

Transferable Calibration with Lower Bias and Variance in Domain Adaptation

Ximei Wang, Mingsheng Long*, Jianmin Wang, and Michael I. Jordan[#]

School of Software, KLiss, BNRist, Tsinghua University [#]University of California, Berkeley
 wxm17@mails.tsinghua.edu.cn {mingsheng, jimwang}@tsinghua.edu.cn
 jordan@cs.berkeley.edu

Abstract

Domain Adaptation (DA) enables transferring a learning machine from a labeled source domain to an unlabeled target domain. While remarkable advances have been made, most of the existing DA methods focus on improving the target accuracy at inference. How to estimate the predictive uncertainty of DA models is vital for decision-making in safety-critical scenarios but remains the boundary to explore. In this paper, we delve into the open problem of *Calibration in DA*, which is extremely challenging due to the coexistence of domain shift and the lack of target labels. We first reveal the dilemma that DA models learn higher accuracy at the expense of well-calibrated probabilities. Driven by this finding, we propose Transferable Calibration (TransCal) to tackle this dilemma, achieving accurate calibration with lower bias and variance in a unified hyperparameter-free optimization framework. As a general post-hoc calibration method, TransCal can be easily applied to recalibrate existing DA methods. Its efficacy has been justified both theoretically and empirically.

1 Introduction

Deep neural networks (DNNs) achieve the state of the art predictive accuracy in machine learning tasks with the benefit of powerful ability to learn discriminative representations [30, 8, 50]. However, in real-world scenarios, it is hard (intolerably time-consuming and labor-expensive) to collect sufficient labeled data through manual labeling, causing DNNs to confront with challenges when generalizing the pre-trained model to a different domain with unlabeled data. To tackle this challenge, researchers propose to transfer knowledge from a different but related domain by leveraging the readily-available labeled data, a.k.a. domain adaptation (DA) [39].

Early domain adaptation methods bridge the source and target domains mainly by learning domain-invariant representations [32, 12] or instance importances [19, 11]. After breakthrough in deep neural networks (DNNs) has been achieved, they are widely believed to be able to learn more transferable features [30, 8, 50, 54], since they disentangle explanatory factors of variations. Recent works in deep domain adaptation can be mainly grouped into two categories: 1) *moment matching*. These methods align representations across domains by minimizing the discrepancy between feature distributions [46, 24, 26, 27, 23]; 2) *adversarial training*. These methods adversarially learn transferable feature representations by confusing a domain discriminator in a two-player game [10, 45, 25, 53].

While numerous domain adaptation methods have been proposed, most of them mainly focus on improving the accuracy in the target domain but fail to estimate the predictive uncertainty, falling short of a miscalibration problem [16]. The accuracy of a deep adapted model constitutes only one side of the coin, here we delve into the other side of the coin, i.e. *the calibration of accuracy and confidence*, which requires the model to output a probability that reflects the true frequency of an event. For

*Corresponding author: Mingsheng Long (mingsheng@tsinghua.edu.cn)

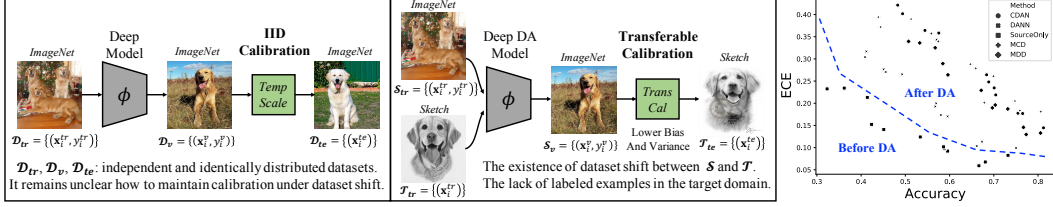


Figure 1: **Left:** A comparison between IID Calibration with TransCal, where ϕ denotes the deep model; **Right:** an observation on the accuracy and ECE of various DA methods (12 transfer tasks of Office-Home [47] with ResNet-50 [18]), indicating that DA models learn higher accuracy than the SourceOnly ones *at the expense of* well-calibrated probabilities. See more results in D.1 of Appendix.

example, if an automated diagnosis system says 1,000 patients have lung cancer with probability 0.1, approximately 100 of them should indeed have lung cancer. Calibration is fundamental to deep neural models and of great significance for decision-making in safety-critical scenarios. With built-in [9, 21] or post-hoc [37, 16] recalibration methods, the confidence and accuracy of deep models can be well-calibrated in the independent and identically distributed (IID) scenarios. However, it remains unclear how to maintain calibration under dataset shifts, especially when we do not have labels from the target dataset, as in the general setting of Unsupervised Domain Adaptation (UDA). We identify two obstacles in the way of applying calibration to UDA:

- *The lack of labeled examples in the target domain.* We know that the existing successful post-hoc IID recalibration methods mostly rely on ground-truth labels in the validation set to select the optimal temperature [37, 16]. However, since ground-truth labels are not available in the target domain, it is not feasible to directly apply IID calibration methods to UDA.
- *Dataset shift entangled with the miscalibration of DNNs.* Since DNNs are believed to learn more transferable features [30, 50], many domain adaptation methods embed DNNs to implicitly close the domain shift and rely on DNNs to achieve higher classification accuracy. However, DNNs are prone to over-confidence [16], falling short of a miscalibration problem.

To this end, we study the open problem of *Calibration in DA*, which is extremely challenging due to the coexistence of domain gap and the lack of target labels. To figure out the calibration error on the target domain of domain adaptation models, we first delve into the predictions and confidences of the target dataset. By calculating the target accuracy and ECE [16] (a calibration error measure defined in 3.1) with various domain adaptation models before calibration, we found something interesting. As shown in the right panel of Figure 1, the accuracy increases from the weakest SourceOnly [18] model to the latest state-of-the-art MDD [53] model, while the ECE becomes larger as well. That is, after applying domain adaptation methods, miscalibration phenomena become severer compared with SourceOnly model, indicating that the domain adaptation models learn higher classification accuracy *at the expense of* well-calibrated probabilities. This dilemma is unacceptable in safety-critical scenarios, as we need higher accuracy while maintaining calibration. Worse still, the well-performed calibration methods in the IID setting cannot be directly applied to DA due to the domain shift.

To tackle the dilemma between accuracy and calibration, we propose a new Transferable Calibration (TransCal) method in DA, achieving accurate calibration with lower bias and variance in a unified hyperparameter-free optimization framework, while a comparison with IID calibration is shown in the left panel of Figure 1. Specifically, we first define a new calibration measure, *Importance Weighted Expected Calibration Error* (IWECE) to estimate the calibration error in the target domain in a transferable calibration framework. Next, we propose a *learnable meta parameter* to further reduce the estimation bias from the perspective of theoretical analysis. Meanwhile, we develop a *serial control variate* method to further reduce the variance of the estimated calibration error. As a general post-hoc calibration method, TransCal can be easily applied to recalibrate existing DA methods. This paper has the following contributions:

- We uncover a dilemma in the open problem of Calibration in DA: existing domain adaptation models learn higher classification accuracy *at the expense of* well-calibrated probabilities.
- We propose a Transferable Calibration (TransCal) method, achieving accurate calibration with lower bias and variance in a unified hyperparameter-free optimization framework.

- We conduct extensive experiments on various DA methods, datasets, and calibration metrics, while the effectiveness of our method has been justified both theoretically and empirically.

2 Related Work

2.1 Domain Adaptation

Existing domain adaptation methods can be mainly grouped into two categories: *moment matching* and *adversarial training*. Moment matching methods align feature distributions across domains by minimizing the distribution discrepancy, in which Maximum Mean Discrepancy [15] is adopted by DAN [24] and DDC [46], and Joint Maximum Mean Discrepancy is utilized by JAN [27]. Motivated by Generative Adversarial Networks (GAN) [13], DANN [10] introduces a domain discriminator to distinguish the source features from the target ones, which are generated by the feature extractor; the domain discriminator and feature extractor are competed in a two-player minimax game. Further, CDAN [25] conditions the adversarial domain adaptation models on discriminative information conveyed in the classifier predictions. MADA [34] uses multiple domain discriminators to capture multimodal structures for fine-grained domain alignment. ADDA [45] adopts asymmetric feature extractors while MCD [42] employs two classifiers consistent across domains. MDD [53] proposes a new domain adaptation margin theory and achieves an impressive performance. Though numerous DA methods have been proposed, most of them focus on improving target accuracy and rare attention has been paid to the predictive uncertainty, causing a miscalibration between accuracy and confidence.

Table 1: Comparisons among calibration methods for unsupervised domain adaptation (UDA).

Method	works with domain shift	works without target label	Bias Reduction	Variance Reduction
Platt Scaling [37]	✗	✗	✗	✗
Temp. Scaling [16]	✗	✗	✗	✗
CPCS [33]	✓	✓	✗	✗
TransCal (proposed)	✓	✓	✓	✓

2.2 Calibration

Among *binary* calibration methods, Histogram Binning [51] is a simple non-parametric one with either equal-width or equal-frequency bins; Isotonic Regression [52] is a strict generalization of histogram binning by jointly optimizing the bin boundaries and bin predictions; Differently, Platt Scaling [37] is a parametric one that transforms the logits of a classifier to probabilities. When extended to *multiclass*, there are two types of methods. 1) *built-in methods*: Monte Carlo dropout (MC-dropout) [9] is popular as it simply uses Dropout [43] during testing phase to estimate predictive uncertainty. Later, [21] finds out that the ensembles of neural networks can work. Further, Stochastic Variational Bayesian Inference (SVI) methods for deep learning [2, 28, 49] are shown effective. However, built-in methods require to modify the classifier learning algorithm or training procedure, which are complex to apply in DA. Thus, we prefer 2) *post-hoc approaches*, including various multi-class extensions of Platt scaling [37]: matrix scaling, vector scaling and temperature scaling [16]. Though remarkable advances of IID calibration are witnessed, it remains unclear how to maintain calibration under dataset shifts [20], especially when the target labels are unavailable in UDA case. Recently, [31] finds that traditional post-hoc IID recalibration methods such as temperature scaling fail to maintain calibration under distributional shift. A recent paper (CPCS) [33] considers calibration under dataset shift using domain adaptation as a base tool for alignment, while we focus on how to maintain calibration in DA. A detailed comparison of typical calibration methods is shown in Table 1.

3 Approach

This paper aims at designing a transferable calibration method with lower bias and variance in DA. Let \mathbf{x} denote the input of the network, \mathbf{y} be the label and d be the Bernoulli variable indicating to which domain \mathbf{x} belongs. In our terminology, the source domain distribution is $p(\mathbf{x})$ and the target domain distribution is $q(\mathbf{x})$. We are given a labeled source domain $\mathcal{D}_s = \{(\mathbf{x}_s^i, \mathbf{y}_s^i)\}_{i=1}^{n_s}$ with

n_s samples ($d = 1$), and an unlabeled target domain $\mathcal{D}_t = \{(\mathbf{x}_t^i)\}_{i=1}^{n_t}$ with n_t samples ($d = 0$). Following the logistics of transferable calibration framework, bias reduction, and variance reduction, we show the design philosophy of transferable calibration step by step based on theoretical analysis.

3.1 IID Calibration

Calibration Metrics. Given a deep neural model ϕ (parameterized by θ) which transforms the random variable input X into the class prediction \hat{Y} and its associated confidence \hat{P} , we can define the *perfect calibration* [16] as $\mathbb{P}(\hat{Y} = Y | \hat{P} = c) = c, \forall c \in [0, 1]$ where Y is the ground truth label. There are some typical metrics to measure calibration error: 1) *Negative Log-Likelihood (NLL)* [14], also known as the cross-entropy loss in the field of deep learning, serves as a proper scoring rule to measure a probabilistic model’s quality [17]. 2) *Brier Score (BS)* [3], defined as the squared error between $p(\mathbf{y}|\mathbf{x}, \theta)$ and \mathbf{y} , is also a proper scoring rule. 3) *Expected Calibration Error (ECE)* [29, 16] first partitions the interval of probability predictions into B bins where B_m is the indices of samples falling into the m -th bin, and then computes the weighted absolute difference between accuracy and confidence across bins:

$$\mathcal{L}_{\text{ECE}} = \sum_{m=1}^B \frac{|B_m|}{n} |\mathbb{A}(B_m) - \mathbb{C}(B_m)|, \quad (1)$$

where for each bin m , the accuracy is $\mathbb{A}(B_m) = |B_m|^{-1} \sum_{i \in B_m} \mathbf{1}(\hat{\mathbf{y}}_i = \mathbf{y}_i)$ and its confidence is $\mathbb{C}(B_m) = |B_m|^{-1} \sum_{i \in B_m} p(\hat{\mathbf{p}}_i | \mathbf{x}_i, \theta)$. ECE is easier to interpret and thereby more popular.

Temperature Scaling Calibration. Temperature scaling is one of the simplest, fastest, and effective IID Calibration methods [16]. Fixing the neural model trained on the training set \mathcal{D}_{tr} , temperature scaling first attains the optimal temperature T^* by minimizing the cross-entropy loss between the logit vectors \mathbf{z}_v scaled by temperature T and the ground truth label \mathbf{y}_v on the validation set \mathcal{D}_v as

$$T^* = \arg \min_T \mathbb{E}_{(\mathbf{x}_v, \mathbf{y}_v) \in \mathcal{D}_v} \mathcal{L}_{\text{NLL}}(\sigma(\mathbf{z}_v/T), \mathbf{y}_v), \quad (2)$$

where σ is the *softmax* function formalized as $\hat{y}_j = \exp(z_j) / \sum_{k=1}^K \exp(z_k)$ for K classes. After that, we transform the logit vector \mathbf{z}_{te} on the test set \mathcal{D}_{te} into calibrated probabilities by $p_{te} = \sigma(\mathbf{z}_{te}/T^*)$.

3.2 Transferable Calibration Framework

As mentioned above, the main challenge of extending temperature scaling method into domain adaptation (DA) setup is that the target calibration error $\mathbb{E}_q = \mathbb{E}_{\mathbf{x} \sim q} [\mathcal{L}_{(\cdot)}(\phi(\mathbf{x}), \mathbf{y})]$ is defined over the target distribution q where labels are inaccessible. However, if density ratio (a.k.a. importance weight) $w(\mathbf{x}) = q(\mathbf{x})/p(\mathbf{x})$ is known, we can estimate target calibration error by the source distribution p :

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim q} [\mathcal{L}_{(\cdot)}(\phi(\mathbf{x}), \mathbf{y})] &= \int_q \mathcal{L}_{(\cdot)}(\phi(\mathbf{x}), \mathbf{y}) q(\mathbf{x}) d\mathbf{x} \\ &= \int_p \frac{q(\mathbf{x})}{p(\mathbf{x})} \mathcal{L}_{(\cdot)}(\phi(\mathbf{x}), \mathbf{y}) p(\mathbf{x}) d\mathbf{x} = \mathbb{E}_{\mathbf{x} \sim p} [w(\mathbf{x}) \mathcal{L}_{(\cdot)}(\phi(\mathbf{x}), \mathbf{y})], \end{aligned} \quad (3)$$

which means $\mathbb{E}_{\mathbf{x} \sim p} [w(\mathbf{x}) \mathcal{L}_{(\cdot)}(\phi(\mathbf{x}), \mathbf{y})]$ is an *unbiased* estimator of the target calibration error \mathbb{E}_q . By simply replacing the NLL loss in Eq. (2) with the importance weighted ECE in Eq. (3), we can attain an optimal temperature suitable in the target domain. Since this kind of calibration is trained on the source data but can transfer to the target domain, we call it *transferable calibration*.

As for the specific type of calibration metric $\mathcal{L}_{(\cdot)}$, the calibration method under covariate shift (CPCS) [33] uses the Brier Score \mathcal{L}_{BS} . However, Brier Score conflates accuracy with calibration since it can be decomposed into two components: calibration error and refinement [7], making it insensitive to predicted probabilities associated with infrequent events [31]. Meanwhile, NLL is minimized if and only if the prediction recovers ground truth \mathbf{y} , however, it may over-emphasize tail probabilities [4]. Hence, we adopt ECE \mathcal{L}_{ECE} , an intuitive and informative calibration metric directly quantifying the goodness of calibration. One may concern that ECE is not a proper scoring rule since the optimum score may not correspond to a perfect prediction, however, as a non-invasive post-hoc recalibration method, the temperature scaling we utilize will not degrade the neural model’s prediction accuracy while maintaining calibration.

Previously, we assume that density ratio is known, however, it is not readily accessible in real-world applications. In this paper, we adopt a mainstream discriminative density ratio estimation method: LogReg [38, 1, 5], which uses Bayesian formula to derive the estimated density ratio from a logistic regression classifier that separates examples from the source and the target domains as

$$\hat{w}(\mathbf{x}) = \frac{q(\mathbf{x})}{p(\mathbf{x})} = \frac{v(\mathbf{x}|d=0)}{v(\mathbf{x}|d=1)} = \frac{P(d=1)P(d=0|\mathbf{x})}{P(d=0)P(d=1|\mathbf{x})}, \quad (4)$$

where v is a distribution over $(\mathbf{x}, d) \in \mathcal{X} \times \{0, 1\}$ and $d \sim \text{Bernoulli}(0.5)$ is a Bernoulli variable indicating to which domain \mathbf{x} belongs. With Eq. (4), the estimated density ratio $\hat{w}(\mathbf{x})$ can be decomposed into two parts, in which the first part $P(d=1)/P(d=0)$ is a constant weight factor that can be estimated with the sample sizes of source and target domains as n_s/n_t , and the second part $P(d=0|\mathbf{x})/P(d=1|\mathbf{x})$ is the ratio of target probability to source probability that can be directly estimated with the probabilistic predictions of the logistic regression classifier.

3.3 Bias Reduction by Learnable Meta Parameter

Through above analysis, we can reach the target calibration error by the estimated importance weights. As long as the true importance weight $w(\mathbf{x})$ is known, an *unbiased* estimation of the target calibration error is feasible. However, the gap between the estimated importance weights and ground-truth ones cannot be ignored, causing a *bias* between the estimated calibration error and the ground-truth calibration error in the target domain. We formalize this bias of calibration as

$$\begin{aligned} \left| \mathbb{E}_{\mathbf{x} \sim q} [\mathcal{L}_{\text{ECE}}^{\hat{w}(\mathbf{x})}] - \mathbb{E}_{\mathbf{x} \sim q} [\mathcal{L}_{\text{ECE}}^{w(\mathbf{x})}] \right| &= |\mathbb{E}_{\mathbf{x} \sim p} [\hat{w}(\mathbf{x}) \mathcal{L}_{\text{ECE}}(\phi(\mathbf{x}), \mathbf{y})] - \mathbb{E}_{\mathbf{x} \sim p} [w(\mathbf{x}) \mathcal{L}_{\text{ECE}}(\phi(\mathbf{x}), \mathbf{y})]| \\ &= |\mathbb{E}_{\mathbf{x} \sim p} [(w(\mathbf{x}) - \hat{w}(\mathbf{x})) \mathcal{L}_{\text{ECE}}(\phi(\mathbf{x}), \mathbf{y})]|. \end{aligned} \quad (5)$$

Note that the bias of estimated calibration error in the target domain is highly related to the estimation error of importance weights. Hence, we focus on the bias of importance weights and show that after applying some basic mathematical inequalities, the estimation bias can be bounded by

$$\begin{aligned} &|\mathbb{E}_{\mathbf{x} \sim p} [(w(\mathbf{x}) - \hat{w}(\mathbf{x})) \mathcal{L}_{\text{ECE}}(\phi(\mathbf{x}), \mathbf{y})]| \\ &\leq \sqrt{\mathbb{E}_{\mathbf{x} \sim p} [(w(\mathbf{x}) - \hat{w}(\mathbf{x}))^2] \mathbb{E}_{\mathbf{x} \sim p} [(\mathcal{L}_{\text{ECE}}(\phi(\mathbf{x}), \mathbf{y}))^2]} \quad (\text{Cachy - Schwarz Inequality}) \\ &\leq \frac{1}{2} (\mathbb{E}_{\mathbf{x} \sim p} [(w(\mathbf{x}) - \hat{w}(\mathbf{x}))^2] + \mathbb{E}_{\mathbf{x} \sim p} [(\mathcal{L}_{\text{ECE}}(\phi(\mathbf{x}), \mathbf{y}))^2]) \quad (\text{AM/GM Inequality}) \\ &\leq \frac{1}{2} (\mathbb{E}_{\mathbf{x} \sim p} [(w(\mathbf{x}) - \hat{w}(\mathbf{x}))^2] + \mathbb{E}_{\mathbf{x} \sim p} [\mathcal{L}_{\text{ECE}}(\phi(\mathbf{x}), \mathbf{y})]) \quad (\mathcal{L}_{\text{ECE}}(\phi(\mathbf{x}), \mathbf{y}) \leq 1), \end{aligned} \quad (6)$$

where AM/GM denotes the inequality of arithmetic and geometric means. It is noteworthy that the domain adaptation model ϕ is fixed since we consider transferable calibration as a post-hoc method. Therefore, we can safely bypass the second term of Eq. (6) and focus our attention on the first term. According to the standard *bounded importance weight assumption* [6], for some bound $M > 0$ we have $w(\mathbf{x}) \leq M$. Then for any \mathbf{x} s.t. $p(\mathbf{x}) \neq 0$, the following inequality holds:

$$\frac{1}{M+1} \leq p(\mathbf{x}) \leq 1, \quad \text{since } w(\mathbf{x}) = \frac{q(\mathbf{x})}{p(\mathbf{x})} = \frac{1-p(\mathbf{x})}{p(\mathbf{x})} = \frac{1}{p(\mathbf{x})} - 1. \quad (7)$$

Thus the first term in Eq. (6) can be further bounded by

$$\mathbb{E}_{\mathbf{x} \sim p} [(w(\mathbf{x}) - \hat{w}(\mathbf{x}))^2] = \mathbb{E}_{\mathbf{x} \sim p} \left[\left(\frac{p(\mathbf{x}) - \hat{p}(\mathbf{x})}{p(\mathbf{x})\hat{p}(\mathbf{x})} \right)^2 \right] \leq (M+1)^4 \mathbb{E}_{\mathbf{x} \sim p} [(p(\mathbf{x}) - \hat{p}(\mathbf{x}))^2]. \quad (8)$$

Plugging Eq. (8) into Eq. (6), we conclude that a smaller M can ensure a lower bias for the estimated weight $\hat{w}(\mathbf{x})$, leading to a smaller bias of the estimated target calibration error, which is also supported by the generalization bound for importance weighted domain adaptation (Theorem 1, [6]). To this end, what we should do is to find some techniques to control the upper bound M of importance weights. It seems that we can normalize each weight by the sum of all weights, leading to a smaller M . Still, only with self-normalization, a few bad samples with very large weights will dominate the estimation, and drastically explode the estimator. Further, can we clip those samples with very large weights by a given threshold? It seems feasible, but the threshold is task-specific and hard to preset, which is not an elegant solution that we pursue. Based on the above theoretical analysis, we propose

to introduce a learnable meta parameter λ ($0 \leq \lambda \leq 1$) to adaptively downscale the extremely large weights, which can decrease M and attain a bias-reduced target calibration error. Formally,

$$T^* = \arg \min_{T, \lambda} \mathbb{E}_{\mathbf{x} \sim p} [\tilde{w}(\mathbf{x}) \mathcal{L}_{\text{ECE}}(\phi(\mathbf{x}), \mathbf{y})], \quad \tilde{w}(\mathbf{x}_i) = (\hat{w}(\mathbf{x}_i))^\lambda / \sum_{i=1}^{n_s} (\hat{w}(\mathbf{x}_i))^\lambda. \quad (9)$$

By jointly optimizing the calibration objective in Eq. (9), we can attain an optimal temperature T^* for transferable calibration, along with a task-specific optimal λ^* for bias reduction. [44] also introduced a control value to importance weighting for model selection, but it was used as a hyperparameter. This work further makes itself learnable in a unified hyperparameter-free optimization framework.

3.4 Variance Reduction by Serial Control Variate

Through the above analysis, we enable transferable calibration and further reduce its bias. However, another main drawback of importance weighting is uncontrolled *variance* as the importance weighted estimator can be drastically exploded by a few bad samples with large weights. For simplicity, denote $\tilde{\mathbb{E}}_q = \mathbb{E}_{\mathbf{x} \sim p} [\tilde{w}(\mathbf{x}) \mathcal{L}_{\text{ECE}}(\phi(\mathbf{x}), \mathbf{y})]$ as $\mathbb{E}_{\mathbf{x} \sim p} \mathcal{L}_{\text{ECE}}^{\tilde{w}}$ hereafter. Replacing the target loss from Lemma 2 of [6] with the weighted calibration error, we can conclude that the variance of transferable calibration error can be bounded by Rényi divergence between p and q (A proof is provided in B.1 of Appendix):

$$\begin{aligned} \text{Var}_{\mathbf{x} \sim p} [\mathcal{L}_{\text{ECE}}^{\tilde{w}}] &= \mathbb{E}_{\mathbf{x} \sim p} [(\mathcal{L}_{\text{ECE}}^{\tilde{w}})^2] - (\mathbb{E}_{\mathbf{x} \sim p} [\mathcal{L}_{\text{ECE}}^{\tilde{w}}])^2 \\ &\leq d_{\alpha+1}(q \| p) (\mathbb{E}_{\mathbf{x} \sim p} \mathcal{L}_{\text{ECE}}^{\tilde{w}})^{1-\frac{1}{\alpha}} - (\mathbb{E}_{\mathbf{x} \sim p} \mathcal{L}_{\text{ECE}}^{\tilde{w}})^2, \quad \forall \alpha > 0. \end{aligned} \quad (10)$$

Apparently, lowering the variance of $\tilde{\mathbb{E}}_q$ results in more accurate estimation. First, Rényi divergence [40] between p and q can be reduced by deep domain adaptation methods [25, 10, 53]. Second, we further reduce the variance by the control variate method [22]. It introduces a related unbiased estimator t to the estimator u that we concern, achieving a new estimator $u^* = u + \eta(t - \tau)$ while $\mathbb{E}[t] = \tau$. As proved in A.2 of Appendix, $\text{Var}[u^*] \leq \text{Var}[u]$ and u^* has an optimal solution when $\hat{\eta} = -\text{Cov}(u, t) / \text{Var}[t]$. To reduce $\text{Var}_{\mathbf{x} \sim p} [\mathcal{L}_{\text{ECE}}^{\tilde{w}}]$, we first adopt the importance weight $\tilde{w}(\mathbf{x})$ as the control variate since the expectation of $\tilde{w}(\mathbf{x})$ is fixed: $\mathbb{E}_{\mathbf{x} \sim p} [\tilde{w}(\mathbf{x})] = 1$. Here, regard $\mathbb{E}_{\mathbf{x} \sim p} [\mathcal{L}_{\text{ECE}}^{\tilde{w}}]$ and $\tilde{w}(\mathbf{x})$ as u and t respectively, and we can attain a new unbiased estimator \mathbb{E}_q^* . When η achieves the optimal solution, the estimation of target calibration error with control variate is

$$\mathbb{E}_q^* = \tilde{\mathbb{E}}_q - \frac{1}{n_s} \frac{\text{Cov}(\mathcal{L}_{\text{ECE}}^{\tilde{w}}, \tilde{w}(\mathbf{x}))}{\text{Var}[\tilde{w}(\mathbf{x})]} \sum_{i=1}^{n_s} [\tilde{w}(\mathbf{x}_s^i) - 1]. \quad (11)$$

Further, we can add the prediction on the source domain $r(\mathbf{x}) = \mathbf{1}(\hat{\mathbf{y}} = \mathbf{y})$ as another control variate because its expectation is also fixed: $\mathbb{E}_{\mathbf{x} \sim p} [r(\mathbf{x})] = c$, since the accuracy should be equal to the confidence c on a perfect calibrated source model as defined in Section 3.1. In this way, control variate method can be easily extended into the serial version in which there is a collection of control variables: t_1, t_2 whose corresponding expectations are τ_1, τ_2 respectively. Formally,

$$\begin{aligned} u^* &= u + \eta_1(t_1 - \tau_1), \\ u^{**} &= u^* + \eta_2(t_2 - \tau_2). \end{aligned} \quad (12)$$

Plugging $r(\mathbf{x})$ as the second control variate into the bottom line of Eq. (12), we can further reduce the variance of target calibration error by the serial control variate method as

$$\mathbb{E}_q^{**} = \mathbb{E}_q^* - \frac{1}{n_s} \frac{\text{Cov}(\mathcal{L}_{\text{ECE}}^{\tilde{w}*}, r(\mathbf{x}))}{\text{Var}[r(\mathbf{x})]} \sum_{i=1}^{n_s} [r(\mathbf{x}_s^i) - c], \quad (13)$$

where $\mathcal{L}_{\text{ECE}}^{\tilde{w}*}$ is the estimated target calibration error after applying the control variate to weight $\tilde{w}(\mathbf{x})$. As a summary, the transferable calibration framework (3)–(4) is improved through: 1) *lowering bias* as (9); 2) *lowering variance* by deep adaptation as (10) and by serial control variate as (11) and (13).

4 Experiments

4.1 Setup

We fully verify our methods on five DA datasets: (1) *Office-Home* [47]: a dataset with 65 categories, consisting of 4 domains: *Artistic (A)*, *Clipart (C)*, *Product (P)* and *Real-World (R)*. (2) *VisDA-2017* [36], a *Simulation-to-Real* dataset with 12 categories. (3) *Sketch* [48], a large-scale dataset

transferring from ImageNet (**I**) to Sketch (**S**) with 1000 categories. (4) *DomainNet* [35]: a dataset with 345 categories, including 6 domains: *Infograph* (**I**), *Quickdraw* (**Q**), *Real* (**R**), *Sketch* (**S**), *Clipart* (**C**) and *Painting* (**P**). (5) *Office-31* [41] contains 31 categories from 3 domains: *Amazon* (**A**), *Webcam* (**W**), *DSLRL* (**D**). For each dataset, we randomly split it and use the *first 80 percent* for training and the *remaining 20 percent* data for validation. We run each experiment for 10 times. We denote *Vanilla* as the standard softmax method before calibration, *Oracle* as the temperature scaling method while the target labels are available. Detailed descriptions are included in C.1, C.2 and C.3 of Appendix .

Table 2: ECE (%) before and after various calibration methods on several DA methods and datasets.

Method	Dataset	Office-Home							Sketch	VisDA
	Transfer Task	A→C	A→P	A→R	C→A	C→P	C→R	Avg	I→S	S→R
MDD	Before Cal. (Vanilla)	33.6	18.7	13.0	28.9	22.9	19.0	22.7	19.7	30.5
	IID Cal. (Temp. Scaling)	28.7	16.4	9.3	<u>21.8</u>	16.5	12.1	17.5	14.7	29.1
	CPCS [33]	29.5	17.3	9.6	22.9	16.7	<u>11.8</u>	18.0	14.2	30.4
	TransCal (w/o Bias)	22.8	14.2	9.0	23.4	14.0	12.8	16.1	10.2	23.5
	TransCal (w/o Variance)	<u>20.9</u>	<u>12.1</u>	<u>6.8</u>	21.6	<u>10.2</u>	12.1	<u>13.9</u>	<u>9.7</u>	<u>17.2</u>
	TransCal (ours)	13.5	11.4	4.8	<u>21.8</u>	7.0	11.1	11.6	8.1	16.1
	Oracle	6.8	8.5	4.7	7.0	5.8	4.0	6.1	4.7	7.4
	Before Cal. (Vanilla)	39.4	28.8	20.5	33.9	27.9	20.1	28.4	18.3	25.7
	IID Cal. (Temp. Scaling)	21.8	22.0	15.1	22.5	20.5	9.1	18.5	13.0	23.2
	CPCS [33]	23.1	22.3	15.4	20.6	20.0	<u>9.0</u>	18.4	12.9	22.9
MCD	TransCal (w/o Bias)	18.5	26.3	15.7	<u>19.2</u>	17.1	8.5	17.6	12.5	10.2
	TransCal (w/o Variance)	<u>16.3</u>	19.3	3.6	21.7	9.1	9.1	<u>13.2</u>	<u>11.3</u>	<u>9.8</u>
	TransCal (ours)	13.1	<u>20.2</u>	<u>5.1</u>	15.5	<u>9.3</u>	9.1	12.0	10.2	7.8
	Oracle	5.6	9.4	2.3	7.1	7.4	2.5	5.7	3.6	1.8
	Before Cal. (Vanilla)	40.2	26.4	17.8	35.8	23.5	21.9	27.6	21.8	29.5
CDAN	IID Cal. (Temp. Scaling)	28.3	17.6	10.1	<u>21.2</u>	13.2	8.2	16.4	9.0	26.6
	CPCS [33]	24.0	17.8	7.1	22.6	11.8	8.9	15.4	6.9	26.4
	TransCal (w/o Bias)	<u>20.3</u>	<u>10.3</u>	<u>5.3</u>	20.7	13.2	5.3	<u>12.5</u>	6.1	24.3
	TransCal (w/o Variance)	25.6	16.6	5.5	<u>21.2</u>	<u>8.2</u>	<u>5.6</u>	13.8	<u>5.6</u>	<u>23.9</u>
	TransCal (ours)	13.2	9.9	5.2	<u>21.2</u>	8.1	6.4	10.7	4.9	21.2
	Oracle	5.8	8.1	4.8	10.0	7.7	4.2	6.8	2.5	2.0

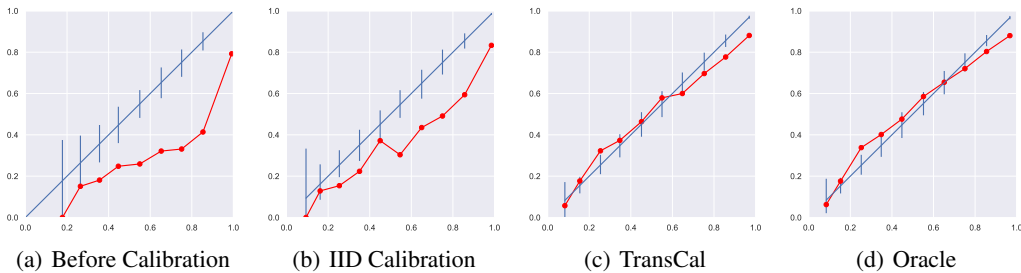


Figure 2: Reliability diagrams from *Clipart* to *Product* with CDAN [25] before and after calibration.

4.2 Results

Quantitative Results. As reported in Table 2, TransCal achieves much lower ECE than competitors (decreases about 30% or more, e.g. when TransCal is used to calibrate MCD on VisDA, the target ECE is reduced from 22.9 to 7.8) on various datasets and domain adaptation methods. Some results of TransCal are even approaching the Oracle ones. Further, the ablation studies on *TransCal (w/o Bias)* and *TransCal (w/o Variance)* verify that both bias reduction term and variance reduction term are effective. TransCal can be generalized to other tasks of Office-Home (D.2.1), to DomainNet and Office-31 (D.2.2), and to more DA methods (D.2.3), all shown in Appendix. Further, the results evaluated by NLL and BS metrics are included in D.2.4 and D.2.5 of Appendix respectively.

Qualitative Results. As shown in Figure 2, the blue lines indicate the distributions for *perfectly* reliable forecasts with standard deviation, and the red lines denote the conditional distributions of the observations. Obviously, If the model is perfectly calibrated, these two lines should be matched. We can see that TransCal is much better and approaches the *Oracle* one on the task: Clipart \rightarrow Product. More reliability diagrams of other tasks to back up this conclusion are shown in D.3 of Appendix .

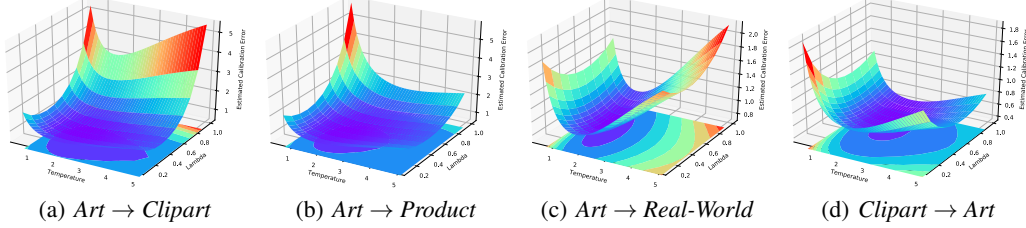


Figure 3: The estimated calibration error with respect to different values of temperature T and meta parameter λ (both are *learnable*), showing that different models achieve optimal values at different λ .

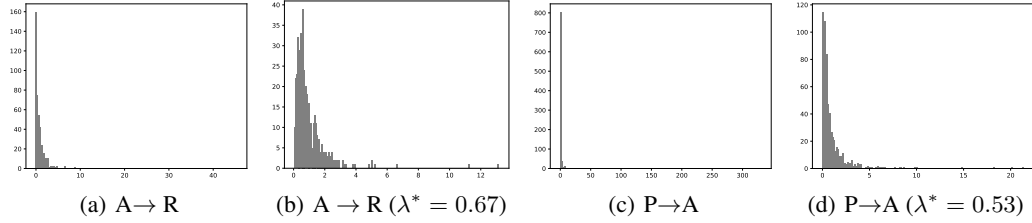


Figure 4: Importance Weight distribution of two DA tasks after transferable calibration with (4(b), 4(d)) and without (4(a), 4(c)) applying the learnable meta parameter, which lowers the value of M .

Table 3: ECE (%) of TransCal with different control variate (CV) methods on MDD [53].

Dataset	Office-Home			Sketch	VisDA
Transfer Task	A \rightarrow C	A \rightarrow P	A \rightarrow R	I \rightarrow S	S \rightarrow R
TransCal (w/o Control Variate)	20.9 \pm 4.68	12.1 \pm 2.46	6.8 \pm 2.22	9.7 \pm 3.17	17.2 \pm 5.74
TransCal (CV via only $w(\mathbf{x})$)	13.9 \pm 4.45	9.6 \pm 1.52	5.9 \pm 1.91	9.3 \pm 1.68	16.4 \pm 5.68
TransCal (CV via only $r(\mathbf{x})$)	13.8 \pm 4.32	10.2 \pm 0.97	5.2 \pm 1.08	8.6 \pm 1.37	16.3 \pm 3.32
TransCal (Parallel Control Variate)	13.6 \pm 4.43	10.6 \pm 1.46	5.2 \pm 1.45	8.7 \pm 1.54	16.3 \pm 3.45
TransCal (Serial Control Variate)	13.5 \pm 3.51	11.4 \pm 0.81	4.8 \pm 0.76	8.1 \pm 1.09	16.1 \pm 1.20

4.3 Insight Analyses

Why Bias Reduction Term Works. From the perspective of optimization, we explore the estimated calibration error with respect to different values of temperature (T) and lambda (λ) in Figure 3, showing that different models achieve optimal values at different λ . Therefore, it is impossible to attain optimal estimated calibration error by presetting a fixed λ . However, with our unified meta-parameter optimization framework, we can adaptively find an optimal λ for each task. From the perspective of importance weight distribution as shown in Figure 4, after applying learnable meta parameter λ , the highest values (M in Section 3.3) of importance weight decrease, leading to a smaller bias in Eq. (5).

Why Serial Control Variate Works. As the theoretical analysis in B.2 of Appendix shows, the variance of \mathbb{E}_q^{**} can be further reduced since $\text{Var}[\mathbb{E}_q^{**}] \leq \text{Var}[\mathbb{E}_q^*] \leq \text{Var}[\mathbb{E}_q]$, but other variants of control variate (CV) method such as Parallel CV may not hold this property. Meanwhile, as shown in Table 3, TransCal (Serial CV) not only achieves better calibration performance but also attains lower calibration variance than other variants of control variate methods.

5 Conclusion

In this paper, we delve into an open and important problem of *Calibration in DA*. We first reveal that domain adaptation models learn higher accuracy *at the expense of* well-calibrated probabilities. Further, we propose a novel transferable calibration (TransCal) approach, achieving more accurate calibration with lower bias and variance in a unified hyperparameter-free optimization framework.

Broader Impact

The open problem of *Calibration in DA* that we delve into is a very promising research direction and important for decision making in safety-critical applications, such as automated diagnosis system for lung cancer. Since our method can be easily applied to recalibrate the existing DA methods and generate more reliable predictions, it will benefit the transfer learning community. If the method fails in some extreme circumstances, it will confuse researchers or engineers who apply our method but it will not bring about any negative ethical or societal consequences. Meanwhile, our method did not leverage biases in the data such as racial discrimination and gender discrimination since we conduct experiments on standard domain adaptation datasets that are more about animals or pieces of equipment in the office. In summary, we hold a positive view of the broader impact on this paper.

References

- [1] S. Bickel and T. Scheffer. Dirichlet-enhanced spam filtering based on biased samples. In *NeurIPS*. MIT Press, 2007.
- [2] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra. Weight uncertainty in neural network. In *ICML*, 2015.
- [3] G. W. BRIER. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78(1):1–3, 1950.
- [4] J. Q. Candela, C. E. Rasmussen, F. H. Sinz, O. Bousquet, and B. Schölkopf. Evaluating predictive uncertainty challenge. In *Machine Learning Challenges*, 2005.
- [5] K. F. Cheng and C. K. Chu. Semiparametric density estimation under a two-sample density ratio model. *Bernoulli*, 10(4):583–604, 2004.
- [6] C. Cortes, Y. Mansour, and M. Mohri. Learning bounds for importance weighting. In *Advances in Neural Information Processing Systems 23*, pages 442–450. Curran Associates, Inc., 2010.
- [7] M. H. DeGroot and S. E. Fienberg. The comparison and evaluation of forecasters. *Journal of the Royal Statistical Society. Series D (The Statistician)*, 32(1/2):12–22, 1983.
- [8] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *ICML*, 2014.
- [9] Y. Gal and Z. Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *ICML*, 2016.
- [10] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 2016.
- [11] B. Gong, K. Grauman, and F. Sha. Connecting the dots with landmarks: Discriminatively learning domain-invariant features for unsupervised domain adaptation. In *ICML*, 2013.
- [12] B. Gong, Y. Shi, F. Sha, and K. Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *CVPR*, 2012.
- [13] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NeurIPS*, 2014.
- [14] I. J. Goodfellow, Y. Bengio, and A. C. Courville. *Deep Learning*. Adaptive computation and machine learning. MIT Press, 2016.
- [15] A. Gretton, D. Sejdinovic, H. Strathmann, S. Balakrishnan, M. Pontil, K. Fukumizu, and B. K. Sriperumbudur. Optimal kernel choice for large-scale two-sample tests. In *NeurIPS*, 2012.
- [16] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On calibration of modern neural networks. In *ICML*, 2017.
- [17] T. Hastie, R. Tibshirani, and J. H. Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, 2nd Edition*. Springer Series in Statistics. Springer, 2009.

- [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [19] J. Huang, A. J. Smola, A. Gretton, K. M. Borgwardt, and B. Schölkopf. Correcting sample selection bias by unlabeled data. In *NeurIPS*, 2006.
- [20] A. Kumar, P. Liang, and T. Ma. Verified uncertainty calibration. In *NeurIPS*, 2019.
- [21] B. Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *NeurIPS*, 2017.
- [22] C. Lemieux. Control variates. In *Wiley StatsRef: Statistics Reference Online*, pages 1–8. American Cancer Society, 2017.
- [23] M. Long, Y. Cao, Z. Cao, J. Wang, and M. I. Jordan. Transferable representation learning with deep adaptation networks. *TPAMI*, 2018.
- [24] M. Long, Y. Cao, J. Wang, and M. I. Jordan. Learning transferable features with deep adaptation networks. In *ICML*, 2015.
- [25] M. Long, Z. Cao, J. Wang, and M. I. Jordan. Conditional adversarial domain adaptation. In *NeurIPS*, 2018.
- [26] M. Long, H. Zhu, J. Wang, and M. I. Jordan. Unsupervised domain adaptation with residual transfer networks. In *NeurIPS*, 2016.
- [27] M. Long, H. Zhu, J. Wang, and M. I. Jordan. Deep transfer learning with joint adaptation networks. In *ICML*, 2017.
- [28] C. Louizos and M. Welling. Multiplicative normalizing flows for variational bayesian neural networks. In *ICML*, 2017.
- [29] M. P. Naeini, G. F. Cooper, and M. Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *AAAI*, 2015.
- [30] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *CVPR*, 2014.
- [31] Y. Ovadia, E. Fertig, and J. Ren. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In *NeurIPS*, 2019.
- [32] S. J. Pan and Q. Yang. A survey on transfer learning. *TKDE*, 22(10):1345–1359, 2010.
- [33] S. Park, O. Bastani, J. Weimer, and I. Lee. Calibrated prediction with covariate shift via unsupervised domain adaptation. 2020.
- [34] Z. Pei, Z. Cao, M. Long, and J. Wang. Multi-adversarial domain adaptation. In *AAAI*, 2018.
- [35] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang. Moment matching for multi-source domain adaptation. *ICCV*, 2019.
- [36] X. Peng, B. Usman, N. Kaushik, J. Hoffman, D. Wang, and K. Saenko. Visda: The visual domain adaptation challenge. *CoRR*, abs/1710.06924, 2017.
- [37] J. C. Platt. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. In *Advances in Linear Large Margin Classifiers*. MIT Press, 1999.
- [38] J. Qin. Inferences for case-control and semiparametric two-sample density ratio models. *Biometrika*, 85(3):619–630, 1998.
- [39] J. Quionero-Candela, M. Sugiyama, A. Schwaighofer, and N. D. Lawrence. *Dataset Shift in Machine Learning*. The MIT Press, 2009.
- [40] A. Rényi. On measures of information and entropy. 1961.

- [41] K. Saenko, B. Kulis, M. Fritz, and T. Darrell. Adapting visual category models to new domains. In *ECCV*, 2010.
- [42] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*, 2018.
- [43] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, 2014.
- [44] M. Sugiyama, M. Krauledat, and K.-R. M  ller. Covariate shift adaptation by importance weighted cross validation. 8:985–1005, 2007.
- [45] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell. Adversarial discriminative domain adaptation. In *CVPR*, 2017.
- [46] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell. Deep domain confusion: Maximizing for domain invariance. *CoRR*, abs/1412.3474, 2014.
- [47] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Panchanathan. Deep Hashing Network for Unsupervised Domain Adaptation. In *CVPR*, 2017.
- [48] H. Wang, S. Ge, E. P. Xing, and Z. C. Lipton. Learning robust global representations by penalizing local predictive power. *NeurIPS*, 2019.
- [49] Y. Wen, P. Vicol, J. Ba, D. Tran, and R. B. Grosse. Flipout: Efficient pseudo-independent weight perturbations on mini-batches. In *ICLR*, 2018.
- [50] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *NeurIPS*, 2014.
- [51] B. Zadrozny and C. Elkan. Obtaining calibrated probability estimates from decision trees and naive bayesian classifiers. In *ICML*, 2001.
- [52] B. Zadrozny and C. Elkan. Transforming classifier scores into accurate multiclass probability estimates. In *KDD*, 2002.
- [53] Y. Zhang, T. Liu, M. Long, and M. Jordan. Bridging theory and algorithm for domain adaptation. In *ICML*, 2019.
- [54] H. Zhao, R. T. des Combes, K. Zhang, and G. J. Gordon. On learning invariant representation for domain adaptation. In *ICML*, 2019.