define the dataset labels as 1 for subsidiary source dataset $\mathcal{D}_{s,n}$ and 0 for original source dataset \mathcal{D}_s and train a linear binary classifier on the features of a frozen ImageNet-pretrained [51] ResNet-50 [18] with a subset of the mixed data, and obtain the classifier error on the other subset as ψ . The DSM is then computed as,

$$d_{\mathcal{A}}(\mathcal{D}_s, \mathcal{D}_{s,n}) = 2\psi(1 - \psi) \quad (14)$$

$$\gamma_{DSM}(\mathcal{D}_s, \mathcal{D}_{s,n}) = 1 - \frac{1}{2} d_{\mathcal{A}}(\mathcal{D}_s, \mathcal{D}_{s,n}) \quad (15)$$

Algorithm 1 Pseudo-code for the proposed approach

Source-side training

- 1: **Input:** source data \mathcal{D}_s , stickered source data $\mathcal{D}_{s,n}$, pseudo-OOS dataset $\mathcal{D}_s^{(od)}$, ImageNet pretrained backbone h (as per [37]), randomly initialized goal classifier f_g and randomly initialized sticker classifier f_n .

Goal task source pre-training

- 2: **for** $iter < MaxIter$ **do**:
 3: Sample batch from $\mathcal{D}_s \cup \mathcal{D}_{s,n}$
 4: Compute $\mathcal{L}_{s,g}$ using Eq. 7 (main paper)
 5: **update** θ_h, θ_{f_g} by minimizing $\mathcal{L}_{s,g}$
 6: **end for**

Sticker task source pre-training

- 7: **for** $iter < MaxIter$ **do**:
 8: Sample batch from $\mathcal{D}_{s,n}$
 9: Sample batch from $\mathcal{D}_s^{(od)}$
 10: Compute $\mathcal{L}_{s,n}$ and $\mathcal{L}_s^{(od)}$ using Eq. 8 (main paper)
 ▷ using samples from $\mathcal{D}_{s,n}$ and $\mathcal{D}_s^{(od)}$ respectively
 11: **update** θ_{f_n} by minimizing $\mathcal{L}_{s,n}, \mathcal{L}_s^{(od)}$ using separate Adam optimizers
 12: **end for**

Target-side training

- 13: **Input:** target data \mathcal{D}_t , stickered target data $\mathcal{D}_{t,n}$, source-side pretrained backbone h , goal classifier f_g and sticker classifier f_n .

Source-free target adaptation

- 14: **for** $iter < MaxIter$ **do**:
 15: Sample batch from \mathcal{D}_t
 16: Sample batch from $\mathcal{D}_{t,n}$
 17: Compute \mathcal{L}_{st} and \mathcal{L}_{div} using Eq. 12, 13 (suppl.)
 ▷ using samples from both \mathcal{D}_t and $\mathcal{D}_{t,n}$
 18: Compute $\mathcal{L}_{t,n}$ using Eq. 9 (main paper)
 ▷ using samples from only $\mathcal{D}_{t,n}$
 19: **update** θ_h, θ_{f_n} by minimizing $\mathcal{L}_{t,n}$
 20: **update** θ_h by minimizing $\mathcal{L}_{st}, \mathcal{L}_{div}$ using separate Adam optimizers
 21: **end for**
-

How to choose the threshold ζ_d ? Insight 2 introduced a threshold ζ_d for DSM to select pretext tasks suitable for subsidiary supervised DA. To choose a threshold, we first consider the \mathcal{A} -distances between the actual source and target domains. These \mathcal{A} -distances are in the range of 1.5 to 2.0 [69] for Office-Home and indicate the range of \mathcal{A} -distances corresponding to realistic domain shifts. This range corresponds to the range of 0 to 0.25 in terms of DSM. In Fig. 7A of the main paper, we observed DSM in a range of 0 to 0.3 for the patch-location and image-rotation subsidiary task samples w.r.t. the original samples, indicating that these tasks induce a realistic domain shift. Contrary to this, our proposed sticker task produced DSM in the range of 0.6 to 0.9, indicating much better

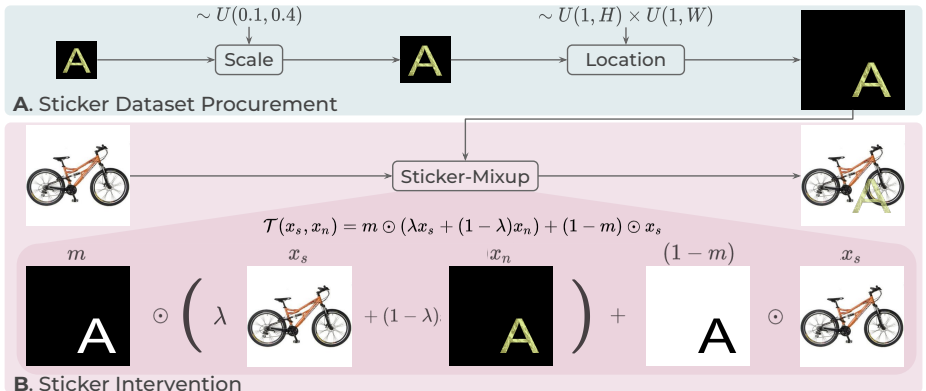


Fig. 8. Illustration of **A.** sticker dataset procurement and **B.** sticker intervention \mathcal{T} (see Sec. C.1). *Best viewed in color.*

domain preservation. Thus, we choose the threshold $\zeta_d = 0.5$ which represents $\sim 70\%$ reduced domain shift w.r.t. realistic domain shifts (*i.e.* w.r.t. 1.5 to 2.0).

B.2.2 Subsidiary-Task Similarity Metric (TSM). γ_{TSM} determines how similar a subsidiary task is to the goal task. TSM is calculated using the basic linear evaluation protocol [62] in self-supervised literature. It illustrates the degree of compatibility between the two tasks. For computing γ_{TSM} , we train a linear classifier f_n on the features $h_{s,g}$ for subsidiary task dataset $\mathcal{D}_{s,n}$ extracted using a frozen source-pretrained ResNet-50 [18] backbone. For the sticker classification task, we randomly select 4 classes to keep the number of classes uniform for the different subsidiary task candidates illustrated in Fig. 1C and Fig. 7B in the main paper. We thus obtain the error for the different subsidiary task classifiers as $\hat{\epsilon}_{s,n}$ and the subsidiary-task similarity metric is computed as:

$$\gamma_{TSM}(\mathcal{D}_s, \mathcal{D}_{s,n}) = 1 - \min_{f_n} \hat{\epsilon}_{s,n}(h_{s,g}) \quad (16)$$

How to choose the threshold ζ_n ? Insight 3 introduced a threshold ζ_n for TSM to select pretext tasks suitable for subsidiary supervised DA. The task similarity of the subsidiary task is dependent on the goal task. For computing the threshold for TSM, we plot the γ_{TSM} for the candidate subsidiary tasks (Fig. 7B) and select the appropriate threshold ζ_n . Based on our observations in Fig. 7B of the main paper, we set ζ_n as 0.6.

Suitability criterion. Definition 1 in the main paper gives the overall suitability criterion for selecting the subsidiary task as:

$$\gamma_{DSM}(\mathcal{D}_s, \mathcal{D}_{s,n}) + \gamma_{TSM}(\mathcal{D}_s, \mathcal{D}_{s,n}) > \zeta \quad (17)$$

Therefore, we set the threshold ζ as a sum of ζ_d and ζ_n *i.e.* 1.1.



Fig. 9. The pseudo-OOS data $\mathcal{D}_s^{(od)}$ contains patch-shuffled versions of source data \mathcal{D}_s . Green circles only highlight the stickers and are not part of the samples.

C Implementation details

C.1 Sticker intervention

We define a sticker as a printed alphabet with a random color and random texture [8] within the alphabet. We scale the sticker randomly and paste it at a random location within a black image (all zeros) with the same size as goal task sample $x_s \in \mathbb{R}^{H \times W}$, yielding $x_n \in \mathbb{R}^{H \times W}$ (see Fig. 8A). The corresponding sticker-task labels y_n , along with x_n , form the sticker dataset \mathcal{D}_n . We also define a pixel-wise mask to perform mixup [81] only at the sticker pixels to avoid the effects of the black background on the rest of the goal task image.

Specifically, $m(u) = \mathbb{1}(x_n(u) \neq 0)$ where $u : [u_x, u_y]$ denotes the spatial index in an $H \times W$ lattice. As shown in Fig. 8B, a goal task sample x , *i.e.* either x_s , $x_s^{(od)}$ or x_t , and a sticker x_n are combined using mixup [81] as,

$$\mathcal{T}(x, x_n) = m \odot (\lambda x + (1 - \lambda)x_n) + (1 - m) \odot x \quad (18)$$

where λ denotes the mixup ratio, \odot represents element-wise multiplication and \mathcal{T} is the sticker intervention (as defined in Insight 4 of main paper).

C.1.1 Hyperparameters

- a) **Sticker shape** is decided by randomly selected alphabets.
- b) **Sticker size** is determined by randomly sampling the size ratio between sticker and goal task images from a uniform distribution over the range $[0.1, 0.4]$.
- c) **Sticker location** for pasting the sticker in the goal task image is sampled from a uniform distribution over the ranges $[1, H]$ and $[1, W]$. The sampled coordinates are rounded down to the nearest integer for pasting the sticker.
- d) **Number of sticker classes** determines the difficulty level of the subsidiary supervised DA problem.
- e) **Mixup ratio** determines the visibility of the sticker w.r.t. the goal task image. We use a constant mixup ratio of 0.4.

We provide ablations for these hyperparameters in Sec. D.2.

C.1.2 Usage. The intervention is applied in the same manner to both source samples x_s as well as target samples x_t , yielding sticker labels y_n for the sticker classifier. Mitsuzumi *et al.* [44] show that, beyond a certain grid size (4x4), shuffling the grid patches makes the domain unrecognizable. Inspired by this, we generate the pseudo-OOS dataset by randomly shuffling the grid patches with a grid size of (6x6) as shown in Fig. 9. The sticker intervention is also applied to the pseudo-OOS samples in order to emphasize the difference between source and pseudo-OOS samples even when stickers are present. However, for pseudo-OOS samples, the sticker label is treated as $y_s^{(od)}$, for the OOS node to act as an implicit domain discriminator, leading to improved source-target alignment.

Enabling source-free DA. The proposed sticker intervention can be used within source-free constraints. This is because, the alphabet font can be shared between source-side and target-side while the texture dataset [8] is open-source.

C.2 Experimental settings

Architecture details. We use a ResNet-50 [18] backbone for Office-Home, Office-31 and DomainNet, and ResNet-101 for VisDA, for a fair comparison with prior works. We employ the same network design as SHOT [37], *i.e.* replacing the classifier with a fully connected layer with batch norm [21] and another fully connected layer with weight normalization [61]. For the subsidiary classifier, we use the same architecture after ResLayer-3.

Optimization details. We employ multiple Adam optimizers during training to avoid loss weighting hyperparameters. Specifically, we use a distinct optimizer for each loss term. In each training iteration, we optimize only one of the losses (round robin method). Each optimizer uses a learning rate of 1e-3. Intuitively, each Adam optimizer’s moment parameters adaptively scale the associated gradients, eliminating the requirement for loss-scaling hyperparameter tuning. For source model training, following [37], we set the maximum number of epochs to 100 and 30 for Office-31 and Office-Home, whereas it is set to 10 and 15 for VisDA and DomainNet respectively. For adaptation, the maximum number of epochs is set to 15 for all datasets, following [37].

D Analysis

We provide more comparisons with prior state-of-the-art methods and report hyperparameter sensitivity analyses.

D.1 Extended comparisons and ablations

a) Single-Source DA for Office-31 and VisDA. Our approach outperforms source-free NRC [78] and SHOT++ [38] by 1.5% and 1.7% respectively on Office-31 (Table 8), and gives comparable performance to non-source-free works. On the larger and more challenging VisDA dataset, our approach surpasses NRC by 1.6% and SHOT++ by 1% (Table 8).

Table 8. Single-Source Domain Adaptation (SSDA) on Office-31 and VisDA benchmarks with mean and standard deviation over 5 runs. The last row indicates the variance over different sets of sticker shapes while others indicate variance over different random seeds. SF indicates *source-free* DA.

Method	SF	Office-31						VisDA	
		A→D	A→W	D→W	W→D	D→A	W→A	Avg	S → R
FAA [20]	✗	94.4	92.3	99.2	99.7	80.5	78.7	90.8	-
RFA [3]	✗	93.0	92.8	99.1	100.0	78.0	77.7	90.2	79.4
SCDA [35]	✗	95.4	95.3	99.0	100.0	77.2	75.9	90.5	-
DMRL [73]	✗	93.4±0.5	90.8±0.3	99.0±0.2	100.0±0.0	73.0±0.3	71.2±0.3	87.9	-
MCC [24]	✗	98.6±0.1	95.5±0.2	98.6±0.1	100.0±0.0	72.8±0.3	74.9±0.3	89.4	-
CAN [25]	✗	95.0±0.3	94.5±0.3	99.1±0.2	99.8±0.2	78.0±0.2	77.0±0.3	90.6	87.2
RWOT [76]	✗	94.5±0.2	95.1±0.2	99.5±0.2	100.0±0.0	77.5±0.1	77.9±0.3	90.8	-
FixBi [45]	✗	95.0±0.4	96.1±0.2	99.3±0.2	100.0±0.0	78.7±0.5	79.4±0.3	91.4	87.2
CDAN+RADA [23]	✗	96.1±0.4	96.2±0.4	99.3±0.1	100.0±0.0	77.5±0.1	77.4±0.3	91.1	76.3
SHOT [37]	✓	94.0	90.1	98.4	99.9	74.7	74.3	88.6	82.9
CPGA [56]	✓	94.4	94.1	98.4	99.8	76.0	76.6	89.9	84.1
HCL [19]	✓	90.8	91.3	98.2	100.0	72.7	72.7	87.6	83.5
VDM-DA [68]	✓	93.2	94.1	98.0	100.0	75.8	77.1	89.7	85.1
A ² Net [75]	✓	94.5	94.0	99.2	100.0	76.7	76.1	90.1	84.3
NRC [78]	✓	96.0	90.8	99.0	100.0	75.3	75.0	89.4	85.9
SHOT++ [38]	✓	94.3	90.4	98.7	99.9	76.2	75.8	89.2	87.3
3C-GAN [34]	✓	92.7±0.4	93.7±0.2	98.5±0.1	99.8±0.2	75.3±0.5	77.8±0.1	89.6	-
SFDA [27]	✓	92.2±0.2	91.1±0.3	98.2±0.3	99.5±0.2	71.0±0.2	71.2±0.2	87.2	-
<i>Ours (random seed)</i>	✓	95.6±0.2	94.6±0.2	99.2±0.1	99.8±0.2	77.0±0.3	77.7±0.3	90.7	88.2±0.4
<i>Ours (random sticker)</i>	✓	95.5±0.1	94.2±0.2	98.9±0.2	99.9±0.1	77.2±0.1	76.3±0.2	90.3	88.0±0.3

b) Multi-Source DA for Office-31. To analyze our performance on closed-set MSDA, we compare our approach with source-free and non-source-free prior arts in Table 12. Even without domain labels, our approach achieves *state-of-the-art* results on the Office-31 benchmark, even for the non-source-free setting.

c) Variance across random seeds. We highlight the significance of our results by reporting the mean and standard deviation of accuracy for 5 runs with different random seeds (2nd last row of Table 8) for SSDA. We observe low variance even w.r.t. prior non-source-free works.

d) Ablations for target adaptation. We present ablations on the goal task objectives for the target-side training (\mathcal{L}_{st} and \mathcal{L}_{div}) in Table 9. First, we compare the baseline *i.e.* source-trained model (#1) with the \mathcal{L}_{div} based DA model (#2). It is interesting to note that only using the diversity objective with subsidiary supervision improves SSDA and MSDA by 2.4% and 5.5% respectively over the baseline (#2 vs. #1), highlighting the relevance of diversity promotion.

The neighborhood clustering based self-training loss \mathcal{L}_{st} improves target clustering in the latent \mathcal{Z} space by bringing the backbone features $h(x)$ closer to their respective nearest neighbors. Using \mathcal{L}_{st} in conjunction with the subsidiary DA loss $\mathcal{L}_{t,n}$ enhances the goal task adaptation by 10.5% and 5.2% for SSDA and MSDA respectively, compared to not using \mathcal{L}_{st} (#4 vs. #2). We observe that employing both \mathcal{L}_{div} and \mathcal{L}_{st} further improves the performance by 3.8% and 1.9% for SSDA and MSDA respectively (#4 vs. #3), demonstrating that the two losses are complementary for goal task DA.

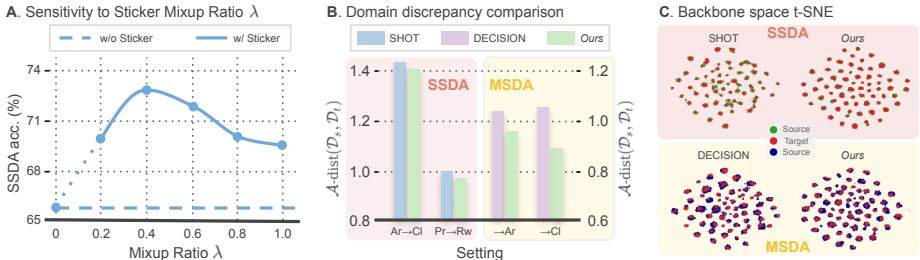


Fig. 10. **A.** Sensitivity to sticker mixup ratio λ for SSDA on Office-Home. **B.** \mathcal{A} -distance between source and target data on Office-Home. **C.** Backbone feature space t-SNE comparisons with SHOT [37] on Rw→Pr (SSDA), DECISION [2] on →Ar (MSDA) from Office-Home.

D.2 Hyperparameter sensitivity analysis

a) Sticker shape. We randomly selected 10 alphabets and used them consistently to report all the results in the main paper. However, to test the variance of our approach w.r.t. sticker shape, we report the mean and standard deviation over 5 runs of SSDA experiments on Office-31 (last row of Table 8), randomly sampling the 10 alphabets (*i.e.* sticker shapes) for each run. We observe a low standard deviation indicating low sensitivity to the sticker shapes.

b) Sticker size. We select this scale range based on empirical evidence (Table 10). We observe that adaptation performance suffers with sticker scale less than 0.1, since the sticker is hardly visible, making it difficult for the sticker classifier to receive meaningful supervision. The performance with larger sized stickers (more than 0.7) also drops as the sticker may occlude goal task content significantly.

c) Sticker location. We observe that our approach is only mildly sensitive to this hyperparameter (Table 10). We restrict the sticker location to regions far from the image centre and observe slightly lower accuracy. On the other hand, pasting the sticker near the image centre area further decreases performance as the sticker may occlude a larger part of the goal task content. Allowing the sticker to be pasted uniformly across the image yields the best performance.

d) Number of sticker classes. We perform a sensitivity analysis for the number of sticker categories $|\mathcal{C}_n|$ for MSDA on Office-Home (Fig. 11). We observe that performance improves with increasing number of classes upto 10 and reduces slightly for higher $|\mathcal{C}_n|$. Overall, we observe consistent gains over the baseline.

e) Mixup ratio λ . In Fig. 10A, we observe consistent gains over the baseline (mixup ratio $\lambda = 0$ *i.e.* sticker classifier and losses not used) for a wide range of λ values. The best performance is observed for $\lambda = 0.4$. Intuitively, higher mixup ratios imply very low sticker visibility while lower mixup ratios imply more occlusion of goal task content, both yielding slightly lower performance.

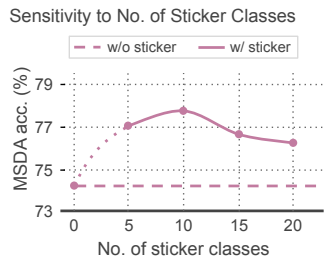


Fig. 11. Sensitivity to no. of sticker classes $|\mathcal{C}_n|$ for Office-Home MSDA.

Intuitively, higher mixup ratios imply very low sticker visibility while lower mixup ratios imply more occlusion of goal task content, both yielding slightly lower performance.

Table 9. Ablation study on Office-Home. SF, SSDA and MSDA indicate source-free, single-source DA and multi-source DA.

#	Target-side			SF	Office-Home	
	\mathcal{L}_{st}	\mathcal{L}_{div}	$\mathcal{L}_{t,n}$		SSDA	MSDA
1	\times	\times	\times	-	60.2	66.9
2	\times	\checkmark	\checkmark	\checkmark	62.6	72.4
3	\checkmark	\times	\checkmark	\checkmark	69.3	75.7
4	\checkmark	\checkmark	\checkmark	\checkmark	73.1	77.6

Table 10. Sensitivity analysis for sticker scale and location on the single-source DA (SSDA) benchmark of Office-Home dataset.

Sticker scale Acc.		Sticker location	
0.05 – 0.1	71.8	Central region	71.5
0.1 – 0.4	72.2	Except central region	72.0
0.4 – 0.7	73.1	Entire image	73.1
0.7 – 1.0	72.0		

D.3 Domain discrepancy analysis

In Fig. 10B, we report \mathcal{A} -distance as a measure of the domain discrepancy $d_{\mathcal{H}}(p_s, p_t)$ across different source-target pairings in the backbone feature space \mathcal{Z} for our approach and prior source-free state-of-the-art SSDA [37] and MSDA [2] works. A lower value for \mathcal{A} -distance indicates lower domain discrepancy. In comparison to prior works, our technique clearly achieves lower \mathcal{A} -distance between source and target for both settings. This implies that our backbone learns domain-agnostic features that are more generalized to the target domain. This corresponds to an increase in target performance and demonstrates that subsidiary supervised adaptation efficiently minimizes the latent space distribution shift, $d_{\mathcal{H}}(p_s, p_t)$, consistent with Insight 1 of the main paper.

D.4 Domain alignment analysis

In Fig. 10C, we present t-SNE [42] visualizations of backbone features learned by SHOT [37] and our approach for SSDA, and DECISION [2] and our approach for MSDA. As expected, all three approaches aid the formation of target clusters but source-target alignment for prior arts is weaker compared to our approach. We also observe that our method better preserves the source clusters (*green* in SSDA and *blue* in MSDA) while producing dense clusters for the target features (*red* in both settings) that are better aligned with the source clusters. This improved source-target alignment can be attributed to the OOS node in the sticker classifier, consistent with Insight 6 presented in the main paper.

Table 11. Training and inference time comparison w.r.t. NRC [78] and SHOT++ [38]. All timings are obtained using a single 1080Ti GPU.

Method	Training time (in sec), Ar→Cl				Inference time (in millisecond)	Office-Home	
	Source pretrain	Sticker pretrain	Target adapt	Total		SSDA Avg.	MSDA Avg.
NRC	282	-	1060	1342	1.9	72.2	74.7
SHOT++	306	-	10043	10349	1.9	73.0	75.7
<i>Ours</i>	282	643	284	1209	1.9	74.0	77.6

Table 12. Multi-Source DA (MSDA) comparisons on Office-31.

Method	SF	Office-31			
		→A	→W	→D	Avg.
PFSA [12]	✗	57.0	97.4	99.7	84.7
SimpAl [69]	✗	70.6	97.4	99.2	89.0
WAMDA [1]	✗	72.0	98.6	99.6	90.0
MIAN [49]	✗	76.2	98.4	99.2	91.3
MLAN [77]	✗	75.7	98.8	99.6	91.4
Source-combine	✗	65.2	94.6	98.4	86.1
SHOT[37]-Ens	✓	75.0	94.9	97.8	89.3
DECISION [2]	✓	75.4	98.4	99.6	91.1
CAiDA [10]	✓	75.8	98.9	99.8	91.6
Ours	✓	78.3	99.1	99.7	92.4

Table 13. Combining multiple subsidiary tasks (SSDA on Office-Home).

	SSDA
Baseline (B)	66.2
B + patch-loc	67.6
B + rotation	67.9
B + rotation + patch-loc	68.0
B + sticker-rot	69.0
B + sticker-clsf	69.7
B + sticker-rot + sticker-clsf	69.5

D.5 Efficiency analysis

We provide detailed training time comparisons of our work w.r.t. NRC [78] and SHOT++ [38] in Table 11. We make certain observations: **1)** We achieve superior target adaptation efficiency with the fastest training (4th column) and the best performance (last 2 columns). Note that we use same learning rate and scheduler as in NRC and SHOT++. **2)** Inference complexity (6th column) is same for all as we do not require the subsidiary classifier during inference.

D.6 Combining subsidiary tasks

Introducing multiple subsidiary tasks in the same framework brings up additional challenges like multi-task balancing. For instance, consider a combination of rotation (Rot) and patch-location (PL). From Fig. 1C in the main paper, Rot has high TSM while PL has high DSM. This does not imply that combining Rot and PL would yield a better overall TSM+DSM, and may rather have a detrimental impact. Thus, one should aim for a subsidiary task having both TSM and DSM greater than those of Rot and PL. Empirically, we do not find any conclusive result. In Table 13, we observe that while Rot+PL shows marginal gains, combining Sticker-rot and Sticker-clsf shows degraded performance.

Table 14. Comparisons w.r.t. pretext task based DA works.

Method	Pretext Task	High DSM+TSM	Additional regularization
SS-DA [17]	Rotation, Rot. Patch Jigsaw	✗	Adv. alignment, AdaBN
JiGen [5]	Jigsaw	✗	Augmentations
PAC [34]	Rotation	✗	Aug. consistency
<i>Ours</i>	Sticker	✓	None

Table 15. Comparisons w.r.t. prior source-free DA works.

Method	Key insights (differences)	Common
SHOT	Info-max. for implicit feature alignment	\mathcal{L}_{div}
SHOT++	Easy-hard target split for better adaptation	\mathcal{L}_{div}
CPGA	Contrastive prototypes for better pseudo-labels	\mathcal{L}_{st}
GSFDA	Local struct. clustering for better repr. learning	$\mathcal{L}_{st}, \mathcal{L}_{div}$
NRC	Cluster assumption for better pseudo-labels	$\mathcal{L}_{st}, \mathcal{L}_{div}$
A ² Net	Dual classifiers to find src-similar tgt samples and contrastive matching for category-wise alignm.	-
Ours	<ol style="list-style-type: none"> How and when subsidiary task is DA-assistive? Criteria for DA-assistive subsidiary tasks Process of sticker intervention 	$\mathcal{L}_{st}, \mathcal{L}_{div}$

D.7 Differences and relationships with prior-arts

These are discussed in Table 14 and 15. Our method is free from additional regularization unlike prior works (Table 14). While our key contributions are unique, the common losses are widely used (*e.g.* GSFDA, NRC in Table 15).

References

- Aggarwal, S., Kundu, J.N., Babu, R.V., Chakraborty, A.: WAMDA: Weighted alignment of sources for multi-source domain adaptation. In: BMVC (2020) 12, 25
- Ahmed, S.M., Raychaudhuri, D.S., Paul, S., Oymak, S., Roy-Chowdhury, A.K.: Unsupervised multi-source domain adaptation without access to source data. In: CVPR (2021) 11, 12, 23, 24, 25
- Awais, M., Zhou, F., Xu, H., Hong, L., Luo, P., Bae, S.H., Li, Z.: Adversarial robustness for unsupervised domain adaptation. In: ICCV (2021) 13, 22
- Ben-David, S., Blitzer, J., Crammer, K., Pereira, F.: Analysis of representations for domain adaptation. In: NeurIPS (2006) 1, 2, 5, 17
- Bousmalis, K., Silberman, N., Dohan, D., Erhan, D., Krishnan, D.: Unsupervised pixel-level domain adaptation with generative adversarial networks. In: CVPR (2017) 4
- Carlucci, F.M., D’Innocente, A., Bucci, S., Caputo, B., Tommasi, T.: Domain generalization by solving jigsaw puzzles. In: CVPR (2019) 1, 3, 4, 7, 8
- Chen, Y.H., Chen, W.Y., Chen, Y.T., Tsai, B.C., Frank Wang, Y.C., Sun, M.: No more discrimination: Cross city adaptation of road scene segmenters. In: ICCV (2017) 1
- Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., Vedaldi, A.: Describing textures in the wild. In: CVPR (2014) 20, 21
- Doersch, C., Gupta, A., Efros, A.A.: Unsupervised visual representation learning by context prediction. In: ICCV (2015) 4
- Dong, J., Fang, Z., Liu, A., Sun, G., Liu, T.: Confident anchor-induced multi-source free domain adaptation. In: NeurIPS (2021) 12, 25
- Feng, Z., Xu, C., Tao, D.: Self-supervised representation learning from multi-domain data. In: ICCV (2019) 4

12. Fu, Y., Zhang, M., Xu, X., Cao, Z., Ma, C., Ji, Y., Zuo, K., Lu, H.: Partial feature selection and alignment for multi-source domain adaptation. In: CVPR (2021) [25](#)
13. Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., Lempitsky, V.: Domain-adversarial training of neural networks. *The Journal of Machine Learning Research* **17**(1), 2096–2030 (2016) [1](#), [4](#), [5](#), [8](#), [14](#), [17](#)
14. Ghifary, M., Kleijn, W.B., Zhang, M., Balduzzi, D.: Domain generalization for object recognition with multi-task autoencoders. In: ICCV (2015) [4](#)
15. Ghifary, M., Kleijn, W.B., Zhang, M., Balduzzi, D., Li, W.: Deep reconstruction-classification networks for unsupervised domain adaptation. In: ECCV (2016) [4](#)
16. Gidaris, S., Singh, P., Komodakis, N.: Unsupervised representation learning by predicting image rotations. In: ICLR (2018) [1](#), [2](#), [4](#)
17. Gomes, R., Krause, A., Perona, P.: Discriminative clustering by regularized information maximization. In: NeurIPS (2010) [17](#)
18. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016) [11](#), [17](#), [19](#), [21](#)
19. Huang, J., Guan, D., Xiao, A., Lu, S.: Model adaptation: Historical contrastive learning for unsupervised domain adaptation without source data. In: NeurIPS (2021) [13](#), [22](#)
20. Huang, J., Guan, D., Xiao, A., Lu, S.: RDA: Robust domain adaptation via fourier adversarial attacking. In: ICCV (2021) [22](#)
21. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. In: ICML (2015) [11](#), [21](#)
22. Jiaolong, X., Liang, X., López, A.M.: Self-supervised domain adaptation for computer vision tasks. *IEEE Access* **7**, 156694–156706 (2019) [4](#), [8](#)
23. Jin, X., Lan, C., Zeng, W., Chen, Z.: Re-energizing domain discriminator with sample relabeling for adversarial domain adaptation. In: ICCV (2021) [13](#), [22](#)
24. Jin, Y., Wang, X., Long, M., Wang, J.: Minimum class confusion for versatile domain adaptation. In: ECCV (2020) [22](#)
25. Kang, G., Jiang, L., Yang, Y., Hauptmann, A.G.: Contrastive adaptation network for unsupervised domain adaptation. In: CVPR (2019) [13](#), [22](#)
26. Kim, D., Saito, K., Oh, T.H., Plummer, B.A., Sclaroff, S., Saenko, K.: CDS: Cross-domain self-supervised pre-training. In: ICCV (2021) [4](#)
27. Kim, Y., Cho, D., Han, K., Panda, P., Hong, S.: Domain adaptation without source data. *IEEE Transactions on Artificial Intelligence* **2**(6), 508–518 (2021) [22](#)
28. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014) [11](#)
29. Kundu, J.N., Kulkarni, A., Bhambri, S., Mehta, D., Kulkarni, S., Jampani, V., Babu, R.V.: Balancing discriminability and transferability for source-free domain adaptation. In: ICML (2022) [3](#)
30. Kundu, J.N., Kulkarni, A., Singh, A., Jampani, V., Babu, R.V.: Generalize then adapt: Source-free domain adaptive semantic segmentation. In: ICCV (2021) [4](#)
31. Kundu, J.N., Venkat, N., M V, R., Babu, R.V.: Universal source-free domain adaptation. In: CVPR (2020) [4](#), [9](#)
32. Kundu, J.N., Venkat, N., Revanur, A., V, R.M., Babu, R.V.: Towards inheritable models for open-set domain adaptation. In: CVPR (2020) [4](#)
33. Larsson, G., Maire, M., Shakhnarovich, G.: Colorization as a proxy task for visual understanding. In: CVPR (2017) [3](#), [4](#), [8](#)
34. Li, R., Jiao, Q., Cao, W., Wong, H.S., Wu, S.: Model adaptation: Unsupervised domain adaptation without source data. In: CVPR (2020) [3](#), [4](#), [9](#), [22](#)
35. Li, S., Xie, M., Lv, F., Liu, C.H., Liang, J., Qin, C., Li, W.: Semantic concentration for domain adaptation. In: ICCV (2021) [11](#), [12](#), [22](#)

36. Li, Y., Yuan, L., Chen, Y., Wang, P., Vasconcelos, N.: Dynamic transfer for multi-source domain adaptation. In: CVPR (2021) [12](#)
37. Liang, J., Hu, D., Feng, J.: Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In: ICML (2020) [3](#), [4](#), [10](#), [11](#), [12](#), [13](#), [17](#), [18](#), [21](#), [22](#), [23](#), [24](#), [25](#)
38. Liang, J., Hu, D., Wang, Y., He, R., Feng, J.: Source data-absent unsupervised domain adaptation through hypothesis transfer and labeling transfer. IEEE Transactions on Pattern Analysis and Machine Intelligence (2021) [4](#), [11](#), [12](#), [13](#), [21](#), [22](#), [24](#), [25](#)
39. Liu, Y., Zhang, W., Wang, J.: Source-free domain adaptation for semantic segmentation. In: CVPR (2021) [4](#)
40. Long, M., Cao, Y., Wang, J., Jordan, M.: Learning transferable features with deep adaptation networks. In: ICML (2015) [4](#)
41. Long, M., Cao, Z., Wang, J., Jordan, M.I.: Conditional adversarial domain adaptation. In: NeurIPS (2017) [1](#), [14](#)
42. Maaten, L.v.d., Hinton, G.: Visualizing data using t-sne. Journal of machine learning research **9**(Nov), 2579–2605 (2008) [24](#)
43. Mishra, S., Saenko, K., Saligrama, V.: Surprisingly simple semi-supervised domain adaptation with pretraining and consistency. In: BMVC (2021) [1](#), [2](#), [3](#), [8](#), [12](#)
44. Mitsuzumi, Y., Irie, G., Ikami, D., Shibata, T.: Generalized domain adaptation. In: CVPR (2021) [7](#), [10](#), [21](#)
45. Na, J., Jung, H., Chang, H.J., Hwang, W.: FixBi: Bridging domain spaces for unsupervised domain adaptation. In: CVPR (2021) [11](#), [13](#), [14](#), [22](#)
46. Nguyen, V.A., Nguyen, T., Le, T., Tran, Q.H., Phung, D.: STEM: An approach to multi-source domain adaptation with guarantees. In: ICCV (2021) [12](#)
47. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: ECCV (2016) [1](#), [2](#), [4](#)
48. Oord, A.v.d., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748 (2018) [4](#)
49. Park, G.Y., Lee, S.W.: Information-theoretic regularization for multi-source domain adaptation. In: ICCV (2021) [25](#)
50. Parkhi, O.M., Vedaldi, A., Zisserman, A., Jawahar, C.V.: Cats and dogs. In: CVPR (2012) [2](#)
51. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: NeurIPS (2019) [17](#)
52. Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., Efros, A.A.: Context encoders: Feature learning by inpainting. In: CVPR (2016) [3](#), [4](#), [8](#)
53. Peng, X., Bai, Q., Xia, X., Huang, Z., Saenko, K., Wang, B.: Moment matching for multi-source domain adaptation. In: ICCV (2019) [11](#)
54. Peng, X., Usman, B., Kaushik, N., Hoffman, J., Wang, D., Saenko, K.: VisDA: The visual domain adaptation challenge. arXiv preprint arXiv:1710.06924 (2017) [11](#)
55. Prabhu, V., Khare, S., Kartik, D., Hoffman, J.: SENTRY: Selective entropy optimization via committee consistency for unsupervised domain adaptation. In: ICCV (2021) [11](#)
56. Qiu, Z., Zhang, Y., Lin, H., Niu, S., Liu, Y., Du, Q., Tan, M.: Source-free domain adaptation via avatar prototype generation and adaptation. In: IJCAI (2021) [11](#), [13](#), [22](#)

57. Ren, Z., Lee, Y.J.: Cross-domain self-supervised multi-task feature learning using synthetic imagery. In: CVPR (2018) [4](#)
58. Saenko, K., Kulis, B., Fritz, M., Darrell, T.: Adapting visual category models to new domains. In: ECCV (2010) [11](#)
59. Saito, K., Kim, D., Sclaroff, S., Saenko, K.: Universal domain adaptation through self supervision. In: NeurIPS (2020) [4](#)
60. Saito, K., Yamamoto, S., Ushiku, Y., Harada, T.: Open set domain adaptation by backpropagation. In: ECCV (2018) [4](#)
61. Salimans, T., Kingma, D.P.: Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In: NeurIPS (2016) [11](#), [21](#)
62. Salman, H., Ilyas, A., Engstrom, L., Kapoor, A., Madry, A.: Do adversarially robust imagenet models transfer better? In: NeurIPS (2020) [5](#), [19](#)
63. Scalbert, M., Vakalopoulou, M., Couzini'e-Devy, F.: Multi-source domain adaptation via supervised contrastive learning and confident consistency regularization. In: BMVC (2021) [12](#), [14](#)
64. Sivaprasad, P.T., Fleuret, F.: Uncertainty reduction for model adaptation in semantic segmentation. In: CVPR (2021) [4](#)
65. Sun, B., Saenko, K.: Deep CORAL: Correlation alignment for deep domain adaptation. In: ECCV (2016) [4](#), [8](#)
66. Sun, Y., Tzeng, E., Darrell, T., Efros, A.A.: Unsupervised domain adaptation through self-supervision. arXiv preprint arXiv:1909.11825 (2019) [1](#), [3](#), [4](#), [8](#), [12](#)
67. Tang, H., Chen, K., Jia, K.: Unsupervised domain adaptation via structurally regularized deep clustering. In: CVPR (2020) [14](#)
68. Tian, J., Zhang, J., Li, W., Xu, D.: VDM-DA: Virtual domain modeling for source data-free domain adaptation. IEEE Transactions on Circuits and Systems for Video Technology (2021) [13](#), [22](#)
69. Venkat, N., Kundu, J.N., Singh, D.K., Revanur, A., Babu, R.V.: Your classifier can secretly suffice multi-source domain adaptation. In: NeurIPS (2020) [12](#), [18](#), [25](#)
70. Venkateswara, H., Eusebio, J., Chakraborty, S., Panchanathan, S.: Deep hashing network for unsupervised domain adaptation. In: CVPR (2017) [11](#)
71. Wallace, B., Hariharan, B.: Extending and analyzing self-supervised learning across domains. In: ECCV (2020) [2](#)
72. Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., Summers, R.: ChestX-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In: CVPR (2017) [2](#)
73. Wu, Y., Inkpen, D., El-Roby, A.: Dual mixup regularized learning for adversarial domain adaptation. In: ECCV (2020) [22](#)
74. Wu, Z., Xiong, Y., Yu, S.X., Lin, D.: Unsupervised feature learning via non-parametric instance discrimination. In: CVPR (2018) [4](#)
75. Xia, H., Zhao, H., Ding, Z.: Adaptive adversarial network for source-free domain adaptation. In: ICCV (2021) [11](#), [13](#), [22](#)
76. Xu, R., Liu, P., Wang, L., Chen, C., Wang, J.: Reliable weighted optimal transport for unsupervised domain adaptation. In: CVPR (2020) [22](#)
77. Xu, Y., Kan, M., Shan, S., Chen, X.: Mutual learning of joint and separate domain alignments for multi-source domain adaptation. In: WACV (2022) [25](#)
78. Yang, S., Wang, Y., van de Weijer, J., Herranz, L., Jui, S.: Exploiting the intrinsic neighborhood structure for source-free domain adaptation. In: NeurIPS (2021) [4](#), [11](#), [12](#), [13](#), [16](#), [21](#), [22](#), [24](#), [25](#)
79. Yang, S., Wang, Y., van de Weijer, J., Herranz, L., Jui, S.: Generalized source-free domain adaptation. In: ICCV (2021) [11](#)

80. You, K., Long, M., Cao, Z., Wang, J., Jordan, M.I.: Universal domain adaptation. In: CVPR (2019) 4
81. Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. In: ICLR (2018) 20
82. Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: ECCV (2016) 4, 8
83. Zhang, R., Isola, P., Efros, A.A.: Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In: CVPR (2017) 4
84. Zhao, H., Combes, R.T.D., Zhang, K., Gordon, G.: On learning invariant representations for domain adaptation. In: ICML (2019) 2, 5