
Gradual Domain Adaptation without Indexed Intermediate Domains

Hong-You Chen

The Ohio State University, USA
chen.9301@osu.edu

Wei-Lun Chao

The Ohio State University, USA
chao.209@osu.edu

Abstract

The effectiveness of unsupervised domain adaptation degrades when there is a large discrepancy between the source and target domains. Gradual domain adaption (GDA) is one promising way to mitigate such an issue, by leveraging additional unlabeled data that gradually shift from the source to the target. Through sequentially adapting the model along the “indexed” intermediate domains, GDA substantially improves the overall adaptation performance. In practice, however, the extra unlabeled data may not be separated into intermediate domains and indexed properly, limiting the applicability of GDA. In this paper, we investigate how to discover the sequence of intermediate domains when it is not already available. Concretely, we propose a coarse-to-fine framework, which starts with a *coarse* domain discovery step via progressive domain discriminator training. This coarse domain sequence then undergoes a *fine* indexing step via a novel cycle-consistency loss, which encourages the next intermediate domain to preserve sufficient discriminative knowledge of the current intermediate domain. The resulting domain sequence can then be used by a GDA algorithm. On benchmark data sets of GDA, we show that our approach, which we name **Intermediate D**omain **L**abeler (**IDOL**), can lead to comparable or even better adaptation performance compared to the pre-defined domain sequence, making GDA more applicable and robust to the quality of domain sequences. Codes are available at <https://github.com/hongyouc/IDOL>.

1 Introduction

The distributions of real-world data change dynamically due to many factors like time, locations, environments, etc. Such a fact poses a great challenge to machine-learned models, which implicitly assume that the test data distribution is covered by the training data distribution. To resolve this generalization problem, unsupervised domain adaption (UDA), which aims to adapt a learned model to the test domain given its unlabeled data [13, 15], has been an active sub-field in machine learning.

Typically, UDA assumes that the “source” domain, in which the model is trained, and the “target” domain, in which the model is deployed, are discrepant but sufficiently related. Concretely, Ben-David et al. [1], Zhao et al. [75] show that the generalization error of UDA is bounded by the discrepancy of the marginal or conditional distributions between domains. Namely, the effectiveness of UDA may degrade along with the increase in domain discrepancy. Take one popular algorithm, self-training [38, 48, 49], for example. Self-training adapts the source model by progressively labeling the unlabeled target data (i.e., pseudo-labels) and using them to fine-tune the model [2, 28, 30, 31, 41, 44, 77, 78]. Self-training works if the pseudo-labels are accurate (as it essentially becomes supervised learning), but is vulnerable if they are not, which occurs when there exists a large domain gap [35].

To address this issue, several recent works investigate *gradual domain adaption (GDA)* [12, 25, 35, 71], in which beyond the source and target domains, the model can access additional unlabeled data from the intermediate domains that shift gradually from the source to the target. By adapting

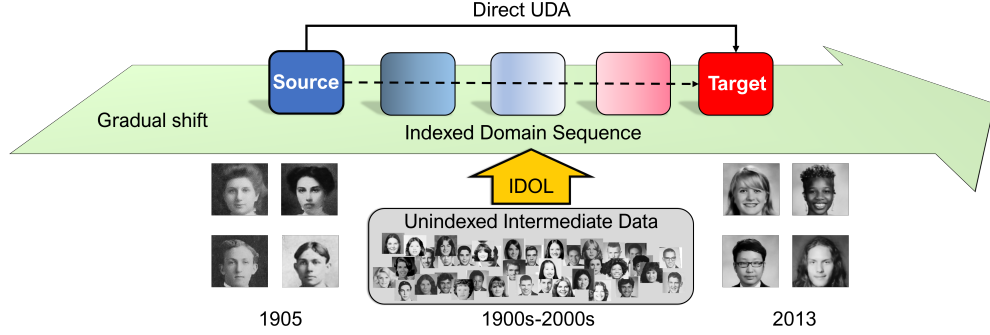


Figure 1: **Gradual domain adaption (GDA) without indexed intermediate domains.** In this setting, one is provided with labeled source, unlabeled target, and additional unlabeled intermediate data that have not been grouped and indexed into a domain sequence. Our approach **intermediate domain labeler (IDOL)** can successfully discover the domain sequence, which can then be leveraged by a GDA algorithm to achieve a higher target accuracy than direct unsupervised domain adaptation (UDA). Images are from the Portraits dataset [14].

the model along the sequence of intermediate domains — i.e., from the ones close to the source to the ones close to the target — the large domain gap between source and target is chipped away by multiple sub-adaptation problems (between consecutive intermediate domains) whose domain gaps are smaller. Namely, every time the model moves a step closer to the target, instead of taking a huge jump that can significantly decrease the performance (see Figure 1 for an illustration). GDA makes sense, as real-world data change gradually more often than abruptly [11, 60, 66]. The recent work by Kumar et al. [35] further demonstrates the strength of GDA both empirically and theoretically.

One potential drawback of GDA is the need for a well-defined sequence of intermediate domains. That is, prior to adaptation, the additional unlabeled data must be grouped into multiple domains and indexed with intermediate domain labels that reflect the underlying data distribution shift from the source to the target. Such information, however, may not be available directly. Existing methods usually leverage side information like time tags [25, 35, 71] to define the sequence, which may be sub-optimal. In some applications, even the side information may not be accessible (e.g., due to privacy concerns), greatly limiting the applicability of GDA.

In this paper, we therefore study GDA in the extreme case — *the additional unlabeled data are neither grouped nor indexed*. Specifically, we investigate how to discover the “domain sequence” from data, such that it can be used to drive GDA. We propose a two-stage coarse-to-fine framework named **Intermediate D**omain **L**abeler (**IDOL**). In the first stage, IDOL labels each intermediate data instance with a *coarse* score that reflects how close it is to the source or target. We study several methods that estimate the distance from a data instance to a domain. We find that a progressively trained domain discriminator — which starts with source vs. target data but gradually adds data close to either of them into training — performs the best, as it better captures the underlying data manifold.

In the second stage, IDOL then builds upon the *coarse* scores to group data into domains by further considering the discriminative (e.g., classification) knowledge the model aims to preserve from the source to the target. Since the additional unlabeled data are fully unlabeled, we take a greedy approach to identify the next intermediate domain (i.e., a group of data) that can best preserve the discriminative knowledge in the current intermediate domain. Concretely, we employ self-training [38] along the domain sequence discovered so far to provide pseudo-labels for the current domain, and propose a novel *cycle-consistency loss* to discover the next domain, such that the model adapted to the next domain can be “adapted back” to the current domain and predict the same pseudo-labels. The output of IDOL is a sequence of intermediate domains that can be used by any GDA algorithms [25, 35, 71].

We validated IDOL on two data sets studied in [35], including Rotated MNIST [36] and Portraits over years [14]. IDOL can successfully discover the domain sequence that leads to comparable GDA performance to using the pre-defined sequence (i.e., by side information). More importantly, IDOL is compatible with pre-defined sequences — by treating them as the coarse sequences — to further improve upon them. We also investigate IDOL in scenarios where the additional unlabeled data may contain outliers and demonstrate IDOL’s effectiveness even in such challenging cases. To our knowledge, our work is the first to tackle GDA without grouped and indexed intermediate domains. The success of IDOL opens up broader application scenarios that GDA can contribute to.

2 Background: gradual domain adaptation (GDA)

2.1 Setup of UDA and GDA

We consider an unsupervised domain adaptation (UDA) problem, in which a classifier is provided with labeled data from the source domain $\mathcal{S} = \{(\mathbf{x}_i^{\mathcal{S}}, y_i^{\mathcal{S}})\}_{i=1}^{|\mathcal{S}|}$ and unlabeled data from the target domain $\mathcal{T} = \{\mathbf{x}_i^{\mathcal{T}}\}_{i=1}^{|\mathcal{T}|}$. Both \mathcal{S} and \mathcal{T} are sampled IID from the underlying joint distributions of the source $P_{\mathcal{S}}$ and target $P_{\mathcal{T}}$, respectively. A source model $\theta_{\mathcal{S}}$ is typically produced by training on \mathcal{S} directly. The goal of UDA is to learn a target model $\theta_{\mathcal{T}}$ such that it can perform well on the target, measured by a loss $\mathcal{L}(\theta_{\mathcal{T}}, P_{\mathcal{T}})$. We focus on the standard UDA setup that assumes no label shifts [6].

Gradual domain adaptation (GDA), on top of UDA, assumes the existence of a sequence of domains $P_0, P_1, \dots, P_{M-1}, P_M$, where P_0 and P_M are the source domain and target domain, respectively, along with $M - 1$ intermediate domains P_1, \dots, P_{M-1} . The data distributions of domains are assumed to change gradually from the source to the target along the sequence. One can access unlabeled data $\mathcal{U}_m = \{\mathbf{x}_i^{\mathcal{U}_m}\}_{i=1}^{|\mathcal{U}_m|}$ from each of the intermediate domain, forming a sequence of unlabeled data $\mathcal{U}_1, \dots, \mathcal{U}_{M-1}$. We denote the union of them by $\mathcal{U} = \{\mathbf{x}_i^{\mathcal{U}}\}_{i=1}^{|\mathcal{U}|}$. For brevity, we assume that each \mathcal{U}_m is disjoint from the others but the size $|\mathcal{U}_m|$ is the same $\forall m \in \{1, \dots, M-1\}$.

2.2 GDA with pre-defined domain sequences using self-training [35]

We summarize the theoretical analysis of GDA using self-training by Kumar et al. [35] as follows.

Self-training for UDA. Self-training [38] takes a pre-trained (source) model θ and the target unlabeled data \mathcal{T} as input. Denote by $\mathbf{z} = f(\mathbf{x}; \theta)$ the predicted logits \mathbf{z} , self-training first applies $f(\cdot; \theta)$ to target data points $\mathbf{x} \in \mathcal{T}$, sharpens the predictions by argmax into pseudo-labels, and then updates the current model θ by minimizing the loss ℓ (e.g., cross entropy) w.r.t the pseudo-labels,

$$\theta_{\mathcal{T}} = \text{ST}(\theta, \mathcal{T}) = \arg \min_{\theta' \in \Theta} \frac{1}{|\mathcal{T}|} \sum_{i=1}^{|\mathcal{T}|} \ell(f(\mathbf{x}_i; \theta'), \text{sharpen}(f(\mathbf{x}_i; \theta))), \quad (1)$$

where $\theta_{\mathcal{T}}$ denotes the resulting target model and $\text{ST}(\theta, \mathcal{T})$ denotes the self-training process. A baseline use of self-training for UDA is to set θ as the source model $\theta_{\mathcal{S}}$. However, when the domain gap between \mathcal{S} and \mathcal{T} is large, $\theta_{\mathcal{S}}$ cannot produce accurate pseudo-labels for effective self-training.

Gradual self-training for GDA. In the GDA setup, self-training is performed along the domain sequence, from the current model θ_m (updated using \mathcal{U}_m already) to the next one θ_{m+1} using \mathcal{U}_{m+1}

$$\theta_{m+1} = \text{ST}(\theta_m, \mathcal{U}_{m+1}), \quad (2)$$

starting from the source model $\theta_0 = \theta_{\mathcal{S}}$. By doing so, the model gradually adapts from \mathcal{S} to \mathcal{T} through the sequence $\mathcal{U}_1, \dots, \mathcal{U}_M$ to produce the final target model $\theta_{\mathcal{T}}$.

$$\theta_{\mathcal{T}} = \theta_M = \text{ST}(\theta_{\mathcal{S}}, (\mathcal{U}_1, \dots, \mathcal{U}_M)). \quad (3)$$

As the domain gap between consecutive domains is smaller than between \mathcal{S} and \mathcal{T} , the pseudo-labels will be more accurate in each step. As a result, self-training in GDA can lead to a lower target error than in UDA. We note that the domain sequences in [35] are pre-defined using side information.

Theoretical analysis. Kumar et al. [35] make three assumptions for theoretical analysis:

- Separation: for every domain P_m , the data are separable such that there exists an R -bounded (linear) classifier $\theta_m \in \Theta$ with $\|\theta_m\|_2 \leq R$ that achieves a low loss of $\mathcal{L}(\theta_m, P_m)$.
- Gradual shift: for the domain sequence P_0, \dots, P_m with no label shifts, the maximum per-class Wasserstein-infinity distance $\rho(P_m, P_{m+1})$ between consecutive domains should be $\leq \rho < \frac{1}{R}$.
- Bounded data: finite samples X from each domain P_m are bounded, i.e., $E_{X \sim P_m}[\|X\|_2^2] \leq B^2$.

Under these assumptions of GDA, if the source model θ_0 has a low loss $\mathcal{L}(\theta_0, P_0)$ on P_0 , then with a probability δ , the resulting target model θ_M through gradual self-training has

$$\mathcal{L}(\theta_M, P_M) \leq \beta^{M+1} \left(\mathcal{L}(\theta_0, P_0) + \frac{4BR + \sqrt{2 \log 2M/\delta}}{\sqrt{|\mathcal{U}|}} \right), \quad \text{where } \beta = \frac{2}{1 - \rho R}. \quad (4)$$

That is, the target error is controlled by the intermediate domain shift ρ . If we can sample infinite data (i.e., $|\mathcal{U}| \rightarrow \infty$) and the source classifier is perfect (i.e., $\mathcal{L}(\theta_0, P_0) = 0$), with a small distance $\rho(P_m, P_{m+1}) \leq \rho$ between consecutive domains, gradual self-training achieves a zero target error.

3 GDA without indexed intermediate domains

3.1 Setup and motivation

In this paper, we study the case of GDA where the additional unlabeled data \mathcal{U} are not readily divided into domain sequences $\mathcal{U}_1, \dots, \mathcal{U}_{M-1}$. In other words, we are provided with a single, gigantic unlabeled set that has a wide range of support from the source to the target domain. One naive way to leverage it is to adapt from \mathcal{S} to \mathcal{U}^1 and then to \mathcal{T} , which however leads to a much larger target error than applying self-training along the properly indexed intermediate domains [35]. We attribute this to the large $\rho(\mathcal{S}, \mathcal{U})$ or $\rho(\mathcal{U}, \mathcal{T})$ according to Equation 4. To take advantage of \mathcal{U} , we must separate it into intermediate domains such that $\rho(P_m, P_{m+1})$ is small enough for every $m \in \{0, \dots, M-1\}$.

3.2 Overview of our approach: intermediate domain labeler (IDOL)

Given the source data \mathcal{S} , the target data $\mathcal{U}_M = \mathcal{T}$, and the unlabeled intermediate data $\mathcal{U} = \{x_i^{\mathcal{U}}\}_{i=1}^{|\mathcal{U}|}$, we propose the intermediate domain labeler (IDOL), whose goal is to sort intermediate data instances in \mathcal{U} and chunk them into $M-1$ domains $\mathcal{U}_1, \dots, \mathcal{U}_{M-1}$ (with equal sizes) for GDA to succeed.

Taking gradual self-training in Equation 3 as an example, IDOL aims to solve the following problem

$$\begin{aligned} \min \mathcal{L}(\theta_{\mathcal{T}}, P_{\mathcal{T}}), \\ \text{s.t. } \theta_{\mathcal{T}} = \text{ST}(\theta_{\mathcal{S}}, (\mathcal{U}_1, \dots, \mathcal{U}_M)) \quad \text{and} \quad (\mathcal{U}_1, \dots, \mathcal{U}_{M-1}) = \text{IDOL}(\mathcal{S}, \mathcal{T}, \mathcal{U}; M-1), \end{aligned} \quad (5)$$

where the target model $\theta_{\mathcal{T}}$ is produced by applying gradual self-training to the domain sequence discovered by IDOL. IDOL accepts the number of intermediate domains $M-1$ as a hyper-parameter.

Solving Equation 5 is hard, as evaluating any output sequence needs running through the entire gradual self-training process, not to mention that we have no labeled target data to estimate $\mathcal{L}(\theta_{\mathcal{T}}, P_{\mathcal{T}})$. In this paper, we propose to solve Equation 5 approximately via a coarse-to-fine procedure.

The coarse stage. We aim to give each $x_i^{\mathcal{U}} \in \mathcal{U}$ a score q_i , such that it tells x 's position in between the source and target. Specifically, a higher/lower q_i indicates that $x_i^{\mathcal{U}}$ is closer to the source/target domain. With these scores, we can already obtain a coarse domain sequence by *sorting them in the descending order and dividing them into $M-1$ chunks*.

The fine stage. Instead of creating the domain sequence $\mathcal{U}_1, \dots, \mathcal{U}_{M-1}$ by looking at the score of each individual $x_i^{\mathcal{U}} \in \mathcal{U}$, we further consider how data grouped into the same domain can collectively preserve the discriminative knowledge from the previous domains — after all, the goal of UDA or GDA is to pass the discriminative knowledge from the labeled source domain to the target domain. To this end, we propose a novel cycle-consistency loss to refine the coarse scores progressively.

3.3 The coarse stage: assigning domain scores

In this stage, we assign each $x_i^{\mathcal{U}} \in \mathcal{U}$ a score q_i , such that higher/lower q_i indicates that $x_i^{\mathcal{U}}$ is closer to the source/target domain. We discuss some design choices. For brevity, we omit the superscript \mathcal{U} .

Confidence scores by the classifier $f(x; \theta)$. We first investigate scoring each intermediate data instance by the confidence score of the source model $\theta_{\mathcal{S}}$, i.e., $q_i = \max f(x_i; \theta_{\mathcal{S}})$ across classes. This is inspired by outlier detection [42] and the common practice of self-training, which selects only data with sufficiently large confidence scores to fine-tune upon [35, 38]. Here we employ an advanced version, which is to use the latest θ (i.e., the model adapted to the current intermediate domain) to select the next intermediate domain. That is, every time we select the highest confidence instances from the *remaining* ones in \mathcal{U} to adapt θ . Overall, we can give data selected in the m -th round ($m \in \{1, \dots, M-1\}$) a score $\frac{M-1-m}{M-2} \in [0, 1]$. One drawback of this method is the blindness to the target domain \mathcal{T} . Specifically, we find that using confidence scores tends to select easily classified examples, which do not necessarily follow the gradual shift from the source to the target.

Manifold distance with the source features. To model the flow from the source to the target, we argue that it is crucial to consider both sides as references. We thus extract features of \mathcal{S} , \mathcal{T} , and \mathcal{U} using the source model $\theta_{\mathcal{S}}$ and apply a manifold learning algorithm (e.g., UMAP [50]) to discover the

¹In [35], the authors investigated sampling $\mathcal{U}_1, \dots, \mathcal{U}_{M-1}$ from \mathcal{U} uniformly at random with replacement.

underlying manifold. Denote by $\gamma(\mathbf{x})$ the dimension-reduced features of \mathbf{x} after UMAP, we compute the ratio of its distance to the nearest point in the \mathcal{T} and \mathcal{S} as the score $q_i = \frac{\min_{\mathbf{x}' \in \mathcal{T}} \|\gamma(\mathbf{x}_i) - \gamma(\mathbf{x}')\|_2}{\min_{\mathbf{x}' \in \mathcal{S}} \|\gamma(\mathbf{x}_i) - \gamma(\mathbf{x}')\|_2}$.

Domain discriminator. We investigate another idea to emphasize the roles of the source and target, which is to train a domain discriminator [13]. Concretely, we construct a binary classifier $g(\cdot; \phi)$ using deep neural networks, which is trained to separate the source data \mathcal{S} (class: 1) and the target data \mathcal{T} (class: 0). We use a binary-cross entropy loss to optimize the learnable parameter ϕ

$$\mathcal{L}(\phi) = -\frac{1}{|\mathcal{S}|} \sum_{\mathbf{x}^S \in \mathcal{S}} \log(\sigma(g(\mathbf{x}^S; \phi))) - \frac{1}{|\mathcal{T}|} \sum_{\mathbf{x}^T \in \mathcal{T}} \log(1 - \sigma(g(\mathbf{x}^T; \phi))), \quad (6)$$

where σ is the sigmoid function. We can then assign a score to $\mathbf{x}_i \in \mathcal{U}$ by $q_i = g(\mathbf{x}_i; \phi)$.

Progressive training for the domain discriminator. The domain discriminator $g(\cdot; \phi)$, compared to the manifold distance, better contrasts the source and target domains but does not leverage any of the data $\mathbf{x} \in \mathcal{U}$. That is, it might give a faithful score to an example \mathbf{x} that is very close to \mathcal{S} or \mathcal{T} , but for other examples that are far away and hence out-of-domain from both ends, the score by $g(\cdot; \phi)$ becomes less reliable, not to mention that neural networks are usually poorly calibrated for these examples [21, 42]. We therefore propose to progressively augment the source and target data with $\mathbf{x} \in \mathcal{U}$ that has either a fairly high or low $g(\mathbf{x}; \phi)$, and fine-tune the domain discriminator $g(\mathbf{x}; \phi)$. We perform this process for K rounds, and every time include $\frac{|\mathcal{U}|}{2K}$ new examples into the source and target. By doing so, the domain discriminator $g(\mathbf{x}; \phi)$ is updated with data from \mathcal{U} and thus can provide more accurate scores to distinguish the remaining examples into the source or target side. More specifically, let N denote $|\mathcal{U}|$, we repeat the following steps for k from 1 to K :

1. Train $g(\cdot, \phi)$ using \mathcal{S} and \mathcal{T} , based on the loss in Equation 6.
2. Predict $\hat{q}_i = g(\mathbf{x}_i, \phi), \forall \mathbf{x}_i \in \mathcal{U}$.
3. Rank all \hat{q}_i in the descending order.
4. The data with the $\frac{N}{2K}$ largest \hat{q}_i scores form a new intermediate domain for the source side (their q_i is set to $\frac{2K-k}{2K}$). We remove them from \mathcal{U} and add them into \mathcal{S} .
5. The data with the $\frac{N}{2K}$ smallest \hat{q}_i scores form a new intermediate domain for the target side (their q_i is set to $\frac{k}{2K}$). We remove them from \mathcal{U} and add them into \mathcal{T} .

Overall, we give a data instance $\mathbf{x} \in \mathcal{U}$ selected in the k -th round ($k \in \{1, \dots, K\}$) a score $\frac{2K-k}{2K}$ if it is added into the source side, or $\frac{k}{2K}$ if it is added into the target side.

3.4 The fine stage: cycle-consistency for refinement

The domain scores developed in subsection 3.3 can already give each individual example $\mathbf{x} \in \mathcal{U}$ a rough position of which intermediate domain it belongs to. However, for GDA to succeed, we must consider intermediate data in a collective manner. That is, each intermediate domain should preserve sufficient discriminative (e.g., classification) knowledge from the previous one, such that after M rounds of adaptation through GDA, the resulting $\theta_{\mathcal{T}}$ in Equation 5 will be able to perform well on \mathcal{T} . This is reminiscent of some recent claims in UDA [46, 75]: matching only the marginal distributions between domains (i.e., blind to the class labels) may lead to sub-optimal adaptation results.

But how can we measure the amount of discriminative knowledge preserved by an intermediate domain (or in the target domain) if we do not have its labeled data? To seek for the answer, we revisit Equation 4 and consider the case $M = 1$. If the two domains P_0 and P_1 are sufficiently close and the initial model θ_0 is well-trained, the resulting model θ_1 after self-training will perform well. Reversely, we can treat such a well-performing θ_1 as the well-trained initial model, and apply self-training again — this time from P_1 back to P_0 . The resulting model θ'_0 , according to the same rationale, should perform well on P_0 . In other words, θ'_0 and θ_0 should predict similarly on data sampled from P_0 . The **similarity between θ'_0 and θ_0 in terms of their predictions** therefore can be used as a proxy to measure **how much discriminative knowledge θ_1 , after adapted to P_1 , preserves**. It is worth noting that measuring the similarity between θ'_0 and θ_0 requires no labeled data from P_0 or P_1 .

Now let us consider a more general case that M is not limited to 1. Let us set $\theta_0 = \theta_S$, where θ_S has been trained with the labeled source data \mathcal{S} till convergence. Based on the concept mentioned

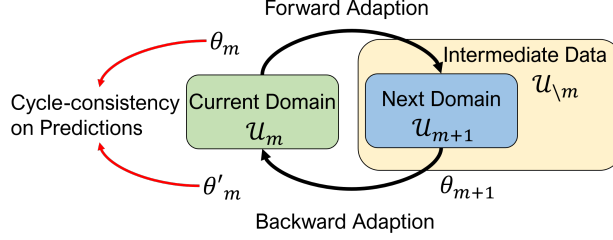


Figure 2: Cycle-consistency (cf. Equation 8).

above, we propose a novel learning framework to discover the domain sequences,

$$\begin{aligned} \arg \min_{(\mathcal{U}_1, \dots, \mathcal{U}_{M-1})} & \mathbb{E}_{\mathbf{x} \sim P_0} [\ell(f(\mathbf{x}; \theta'_0), \text{sharpen}(f(\mathbf{x}; \theta_0)))] , \\ \text{s.t., } & \theta'_0 = \text{ST}(\theta_M, (\mathcal{U}_{M-1}, \dots, \mathcal{U}_0)), \\ & \theta_M = \text{ST}(\theta_0, (\mathcal{U}_1, \dots, \mathcal{U}_M)), \end{aligned} \quad (7)$$

where \mathcal{U}_0 is the source data \mathcal{S} with labels removed, and $\mathcal{U}_M = \mathcal{T}$ is the target data. θ'_0 is a gradually self-trained model from θ_0 that undergoes a cycle $\mathcal{U}_1 \rightarrow \dots \rightarrow \mathcal{U}_M \rightarrow \dots \rightarrow \mathcal{U}_0$. In other words, Equation 7 aims to find the intermediate domain sequences such that performing gradual self-training along it in a forward and then backward fashion leads to cycle-consistency [76].

Optimization. Solving Equation 7 is still not trivial due to its combinatorial nature. We therefore propose to solve it **greedily, starting with discovering \mathcal{U}_1 , and then \mathcal{U}_2 , and then so on.** Concretely, we decompose Equation 7 into a series of sub-problems, as illustrated in Figure 2, each is defined as

$$\begin{aligned} \arg \min_{\mathcal{U}_{m+1} \subset \mathcal{U}_{\setminus m}} & \frac{1}{|\mathcal{U}_m|} \sum_{\mathbf{x} \in \mathcal{U}_m} \ell(f(\mathbf{x}; \theta'_m), \text{sharpen}(f(\mathbf{x}; \theta_m))), \\ \text{s.t., } & \theta'_m = \text{ST}(\theta_{m+1}, \mathcal{U}_m), \\ & \theta_{m+1} = \text{ST}(\theta_m, \mathcal{U}_{m+1}), \end{aligned} \quad (8)$$

where θ_m is the model already adapted to the current domain \mathcal{U}_m , and $\mathcal{U}_{\setminus m} = \mathcal{U} \setminus \cup_{j=1}^m \mathcal{U}_j$ is the remaining data that have not been selected by the m intermediate domains so far. That is, each sub-problem discovers only the next intermediate domain \mathcal{U}_{m+1} via cycle-consistency.

Let $N = |\mathcal{U}_{\setminus m}|$ and $\mathcal{U}_{\setminus m} = \{\mathbf{x}_i\}_{i=1}^N$, selecting $\mathcal{U}_{m+1} \subset \mathcal{U}_{\setminus m}$ is equivalent to selecting a binary indicator vector $\mathbf{q} \in \{0, 1\}^N$, in which $q_i = 1$ ² means that $\mathbf{x}_i \in \mathcal{U}_m$ is included into \mathcal{U}_{m+1} . This representation turns the self-training step $\theta_{m+1} = \text{ST}(\theta_m, \mathcal{U}_{m+1})$ in Equation 8 into

$$\text{ST}(\theta_m, \mathbf{q}) = \arg \min_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N q_i \times \ell(f(\mathbf{x}_i; \theta), \text{sharpen}(f(\mathbf{x}_i; \theta_m))), \quad (9)$$

We choose to solve Equation 8 approximately by relaxing the binary vector $\mathbf{q} \in \{0, 1\}^N$ into a real vector $\mathbf{q} \in \mathbb{R}^N$, which fits perfectly into the meta-reweighting framework [29, 55]. Concretely, meta-reweighting treats $\mathbf{q} \in \mathbb{R}^N$ as learnable differentiable parameters associated to training examples (in our case, the data in $\mathcal{U}_{\setminus m}$). Meta-reweighting for Equation 8 can be implemented via the following six steps for multiple iterations.

1. Detach: $\theta \leftarrow \theta_m$,
2. Forward: $\theta(\mathbf{q}) \leftarrow \theta - \frac{\eta_\theta}{|\mathcal{U}_{\setminus m}|} \times \frac{\partial \sum_{i \in \mathcal{U}_{\setminus m}} q_i \times \ell(f(\mathbf{x}_i; \theta), \text{sharpen}(f(\mathbf{x}_i; \theta_m)))}{\partial \theta}$,
3. Detach: $\theta' \leftarrow \theta(\mathbf{q})$,
4. Backward: $\theta(\mathbf{q}) \leftarrow \theta(\mathbf{q}) - \frac{\eta_\theta}{|\mathcal{U}_m|} \times \frac{\partial \sum_{j \in \mathcal{U}_m} \ell(f(\mathbf{x}_j; \theta(\mathbf{q})), \text{sharpen}(f(\mathbf{x}_j; \theta'))) }{\partial \theta(\mathbf{q})}$,
5. Update: $\mathbf{q} \leftarrow \mathbf{q} - \frac{\eta_q}{|\mathcal{U}_m|} \times \frac{\partial \sum_{j \in \mathcal{U}_m} \ell(f(\mathbf{x}_j; \theta(\mathbf{q})), \text{sharpen}(f(\mathbf{x}_j; \theta_m)))}{\partial \mathbf{q}}$,
6. Update: $q_i \leftarrow \max\{0, q_i\}$.

²Here we use the same notation as the coarse scores defined in subsection 3.3, since later we will initialize these values indeed by the coarse scores.

The aforementioned updating rule gives x_i a higher value of q_i if it helps preserve the discriminative knowledge. After obtaining the updated $\mathbf{q} \in \mathbb{R}^N$, we then sort it in the descending order and select the top $\frac{|\mathcal{U}|}{M-1}$ examples to be \mathcal{U}_{m+1} (i.e., every intermediate domain has an equal size).

Discussion. The way we relax the binary-valued vector $\mathbf{q} \in \{0, 1\}^N$ by real values is related to the linear programming relaxation of an integer programming problem. It has been widely applied, for example, in discovering sub-domains within a dataset [16]. Theoretically, we should constrain each element of \mathbf{q} to be within $[0, 1]$. Empirically, we found that even without the clipping operation to upper-bound q_i (i.e., just performing $\max\{0, q_i\}$), the values in q_i do not explode and the algorithm is quite stable: we see a negligible difference of using upper-bounded q_i or not in the resulting accuracy of GDA. This is also consistent with the common practice of using meta-reweighting [29, 55].

Initialization. We initialize q_i by the coarse scores defined in subsection 3.3. This is for one important reason: Equation 8 searches for the next intermediate domain \mathcal{U}_{m+1} only based on the current intermediate domain \mathcal{U}_m . Thus, it is possible that \mathcal{U}_{m+1} will include some data points that help preserve the knowledge in \mathcal{U}_m but are distributionally deviated from the gradual transition between \mathcal{U}_m and the target domain \mathcal{T} . The scores defined in subsection 3.3 thus serve as the guidance; we can also view Equation 8 as a refinement step of the coarse scores. As will be seen in section 4, the qualities of the initialization and the refinement steps both play important roles in IDOL’s success.

4 Experiment

4.1 Setup

We mainly study two benchmark datasets used in [35]. Rotated MNIST [36] is a popular synthetic dataset to study distributional shifts. It gradually rotates the original 50,000 images in the MNIST [37] dataset with $[0, 60]$ degrees uniformly, where $[0, 5)/[5, 55)/[55, 60]$ degrees are for the source/intermediate/target domains, respectively. However, since original images have no ground truths of the rotation, the rotation annotations are not expected to be perfect. Portraits dataset [14] is a real-world gender classification dataset of a collection of American high school seniors from the year 1905 to 2013. Naturally, the portrait styles change along years (e.g., lip curvature [14], fashion styles). The images are first sorted by the years and split. For both datasets, each domain contains 2000 images, and 1,000 images are reserved for both the source and target domains for validation.

We follow the setup in [35]: each model is a convolutional neural network trained for 20 epochs for each domain consequently (including training on the source data), using Adam optimizer [32] with a learning rate 0.001, batch size 32, and weight decay 0.02. We use this optimizer as the default if not specified. Hyper-parameters of IDOL include $K = 2M$ rounds for progressive training and 30 epochs of refinement per step (with mini-batch 128), where $M = 19$ for the Rotated MNIST and $M = 7$ for the Portraits. More details are in the supplementary material.

To study how the intermediate domain sequences affect gradual self-training, we consider baselines including “Source only” and “UDA” that directly self-trains on the target data (\mathcal{T}) and all the intermediate data (\mathcal{U}). We compare different domain sequences including *Pre-defined* indexes, e.g., *rotations* and *years*. We further consider the challenging setup that all intermediate data are unindexed. *Random* sequences are by dividing the unlabeled data into $M - 1$ domains randomly, with self-training on the target in the end. For our IDOL, different domain scores introduced in subsection 3.3 are adopted for the coarse domain sequences, and we optionally apply cycle-consistency refinement to make them fine-grained.

4.2 Main study: comparison with pre-defined domain sequences

The main results on Rotated MNIST and Portraits are provided in Table 1. We first verify that without any domain sequences, UDA (on target and intermediate data) and GDA (on random sequences) do not perform well, though they slightly improve over the source model. On the contrary, training with pre-defined indexes that arguably guide the gradual adaptation significantly performs better.

IDOL is competitive to pre-defined sequences. IDOL can produce meaningful domain sequences (see Figure 3), *given only the intermediate unlabeled data without any order*. The domain discriminator with progressive training performs the best and is comparable to pre-defined sequences. The confidence scores by the classifier are not enough to capture domain shifts.

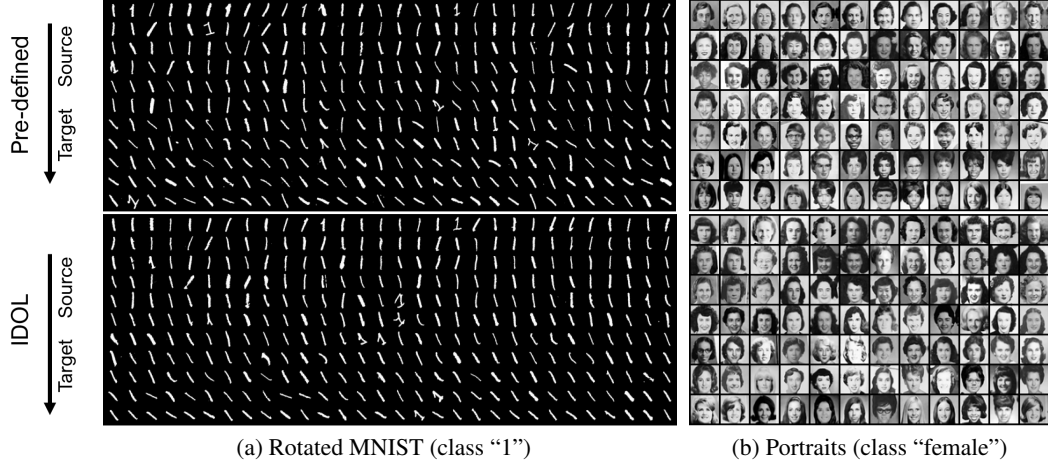


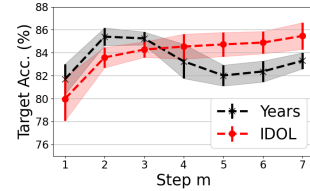
Figure 3: Random samples (the most representative class) from pre-defined indexes or IDOL sequences.

Table 1: Gradual self-training on Rotated MNIST and Portraits. Bottom section: IDOL with domain scores.

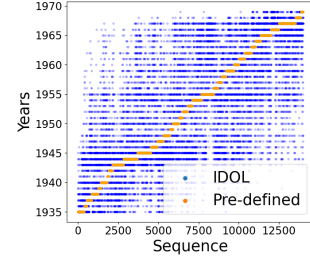
Coarse scores	Indexed?	Adaptation	Refined?	Rotated MNIST	Portraits
None	✗	Source only	-	31.9±1.7	75.3±1.6
		UDA (\mathcal{T})	-	33.0±2.2	76.9±2.1
		UDA ($\mathcal{T} + \mathcal{U}$)	-	38.0±1.6	78.9±3.0
Pre-defined [35]	✓	GDA	✗ ✓	87.9±1.2 93.3±2.3	83.8±0.8 85.8±0.4
Random	✗	GDA	✗ ✓	39.5±2.0 57.5±2.7	81.1±1.8 82.5±2.2
			✗ ✓	45.5±3.5 72.4±3.1	79.3±1.7 81.9±0.8
Classifier confidence	✗	GDA	✗	82.1±2.7	82.3±0.9
Manifold distance			✗	85.7±2.7	83.4±0.8
Domain discriminator			✗	87.5±2.0	85.5±1.0
Progressive domain discriminator			✓		

Refinement helps and coarse initialization matters. We observe that it is crucial to have high-quality coarse sequences as the initialization; the fine indexes by cycle-consistency refinement could further improve the coarse sequences, validating its effectiveness. We note that the Rotated MNIST dataset treats the original MNIST data as 0-degree rotation and artificially rotates the data given a rotation index. However, for data in the original MNIST dataset, there already exist variations in terms of rotations. This can be seen in Figure 3: in the first row of 0-degree rotation based on the pre-defined indexes, the examples of digit 1 do have slight rotations. Interestingly, in this situation, IDOL could potentially capture the true rotation of each example to further improve GDA.

Analysis. To understand why refinement helps, we take a closer look at the Portraits dataset that is naturally sorted by *years*. Are *years* the best factor to construct the domain sequence? We compare it with the refined sequence by IDOL. In Figure 4, we monitor the target accuracy of each step of adaptation and find that the accuracy fluctuates by using years. Interestingly, learning with IDOL stably and gradually improves the target accuracy and ends up with a higher accuracy, validating that refinement with cycle-consistency indeed produces better sequences. Specifically, the IDOL sequence has a reasonable 0.727 correlation with *years* but they are not perfectly aligned. Several factors such as fashions, hairstyles, eyeglasses, and lip curvatures may not perfectly align with years [14]. Even within the same year, individuals could have variations in these factors as well. We thus hypothesize that IDOL can discover the domain sequence that reflects a smoother transition of these factors.



(a) Target Acc. of gradual ST.



(b) Years vs. IDOL sequence.

Figure 4: Portraits with year indexes or IDOL domain sequence.

We note that the intermediate domains constructed by IDOL (domain scores with or without refinement) may not match the target “class” distributions. We monitor the number of examples of class c in each domain \mathcal{U}_m , i.e., $|\mathcal{U}_{m,c}|$, and compute $\frac{1}{M-1} \sum_{m=1}^{M-1} \frac{\max_c \{|\mathcal{U}_{m,c}|\}}{\min_c \{|\mathcal{U}_{m,c}|\}}$, which ideally should be 1.0, i.e., class-balanced. We observe fine-grained indexes are indeed more class-balanced (coarse vs fine-grained): 1.33 vs 1.23 on Rotated MNIST and 1.27 vs 1.25 on Portraits.

4.3 Case study: learning with partial or outlier information

If partial or outlier information of the intermediate domains is given, is IDOL still helpful? We examine the question by three experiments. First, we consider *unsorted domains*: the data are grouped into $M - 1$ domains according to the pre-defined indexes but the domains are *unordered*. We find that if we sort the domains with the mean domain scores $\frac{1}{|\mathcal{U}_m|} \sum_{x_i \in \mathcal{U}_m} q_i$, then we can perfectly recover the pre-defined order. Second, we investigate if we have *fewer, coarsely separated pre-defined domains*: we double the size of $|\mathcal{U}_m|$. We find that on Rotated MNIST (w/ pre-defined indexes), the accuracy degrades by 11%. IDOL can recover the accuracy since IDOL outputs a sorted sequence of data points, from which we can construct more, fine-grained intermediate domains. Third, we investigate the addition of *outlier domains*, in which the starting/end intermediate domains do not match the source/target domains exactly. For example, on Rotated MNIST, the intermediate domains are expanded to cover $[-30, 90]$ degrees. GDA on such domain sequence with outlier domains degrades to 77.0 ± 2.0 accuracy. Refinement can still improve it to 81.3 ± 1.4 .

4.4 Case study: learning with low-quality intermediate data or indexes

In practice, the intermediate unlabeled data or the pre-defined domain indexes for them might not be perfect. We consider several realistic cases to show that our method can make GDA more robust to the quality of the intermediate data or the domain sequences. First, we investigate *fewer intermediate data*: Equation 4 implies that the target error increases if the intermediate data are too sparse. We repeat the experiments in subsection 4.2 but with only 30% intermediate data. Second, we consider *noisy indexes*: building upon 30% clean indexes, we further add more data with noisy, random indexes.

Figure 5 summarizes the results. We first notice that with only 30% intermediate data, even if the indexes are clean, the accuracy decreases drastically by 14% for MNIST and by 5% for Portraits. For MNIST, adding more noisily-indexed data actually decreases the accuracy. Portraits gains some improvements with the noisily-indexed data, suggesting that the task could benefit from learning with intermediate data (even without accurate indexes). One advantage of IDOL is that we can annotate domain indexes for any unindexed data. On both datasets, IDOL stably improves with more data available, significantly outperforming that without IDOL.

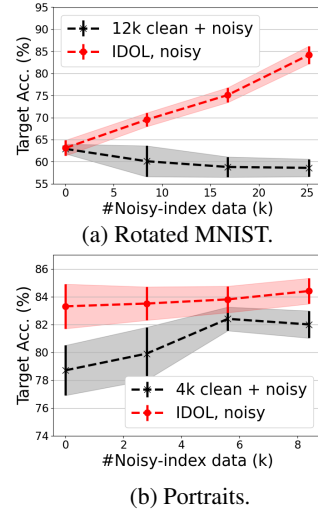


Figure 5: Different numbers of noisily-indexed intermediate data.

4.5 Case study: CIFAR10-STL UDA with additional, open-domain unlabeled data

We examine IDOL on a CIFAR10-to-STL [5, 34] UDA task which is known to be challenging due to the small data size [4]. The STL dataset has an additional unlabeled set, whose data are sub-sampled from ImageNet [9]. We use this set as the intermediate data, which contains unknown or outlier domains and classes. Here, we do not aim to compete with the state of the arts on the CIFAR10-to-STL UDA task, but study a more general application scenario for IDOL.

Method	Target Acc.
Source only	76.6 \pm 0.4
UDA (\mathcal{T}) (lr = 10^{-4})	69.4 \pm 0.4
UDA (\mathcal{T}) (lr = 10^{-5})	75.1 \pm 0.3
UDA ($\mathcal{T} + \mathcal{U}$)	61.1 \pm 0.8
GDA w/ conf.	77.1 \pm 0.5
GDA w/ IDOL	78.1 \pm 0.4

To leverage these data, we use $M = 3$ and filter out 80% of them using the classifier confidence. We first train a ResNet-20 [23] source model. In Table 2, we observe that direct UDA using self-training on STL or unlabeled data are inferior to the source model, even if the learning rate is carefully tuned. GDA with IDOL achieves the highest accuracy, demonstrating its robustness and wide applicability.

5 Related Work

Domain adaption. UDA has been studied extensively [1, 53, 63], and many different approaches have been proposed such as minimizing domain divergence [13, 20, 45, 61, 62, 65], cross-domain [39, 57, 58, 61] and domain-invariant features [18, 26, 54, 67–69, 73]. Self-training recently emerges as a simple yet effective approach [2, 28, 30, 31, 41, 64, 77]. It uses the source model to provide pseudo-labels for unlabeled target data, and approaches UDA via supervised learning in the target domain [38, 48, 49]. Theoretical analyses for applying self-training in UDA are provided in [3, 70].

Gradual domain adaption. Many UDA works show the benefits of gradually bridging the domain gap. With features as input, several work [8, 15, 19] construct Grassmannian manifolds between the source and the target. For deep models, progressive DA [27] proposes to create the synthetic intermediate domain with an image-to-image translation network. Other works [7, 17] instead generate intermediate data with jointly-trained generators. Na et al. [51] augments the intermediate data with Mixup [74]. Unlike all the above work that focus on building synthetic intermediate domains given only the source and target data, gradual domain adaption proposed in [12, 25, 35, 71] studies how to leverage extra real intermediate data with pre-defined domain indexes.

Learning with cycle consistency. The concept of cycle consistency has been successfully applied in many machine learning tasks. On a high level, cycle supervision enforces that the translation to the target should be able to be translated back to match the source. For instance, back-translation [22, 59] is a popular technique to perform unsupervised neural translation in natural language processing. Many computer vision tasks such as image generation [56], image segmentation [40], style transfer [76], etc., can also be improved by matching the structural outputs with cyclic signals. We propose a novel approach that uses cycle consistency to discover intermediate domains for gradual adaption.

Discovering feature subspaces. We further discuss the connections of our work to existing efforts about subspace learning. One powerful technique is subspace clustering [10, 33, 52] that partitions data into many subspaces (e.g., classes) unsupervisedly. However, it may not be applied to GDA directly since it groups with discriminative features which might not capture domain shifts. Specifically for UDA, some works aim to discover sub-spaces in the source/target domains [16, 24, 43, 47, 72]. The assumption is that the datasets might not be collected from one environment but generated from many different distributions, re-formulating it as a multi-domain adaption problem. We note that, the sub-domains discovered in these methods are fundamentally different from the domain sequence we discover. The sub-domains are assumed to have a dedicated adaption mapping to the target domain. The domain sequences in GDA are for constructing a path from the source to the target via intermediate data. How to link these learning techniques to GDA will be interesting future work.

6 Conclusion

Gradual domain adaptation leverages intermediate data to bridge the domain gap between the source and target domains. However, it relies on a strong premise — prior knowledge of domain indexes to sort the intermediate data into a domain sequence. We propose the IDOL algorithm that can produce comparable or even better domain sequences without pre-defined indexes. With a simple progressively-trained domain discriminator, it matches the performance of using pre-defined indexes, and further improves with the proposed cycle-consistency refinement. Essentially, IDOL is a useful tool to augment any GDA algorithms with high-quality domain sequences given unindexed data.

Acknowledgments and funding transparency statement

This research is partially supported by NSF IIS-2107077 and the OSU GI Development funds. We are thankful for the generous support of the computational resources by the Ohio Supercomputer Center.

References

- [1] 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, 10
- [2] Minmin Chen, Kilian Q Weinberger, and John Blitzer. Co-training for domain adaptation. In *NeurIPS*, 2011. 1, 10

- [3] Yining Chen, Colin Wei, Ananya Kumar, and Tengyu Ma. Self-training avoids using spurious features under domain shift. In *NeurIPS*, 2020. 10
- [4] Safa Cicek and Stefano Soatto. Unsupervised domain adaptation via regularized conditional alignment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1416–1425, 2019. 9
- [5] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223. JMLR Workshop and Conference Proceedings, 2011. 9
- [6] Remi Tachet des Combes, Han Zhao, Yu-Xiang Wang, and Geoff Gordon. Domain adaptation with conditional distribution matching and generalized label shift. In *NeurIPS*, 2020. 3
- [7] Shuhao Cui, Shuhui Wang, Junbao Zhuo, Chi Su, Qingming Huang, and Qi Tian. Gradually vanishing bridge for adversarial domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12455–12464, 2020. 10
- [8] Zhen Cui, Wen Li, Dong Xu, Shiguang Shan, Xilin Chen, and Xuelong Li. Flowing on riemannian manifold: Domain adaptation by shifting covariance. *IEEE transactions on cybernetics*, 44(12):2264–2273, 2014. 10
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 9
- [10] Ehsan Elhamifar and René Vidal. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE transactions on pattern analysis and machine intelligence*, 35(11):2765–2781, 2013. 10
- [11] Ali Farshchian, Juan A Gallego, Joseph P Cohen, Yoshua Bengio, Lee E Miller, and Sara A Solla. Adversarial domain adaptation for stable brain-machine interfaces. *arXiv preprint arXiv:1810.00045*, 2018. 2
- [12] Michael Gadermayr, Dennis Eschweiler, Barbara Mara Klinkhammer, Peter Boor, and Dorit Merhof. Gradual domain adaptation for segmenting whole slide images showing pathological variability. In *International Conference on Image and Signal Processing*, pages 461–469. Springer, 2018. 1, 10
- [13] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 17(1): 2096–2030, 2016. 1, 5, 10
- [14] Shiry Ginosar, Kate Rakelly, Sarah Sachs, Brian Yin, and Alexei A Efros. A century of portraits: A visual historical record of american high school yearbooks. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 1–7, 2015. 2, 7, 8
- [15] Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *2012 IEEE conference on computer vision and pattern recognition*, pages 2066–2073. IEEE, 2012. 1, 10
- [16] Boqing Gong, Kristen Grauman, and Fei Sha. Reshaping visual datasets for domain adaptation. *Advances in Neural Information Processing Systems*, 2013. 7, 10
- [17] Rui Gong, Wen Li, Yuhua Chen, and Luc Van Gool. Dlow: Domain flow for adaptation and generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2477–2486, 2019. 10
- [18] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 10
- [19] Raghuraman Gopalan, Ruonan Li, and Rama Chellappa. Domain adaptation for object recognition: An unsupervised approach. In *2011 international conference on computer vision*, pages 999–1006. IEEE, 2011. 10
- [20] Arthur Gretton, Alex Smola, Jiayuan Huang, Marcel Schmittfull, Karsten Borgwardt, and Bernhard Schölkopf. Covariate shift by kernel mean matching. *Dataset shift in machine learning*, 3(4):5, 2009. 10
- [21] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International Conference on Machine Learning*, pages 1321–1330. PMLR, 2017. 5

- [22] Di He, Yingce Xia, Tao Qin, Liwei Wang, Nenghai Yu, Tie-Yan Liu, and Wei-Ying Ma. Dual learning for machine translation. *Advances in neural information processing systems*, 29:820–828, 2016. 10
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 9, 16
- [24] Judy Hoffman, Brian Kulis, Trevor Darrell, and Kate Saenko. Discovering latent domains for multisource domain adaptation. In *European Conference on Computer Vision*, pages 702–715. Springer, 2012. 10
- [25] Judy Hoffman, Trevor Darrell, and Kate Saenko. Continuous manifold based adaptation for evolving visual domains. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 867–874, 2014. 1, 2, 10
- [26] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei A Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *ICML*, 2018. 10
- [27] Han-Kai Hsu, Chun-Han Yao, Yi-Hsuan Tsai, Wei-Chih Hung, Hung-Yu Tseng, Maneesh Singh, and Ming-Hsuan Yang. Progressive domain adaptation for object detection. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020. 10
- [28] Naoto Inoue, Ryosuke Furuta, Toshihiko Yamasaki, and Kiyoharu Aizawa. Cross-domain weakly-supervised object detection through progressive domain adaptation. In *CVPR*, 2018. 1, 10
- [29] Muhammad Abdullah Jamal, Matthew Brown, Ming-Hsuan Yang, Liqiang Wang, and Boqing Gong. Rethinking class-balanced methods for long-tailed visual recognition from a domain adaptation perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7610–7619, 2020. 6, 7, 15
- [30] Mehran Khodabandeh, Arash Vahdat, Mani Ranjbar, and William G Macready. A robust learning approach to domain adaptive object detection. In *ICCV*, 2019. 1, 10
- [31] Seunghyeon Kim, Jaehoon Choi, Taekyung Kim, and Changick Kim. Self-training and adversarial background regularization for unsupervised domain adaptive one-stage object detection. In *ICCV*, 2019. 1, 10
- [32] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 7, 16
- [33] Hans-Peter Kriegel, Peer Kröger, and Arthur Zimek. Clustering high-dimensional data: A survey on subspace clustering, pattern-based clustering, and correlation clustering. *Acm transactions on knowledge discovery from data (tkdd)*, 3(1):1–58, 2009. 10
- [34] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 9
- [35] Ananya Kumar, Tengyu Ma, and Percy Liang. Understanding self-training for gradual domain adaptation. In *International Conference on Machine Learning*, pages 5468–5479. PMLR, 2020. 1, 2, 3, 4, 7, 8, 10, 16, 17, 18
- [36] Hugo Larochelle, Dumitru Erhan, Aaron Courville, James Bergstra, and Yoshua Bengio. An empirical evaluation of deep architectures on problems with many factors of variation. In *Proceedings of the 24th international conference on Machine learning*, pages 473–480, 2007. 2, 7
- [37] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 7
- [38] Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, 2013. 1, 2, 3, 4, 10
- [39] Seungmin Lee, Dongwan Kim, Namil Kim, and Seong-Gyun Jeong. Drop to adapt: Learning discriminative features for unsupervised domain adaptation. In *ICCV*, 2019. 10
- [40] Yunsheng Li, Lu Yuan, and Nuno Vasconcelos. Bidirectional learning for domain adaptation of semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6936–6945, 2019. 10
- [41] Jian Liang, Ran He, Zhenan Sun, and Tieniu Tan. Distant supervised centroid shift: A simple and efficient approach to visual domain adaptation. In *CVPR*, 2019. 1, 10
- [42] Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint arXiv:1706.02690*, 2017. 4, 5

- [43] Ziwei Liu, Zhongqi Miao, Xingang Pan, Xiaohang Zhan, Dahua Lin, Stella X Yu, and Boqing Gong. Open compound domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12406–12415, 2020. 10
- [44] Mingsheng Long, Jianmin Wang, Guiguang Ding, Jianguang Sun, and Philip S Yu. Transfer feature learning with joint distribution adaptation. In *Proceedings of the IEEE international conference on computer vision*, pages 2200–2207, 2013. 1
- [45] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I Jordan. Learning transferable features with deep adaptation networks. *arXiv preprint arXiv:1502.02791*, 2015. 10
- [46] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. Conditional adversarial domain adaptation. *arXiv preprint arXiv:1705.10667*, 2017. 5
- [47] Massimiliano Mancini, Lorenzo Porzi, Samuel Rota Buló, Barbara Caputo, and Elisa Ricci. Boosting domain adaptation by discovering latent domains. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3771–3780, 2018. 10
- [48] David McClosky, Eugene Charniak, and Mark Johnson. Effective self-training for parsing. In *ACL*, 2006. 1, 10
- [49] David McClosky, Eugene Charniak, and Mark Johnson. Reranking and self-training for parser adaptation. In *ACL*, 2006. 1, 10
- [50] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018. 4, 17
- [51] Jaemin Na, Heechul Jung, HyungJin Chang, and Wonjun Hwang. Fixbi: Bridging domain spaces for unsupervised domain adaptation. *arXiv preprint arXiv:2011.09230*, 2020. 10
- [52] Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems*, 14:849–856, 2001. 10
- [53] Sinno Jialin Pan, Ivor W Tsang, James T Kwok, and Qiang Yang. Domain adaptation via transfer component analysis. *IEEE Transactions on Neural Networks*, 22(2):199–210, 2010. 10
- [54] Zhongyi Pei, Zhangjie Cao, Mingsheng Long, and Jianmin Wang. Multi-adversarial domain adaptation. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. 10
- [55] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *International Conference on Machine Learning*, pages 4334–4343. PMLR, 2018. 6, 7, 15, 17
- [56] Paolo Russo, Fabio M Carlucci, Tatiana Tommasi, and Barbara Caputo. From source to target and back: symmetric bi-directional adaptive gan. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8099–8108, 2018. 10
- [57] Kuniaki Saito, Yoshitaka Ushiku, Tatsuya Harada, and Kate Saenko. Adversarial dropout regularization. *arXiv preprint arXiv:1711.01575*, 2017. 10
- [58] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*, 2018. 10
- [59] Rico Sennrich, Barry Haddow, and Alexandra Birch. Improving neural machine translation models with monolingual data. *arXiv preprint arXiv:1511.06709*, 2015. 10
- [60] Tegjyot Singh Sethi and Mehmed Kantardzic. On the reliable detection of concept drift from streaming unlabeled data. *Expert Systems with Applications*, 82:77–99, 2017. 2
- [61] Rui Shu, Hung H Bui, Hirokazu Narui, and Stefano Ermon. A dirt-t approach to unsupervised domain adaptation. In *ICLR*, 2018. 10
- [62] Masashi Sugiyama, Matthias Krauledat, and Klaus-Robert MÅžller. Covariate shift adaptation by importance weighted cross validation. *JMLR*, 8(May):985–1005, 2007. 10
- [63] Baochen Sun, Jiashi Feng, and Kate Saenko. Return of frustratingly easy domain adaptation. In *AAAI*, 2016. 10
- [64] Qingyi Tao, Hao Yang, and Jianfei Cai. Zero-annotation object detection with web knowledge transfer. In *ECCV*, 2018. 10

- [65] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014. 10
- [66] Alexander Vergara, Shankar Vembu, Tuba Ayhan, Margaret A Ryan, Margie L Homer, and Ramón Huerta. Chemical gas sensor drift compensation using classifier ensembles. *Sensors and Actuators B: Chemical*, 166:320–329, 2012. 2
- [67] Riccardo Volpi, Pietro Morerio, Silvio Savarese, and Vittorio Murino. Adversarial feature augmentation for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5495–5504, 2018. 10
- [68] Ximei Wang, Liang Li, Weirui Ye, Mingsheng Long, and Jianmin Wang. Transferable attention for domain adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 5345–5352, 2019.
- [69] Yan Wang, Chen Xiangyu, You Yurong, Erran Li Li, Bharath Hariharan, Mark Campbell, Kilian Q. Weinberger, and Chao Wei-Lun. Train in germany, test in the usa: Making 3d object detectors generalize. In *CVPR*, 2020. 10
- [70] Colin Wei, Kendrick Shen, Yining Chen, and Tengyu Ma. Theoretical analysis of self-training with deep networks on unlabeled data. In *ICLR*, 2021. 10
- [71] Markus Wulfmeier, Alex Bewley, and Ingmar Posner. Incremental adversarial domain adaptation for continually changing environments. In *2018 IEEE International conference on robotics and automation (ICRA)*, pages 4489–4495. IEEE, 2018. 1, 2, 10
- [72] Caiming Xiong, Scott McCloskey, Shao-Hang Hsieh, and Jason Corso. Latent domains modeling for visual domain adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 28, 2014. 10
- [73] Minghao Xu, Jian Zhang, Bingbing Ni, Teng Li, Chengjie Wang, Qi Tian, and Wenjun Zhang. Adversarial domain adaptation with domain mixup. In *AAAI*, 2020. 10
- [74] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018. 10
- [75] Han Zhao, Remi Tachet Des Combes, Kun Zhang, and Geoffrey Gordon. On learning invariant representations for domain adaptation. In *International Conference on Machine Learning*, pages 7523–7532. PMLR, 2019. 1, 5
- [76] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. 6, 10
- [77] Yang Zou, Zhiding Yu, BVK Vijaya Kumar, and Jinsong Wang. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *ECCV*, 2018. 1, 10
- [78] Yang Zou, Zhiding Yu, Xiaofeng Liu, BVK Kumar, and Jinsong Wang. Confidence regularized self-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5982–5991, 2019. 1