# Butterfly: One-step Approach towards Wildly Unsupervised Domain Adaptation

Feng Liu, *Member, IEEE,* Jie Lu, *Fellow, IEEE,*
Bo Han, Gang Niu, Guangquan Zhang and Masashi Sugiyama

**Abstract**—In *unsupervised domain adaptation* (UDA), classifiers for the *target domain* (TD) are trained with *clean* labeled data from the *source domain* (SD) and unlabeled data from TD. However, in the wild, it is difficult to acquire a large amount of perfectly clean labeled data in SD given limited budget. Hence, we consider a new, more realistic and more challenging problem setting, where classifiers have to be trained with *noisy* labeled data from SD and unlabeled data from TD—we name it *wildly UDA* (WUDA). We show that WUDA ruins all UDA methods if taking no care of label noise in SD, and to this end, we propose a *Butterfly* framework, a powerful and efficient solution to WUDA. Butterfly maintains four deep networks simultaneously, where two take care of all adaptations (i.e., noisy-to-clean, labeled-to-unlabeled, and SD-to-TD-distributional) and then the other two can focus on classification in TD. As a consequence, Butterfly possesses all the conceptually necessary components for solving WUDA. Experiments demonstrate that, under WUDA, Butterfly significantly outperforms existing baseline methods. The code of Butterfly can be found at github.com/fengliu90/Butterfly.

**Index Terms**—machine learning, weakly-supervised learning, transfer learning

✦

## 1 INTRODUCTION

DOMAIN adaptation (DA) aims to learn a discriminative classifier in the presence of a shift between training data in source domain and test data in target domain [1], [2], [3], [4], [5], [6]. Currently, DA can be divided into three categories: *supervised DA* [7], *semi-supervised DA* [8], [9], [10], [11], [12] and *unsupervised DA* (UDA) [13], [14], [15], [16], [17], [18], [19], [20]. When the number of labeled data is few in target domain, supervised DA is also known as *few-shot DA* [21]. Since unlabeled data in target domain can be easily obtained, UDA exhibits the greatest potential in the real world [22], [23], [24], [25], [26], [27], [28], [29], [30], [31], [32].

UDA methods train with clean labeled data in a source domain (i.e., clean source data) and unlabeled data in a target domain (i.e., unlabeled target data) to obtain classifiers for the *target domain* (TD), which mainly consist of three orthogonal techniques: *integral probability metrics* (IPM) [14], [33], [34], [35], [36], [37], [38], *adversarial training* [23], [25], [39], [40], [41], [42], [43] and *pseudo labeling* [13]. Compared to IPM- and adversarial-training-based methods, the pseudo-labeling-based method (i.e., *asymmetric tri-training domain adaptation* (ATDA) [13]) can construct a high-quality target-specific representation, providing a better classification performance.

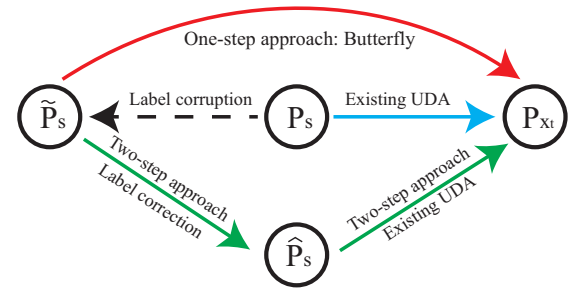However, in the wild, the data volume of the source



Fig. 1. Wildly unsupervised domain adaptation (WUDA). The blue line denotes that UDA transfers knowledge from clean source data ($P_s$) to unlabeled target data ($P_{x_t}$). However, perfectly clean data is hard to acquire. This brings *wildly unsupervised domain adaptation* (WUDA), namely transferring knowledge from noisy source data ($\widetilde{P}_s$) to unlabeled target data ($P_{x_t}$). Note that label corruption process (black dash line) is unknown in practice. To handle WUDA, a compromise solution is a two-step approach (green line), which sequentially combines label-noise algorithms ($\widetilde{P}_s \rightarrow \hat{P}_s$, label correction) and existing UDA ($\hat{P}_s \rightarrow P_{x_t}$). This paper proposes a robust one-step approach called Butterfly (red line, $\widetilde{P}_s \rightarrow P_{x_t}$ directly), which eliminates noise effects from $\widetilde{P}_s$.

domain tends to be large [45]. To avoid the expensive labeling cost, labeled data in the source domain normally come from amateur annotators or the Internet [46], [47], [48]. This brings us a new, more realistic and more challenging problem, *wildy unsupervised domain adaptation* (abbreviated as WUDA, Figure 1). This adaptation aims to transfer knowledge from noisy labeled data in the source domain ($\widetilde{P}_s$, i.e., noisy source data) to unlabeled target data ($P_{x_t}$). Unfortunately, existing UDA methods share an implicit assumption that *there are no noisy source data* [44], [49]. Namely, these methods focus on transferring knowledge from clean source data ($P_s$) to unlabeled target data ($P_{x_t}$). Therefore, these methods cannot well handle WUDA (Figure 2).

- *Feng Liu is with the Centre for Artificial Intelligence (CAI), Faculty of Engineering and Information Technology, University of Technology Sydney, Sydney, NSW, Australia, and with Center for Advanced Intelligence Project, RIKEN, Japan. Jie Lu and Guangquan Zhang are with the Centre for Artificial Intelligence (CAI), Faculty of Engineering and Information Technology, University of Technology Sydney, Sydney, NSW, Australia. Bo Han is with Center for Advanced Intelligence Project, RIKEN, Japan, and with Department of Computer Science, Hong Kong Baptist University, Hong Kong SAR. Gang Niu is with Center for Advanced Intelligence Project, RIKEN, Japan. Masashi Sugiyama is with Center for Advanced Intelligence Project, RIKEN, Japan, and with Graduate School of Frontier Sciences, University of Tokyo, Japan.*

- *E-mail: {Feng.liu, jie.lu}@uts.edu.au, {bo.han, gang.niu}@riken.jp, guangquan.zhang@uts.edu.au and sugi@k.u-tokyo.ac.jp*

(a) Symmetry-flip noise: $S \rightarrow M$ (left), $M \rightarrow S$ (right)  (b) Pair-flip noise: $S \rightarrow M$ (left), $M \rightarrow S$ (right)
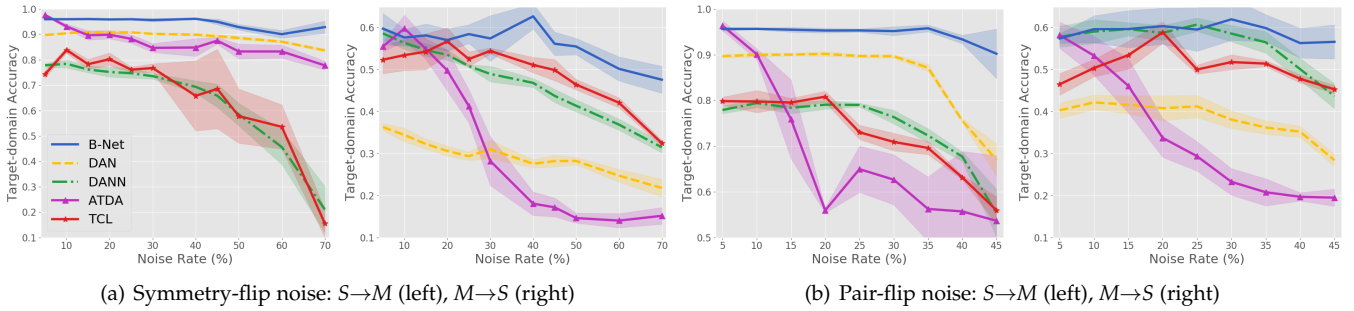
Fig. 2. WUDA ruins representative UDA methods. Representative UDA methods includes *deep adaptation network* (DAN, an IPM based method [36]), *domain-adversarial neural network* (DANN, an adversarial training based method [23]), *asymmetric tri-training domain adaptation* (ATDA, a pseudo-label based method [13]) and *transferable curriculum learning* (TCL, a robust UDA method [44]). B-Net is our proposed WUDA method. We report target-domain accuracy of all methods when the noise rate of source domain changes (a) from $5\%$ to $70\%$ (symmetry-flip noise) and (b) from $5\%$ to $45\%$ (pair-flip noise). Clearly, as the noise rate of source domain increases, the target-domain accuracy of representative UDA methods drops quickly while that of B-Net keeps stable consistently.

To validate this fact, we empirically reveal the deficiency of existing UDA methods (Figure 2, e.g., *deep adaptation network* (DAN) [36] and *domain-adversarial neural network* (DANN) [23]). To improve these methods, a straightforward solution is a two-step approach. In Figure 1, we can first use label-noise algorithms to train a classifier on noisy source data, then leverage this trained classifier to assign pseudo labels for noisy source data. Via UDA methods, we can transfer knowledge from pseudo-labeled source data ($\hat{P}_s$) to unlabeled target data ($P_{x_t}$). Nonetheless, pseudo-labeled source data are still noisy, and such two-step approach may not eliminate noise effects.

To circumvent the issue of two-step approach, we present a robust one-step approach called *Butterfly*. In high level, Butterfly directly transfers knowledge from $\widetilde{P}_s$ to $P_{x_t}$, and uses the transferred knowledge to construct target-specific representations. In low level, Butterfly maintains four networks dividing two branches (Figure 3): Two networks in Branch-I are jointly trained on noisy source data and pseudo-labeled target data (data in *mixture domain* (MD)); while two networks in Branch-II are trained on pseudo-labeled target data. Our ablation study (see Section 9.9) confirms the network design of Butterfly (see Section 7) is the optimal.

The reason why Butterfly can be robust takes root in the *dual-checking principle* (DCP): Butterfly checks high-correctness data out, from not only the data in MD but also the pseudo-labeled target data. After cross-propagating these high-correctness data, Butterfly can obtain high-quality *domain-invariant representations* (DIR) and *target-specific representations* (TSR) simultaneously in an iterative manner. If we only check data in MD (i.e., B-Net-M in Section 9.9), the error existed in pseudo-labeled target data will accumulate, leading to the low-quality DIR and TSR.

We conduct experiments on simulated WUDA tasks, including 4 *MNIST-to-SYND* tasks, 4 *SYND-to-MNIST* tasks and 24 human-sentiment tasks. Besides, we conduct experiments on 3 real-world WUDA tasks. Empirical results demonstrate that Butterfly can robustly transfer knowledge from noisy source data to unlabeled target data. Meanwhile, Butterfly performs much better than existing UDA methods when *source domain* (SD) suffers the extreme (e.g., $45\%$) noise.

## 2 LITERATURE REVIEW

This section reviews the existing UDA methods in detail.

UDA methods train with clean source data and unlabeled target data to classify target-domain data, which mainly consist of three orthogonal techniques: *integral probability metrics* (IPM) [14], [33], [34], [35], [36], *adversarial training* [23], [25], [39], [40], [41], [42] and *pseudo labeling* [13].

IPMs (such as maximum mean discrepancy [34], [50] and Wasserstein distance [35]) are used to measure the discrepancy between distributions of two domains. By minimizing the IPM between two domains, models trained with clean source data can classify unlabeled target data accurately [14], [33], [36]. In this line, representative methods include conditional transferable components [14], scatter component analysis [33] and DAN [36].

Another technique is the adversarial training method inspired by the theory of domain adaptation [1]. This theory suggests that predictions must be based on features, and these features cannot be used to discriminate source and target domains [23], [40], [42]. For example, DANN considers two deep networks: one is used to construct new features that predict labels in the TD; while the other is to make two domains non-distinguishing based on these new features [23]. DANN simultaneously trains two deep networks to find domain-invariant representations between two domains.

The last technique is the pseudo-label method, which regards pseudo labels given by a classifier as true labels [13], [51]. The *joint domain adaptation* (JDA) matches joint distributions of two domains using these pseudo labels [51]. The *asymmetric tri-training domain adaptation* (ATDA) leverages three networks asymmetrically [13]. Specifically, two networks are used to annotate unlabeled target data, namely generating pseudo labels. The other network can obtain target-specific representations based on the pseudo-labeled data. Since pseudo-label UDA methods can effectively reduce the upper bound of expected risk in the TD [13], [52], we also consider using the pseudo-label technique to help address the WUDA problem (like ATDA [13]).

## 3 PRELIMINARY

In this section, we introduce notations used in this paper and two common label-noise generation processes [53], [54].

## 3.1 Notations

The following notations are used to demonstrate theoretical results of this paper.

- a space $\mathcal{X} \subset \mathbb{R}^d$ and $\mathcal{Y} = \{1, 2, \ldots, K\}$ as a label set;
- $f_t(x_t)$ and $\tilde{f}_t(x_t)$ represent the ground-truth and pseudo labeling function of the target domain, where $f_t, \tilde{f}_t : \mathcal{X} \to \mathcal{Y}$;
- $A = \{x : \tilde{f}_t(x) \neq f_t(x)\}$ and $B = \mathcal{X}/A$ represent the area where $\tilde{f}_t(x) \neq f_t(x)$ (the set $A$) and the area where $\tilde{f}_t(x) = f_t(x)$ (the set $B$);
- $\tilde{p}_s(x_s, \tilde{y}_s)$, $p_s(x_s, y_s)$ and $q_s(x_s, y_s)$ represent probability densities of noisy, correct and incorrect *multivariate random variable* (m.r.v.) defined on $\mathcal{X} \times \mathcal{Y}$, respectively, and $\tilde{p}_{x_s}(x_s)$, $p_{x_s}(x_s)$ and $q_{x_s}(x_s)$ are their marginal densities on $\mathcal{X}$;
- $p_{x_t}(x_t)$ represents the probability density of m.r.v. $X_t$ defined on $\mathcal{X}$;
- $q_{x_t}(x) = p_{x_t}(x)1_A(x)/P_{x_t}(A)$ represents the probability density of $X_t$ restricted in $A$;
- $p'_{x_t}(x_t) = p_{x_t}(x_t)1_B(x_t)/P_{x_t}(B)$ represents the probability density of $X_t$ restricted in $B$;
- $\mathcal{H}$ is the class of arbitrary *decision functions* $h : \mathcal{X} \to \mathbb{R}^K$;
- $\ell : \mathbb{R}^K \times \mathcal{Y} \to \mathbb{R}_+$ is the loss function. $\ell(t, y)$ means the loss incurred by predicting an output $t$ (e.g., $h(x)$) when the ground truth is $y$;
- $\mathbb{L}_{\mathcal{H}} = \{\ell(h(x), y)|h \in \mathcal{H}, x \in \mathcal{X}, y \in \mathcal{Y}\}$ is the class of loss functions associated with $\mathcal{H}$;
- expected risks on the noisy m.r.v. and correct m.r.v.:

$$\tilde{R}_s(h) = \mathbb{E}_{\tilde{p}_s(x_s, \tilde{y}_s)}[\ell(h(x_s), \tilde{y}_s)],$$
$$R_s(h) = \mathbb{E}_{p_s(x_s, y_s)}[\ell(h(x_s), y_s)];$$

- expected discrepancy (associated with $\ell$) between an arbitrary decision function $h$ and a ground-truth or pseudo labeling function $f$ ($f$ could be $f_t$ or $\tilde{f}_t$) under different marginal densities:

$$\tilde{R}_s(h, f) = \mathbb{E}_{\tilde{p}_{x_s}(x_s)}[\ell(h(x_s), f(x_s))],$$
$$R_s(h, f) = \mathbb{E}_{p_{x_s}(x_s)}[\ell(h(x_s), f(x_s))],$$
$$R_t(h, f) = \mathbb{E}_{p_{x_t}(x_t)}[\ell(h(x_t), f(x_t))].$$

## 3.2 Generating label-noise via the transition matrix

We assume that there are clean source data denoted by a m.r.v. $(X_s, Y_s)$ defined on $\mathcal{X} \times \mathcal{Y}$ with the probability density $p_s(x_s, y_s)$. However, samples of $(X_s, Y_s)$ cannot be directly obtained and we can only observe noisy source data (denoted by m.r.v. $(X_s, \tilde{Y}_s)$) with the probability density $\tilde{p}_s(x_s, \tilde{y}_s)$ [53]. $\tilde{p}_s(x_s, \tilde{y}_s)$ is generated from $p_s(x_s, y_s)$ and a transition matrix $Q_{ij} = \Pr(\tilde{Y}_s = j|Y_s = i)$. Each element in $Q$, $\Pr(\tilde{Y}_s = j|Y_s = i)$, is a transition probability, i.e., the flip rate from a correct label $i$ to a noisy label $j$.

## 3.3 Generating label-noise via the sample selection

The transition matrix $Q$ is easily estimated in certain situations [53]. However, in more complex situations, such as clothing1M dataset [55], noisy source data is directly generated by selecting data from a pool, which mixes correct data (data with correct labels) and incorrect data (data with incorrect labels). Namely, how the correct label $i$ is corrupted to $j$ ($i \neq j$) is unclear.

Let $(X_s, Y_s, V_s)$ be a m.r.v. defined on $\mathcal{X} \times \mathcal{Y} \times \mathcal{V}$ with the probability density $p_s^{\text{po}}(x_s, y_s, v_s)$, where $\mathcal{V} = \{0, 1\}$ is the *perfect-selection random variable*. $V_s = 1$ means "correct" and $V_s = 0$ means "incorrect". Nonetheless, samples of $(X_s, Y_s, V_s)$ cannot be obtained and we can only observe $(X_s, \tilde{Y}_s)$ from a distribution with the following density.

$$\tilde{p}_s(x_s, \tilde{y}_s) = \sum_{v_s=0}^{1} p_{X_s, Y_s|V_s}^{\text{po}}(x_s, y_s|v_s)p_{V_s}^{\text{po}}(v_s), \qquad (1)$$

where $p_{V_s}^{\text{po}}(v_s) = \int_{\mathcal{X}} \sum_{y_s=1}^{K} p_s^{\text{po}}(x_s, y_s, v_s)dx_s$. Eq. (1) means that we lose the information regarding $V_s$. If we uniformly draw samples from $\tilde{p}_s(x_s, \tilde{y}_s)$, the noise rate of these samples is $p_{V_s}^{\text{po}}(0)$. It is clear that the m.r.v. $(X_s, Y_s|V_s = 1)$ is the m.r.v. $(X_s, Y_s)$ mentioned in Section 3.2. Then, $q_s(x_s, y_s)$ is used to describe the density of incorrect m.r.v. $(X_s, Y_s|V_s = 0)$. Using $p_s(x_s, y_s)$ and $q_s(x_s, y_s)$, $\tilde{p}_s(x_s, \tilde{y}_s)$ is expressed by

$$\tilde{p}_s(x_s, \tilde{y}_s) = (1 - \rho)p_s(x_s, y_s) + \rho q_s(x_s, y_s), \qquad (2)$$

where $\rho = p_{V_s}^{\text{po}}(0)$. To reduce noise effects from incorrect data, researchers aim to recover the information of $V_s$, i.e., to select the correct data [54], [56], [57].

## 4 WILDLY UNSUPERVISED DOMAIN ADAPTATION

In this section, we first define a new, more realistic and more challenging problem setting called *wildly unsupervised domain adaptation (WUDA)*, and explain the nature of WUDA. Then, we empirically show that representative UDA methods cannot handle WUDA well, which motivates us to propose a novel method to address the WUDA problem (Section 7).

**Problem 1** (Wildly Unsupervised Domain Adaptation). *Let $X_t$ be a m.r.v. defined on the space $\mathcal{X}$ with respect to the probability density $p_{x_t}$, $(X_s, \tilde{Y}_s)$ be a m.r.v. defined on the space $\mathcal{X} \times \mathcal{Y}$ with respect to the probability density $\tilde{p}_s$, where $\tilde{p}_s$ is the probability density regarding noisy source data (generated in Section 3.2 or 3.3), and $\mathcal{Y} = \{1, \ldots, K\}$ is the label set. Let $p_{x_s}$ be the marginal density of $\tilde{p}_s$. Given i.i.d. data $\tilde{D}_s = \{(x_{si}, \tilde{y}_{si})\}_{i=1}^{n_s}$ and $D_t = \{x_{ti}\}_{i=1}^{n_t}$ drawn from $\tilde{P}_s$ and $P_{x_t}$ separately, in wildly unsupervised domain adaptation, we aim to train with noisy source data $\tilde{D}_s$ and target data $D_t$ to accurately annotate data drawn from $P_{x_t}$, where $p_{x_s} \neq p_{x_t}$.*

**Remark 1.** In Problem 1, $\tilde{D}_s$ is noisy source data, $D_t$ is unlabeled target data, and $\tilde{P}_s$ and $P_{x_t}$ are two probability measures corresponding to densities $\tilde{p}_s(x_s, \tilde{y}_s)$ and $p_{x_t}(x_t)$.

### 4.1 Nature of WUDA

Specifically, there are five distributions involved in WUDA setting: 1) a marginal distribution on source data, i.e., $p_{x_s}$ in Problem 1; 2) a marginal distribution on target data, i.e., $p_{x_t}$ in Problem 1; 3) an incorrect conditional distribution of label given $x_s$, $q(y_s|x_s)$; 4) a correct conditional distribution of label given $x_s$, $p(y_s|x_s)$ and 5) a correct conditional distribution of label given $x_t$, $p(y_t|x_t)$.

Based on Problem 1 and Section 3.3, noisy source data $\tilde{D}_s$ are drawn from $\tilde{p}_s(x_s, y_s) = p_{x_s}(x_s)((1 - \rho)p(y_s|x_s) + \rho q(y_s|x_s))$, where $\rho$ is the noise rate in source data. Namely,
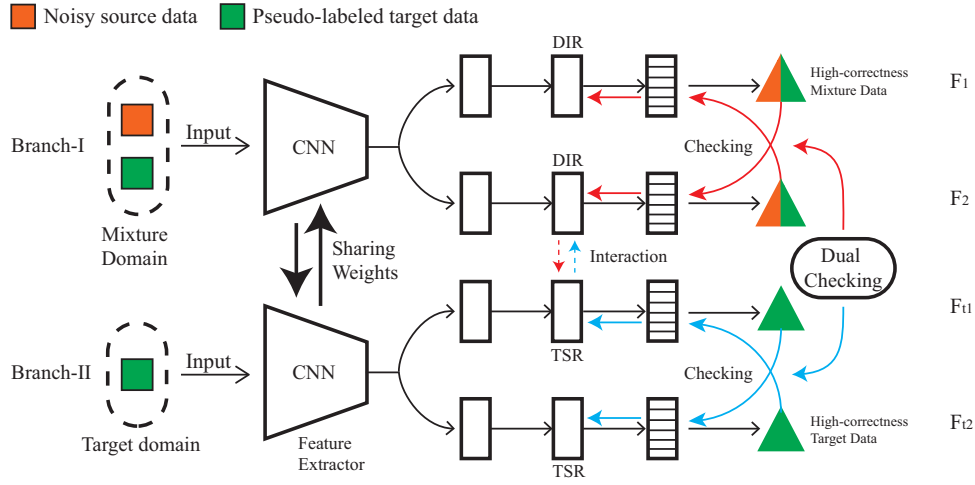
Fig. 3. Butterfly Framework. Two networks ($F_1$ and $F_2$) in Branch-I are jointly trained on noisy source data and pseudo-labeled target data (mixture domain). Two networks in Branch-II ($F_{t1}$ and $F_{t2}$) are trained on pseudo-labeled target data. By using the *dual-checking principle* (DCP), Butterfly checks high-correctness data out from both mixture and pseudo-labeled target data. After cross-propagating checked data, Butterfly can obtain high-quality *domain-invariant representations* (DIR) and *target-specific representations* (TSR) simultaneously in an iterative manner (Algorithms 1 and 2). Note that DIR interacts with TSR via the shared CNN. Besides, in the first training epoch, since we do not have any pseudo-labeled target data, we use noisy source data as the pseudo-labeled target data, which follows [13].

source data $\tilde{D}_s$ are mixture of correct source data from $p_{x_s}(x_s)p(y_s|x_s)$ and incorrect data from $p_{x_s}(x_s)q(y_s|x_s)$. Target data $D_t$ are drawn from $p_{x_t}$. In WUDA setting, we aim to train a classifier with $\tilde{D}_s$ and $D_t$. This classifier is expected to accurately annotate data from $p_{x_t}$, i.e., to accurately simulate distribution 5).

This paper considers WUDA under the common assumption used in the label-noise field, i.e., the $i^{th}$ element in the diagonal of the noise transition matrix is greater than other elements in the $i^{th}$ row or $i^{th}$ column of the noise transition matrix, where $i = 1, 2, \ldots, K$ [53]. Therefore, the proposed approach is able to solve any WUDA problem under the above assumption in principle.

### 4.2 WUDA ruins UDA methods

We take a simple example to illustrate the phenomenon that WUDA ruins representative UDA methods. In Section 5.1, we theoretically analyze the reason of this phenomenon.

We corrupt source data using symmetry flipping [58] and pair flipping [56] that are two representative ways to corrupt true labels. Precise definitions of symmetry flipping ($Q_S$) and pair flipping ($Q_P$) are presented below, where $\rho$ is the noise rate and $K$ is the number of labels.

$$Q_S = \begin{bmatrix} 1-\rho & \frac{\rho}{K-1} & \cdots & \frac{\rho}{K-1} & \frac{\rho}{K-1} \\ \frac{\rho}{K-1} & 1-\rho & \frac{\rho}{K-1} & \cdots & \frac{\rho}{K-1} \\ \vdots & & \ddots & & \vdots \\ \frac{\rho}{K-1} & \cdots & \frac{\rho}{K-1} & 1-\rho & \frac{\rho}{K-1} \\ \frac{\rho}{K-1} & \frac{\rho}{K-1} & \cdots & \frac{\rho}{K-1} & 1-\rho \end{bmatrix}_{K \times K}, \quad (3)$$

$$Q_P = \begin{bmatrix} 1-\rho & \rho & 0 & \cdots & 0 \\ 0 & 1-\rho & \rho & & 0 \\ \vdots & & \ddots & \ddots & \vdots \\ 0 & & & 1-\rho & \rho \\ \rho & 0 & \cdots & 0 & 1-\rho \end{bmatrix}_{K \times K}. \quad (4)$$

For example, if $\rho = 0.4$ and $K = 11$, for the symmetry flipping, the probability that label "0" is corrupted to label "1" is $(1-\rho)/(K-1) = 0.04$. For the pair flipping, the probability that label "0" is corrupted to label "1" is $\rho = 0.4$. To instantiate noisy source data and target data, we leverage *MNIST* and *SYND* (see Figure 4), respectively (i.e., $K = 10$).

We first construct two WUDA tasks with symmetry-flip noise: corrupted *SYND→MNIST* (S→M) and corrupted *MNIST→SYND* (M→S). In Figure 2-(a), we report accuracy of representative UDA methods on unlabeled target data, when the noise rate $\rho$ of SD changes from $5\%$ to $70\%$. It is clear that target-domain accuracy of these representative UDA methods drops quickly when $\rho$ increases. This means that WUDA ruins representative UDA methods. Then, we construct another two WUDA tasks with pair-flip noise. In Figure 2-(b), we report target-domain accuracy, when the noise rate $\rho$ of SD changes from $5\%$ to $45\%$. Again, WUDA still ruins representative UDA methods. Note that, in practice, pair-flip noise is much harder than symmetry-flip noise, the noise rate of pair-flip noise cannot be over $50\%$ [56]. However, the proposed Butterfly network (abbreviated as B-Net, Figure 3) performs robustly when $\rho$ increases (blue lines in Figure 2).

In Section 5, we will analyze the WUDA problem in theory and show why WUDA provably ruins all UDA methods and why the two-step approach is a compromise solution. Then, Section 6 presents how to address the WUDA problem in principle.

## 5 ANALYSIS OF WUDA PROBLEM

In this section, we analyze the WUDA problem from a theoretical view and show the difficulty of the WUDA problem. Completed proofs of lemmas and theorems are demonstrated in the Appendix. In the main content, we provide the main ideas of proving these theoretical results (i.e., *Proof (sketch)*).

## 5.1 WUDA provably ruins UDA methods

Theoretically, we show that existing UDA methods cannot directly transfer useful knowledge from $\tilde{D}_s$ to $D_t$. We first present the relation between $R_s(h)$ and $\tilde{R}_s(h)$.

**Theorem 1.** *If $\tilde{p}_s(x_s, \tilde{y}_s)$ is generated by a transition matrix $Q$ as demonstrated in Section 3.2, we have*

$$\tilde{R}_s(h) = R_s(h) + \mathbb{E}_{p_{x_s}(x_s)}[\boldsymbol{\eta}^T(x_s)(Q - I)\boldsymbol{\ell}(h(x_s))], \quad (5)$$

*where $\boldsymbol{\ell}(h(x_s)) = [\ell(h(x_s), 1), ..., \ell(h(x_s), K)]^T$ and $\boldsymbol{\eta}(x_s) = [p_{Y_s|X_s}(1|x_s), ..., p_{Y_s|X_s}(K|x_s)]^T$. If $\tilde{p}_s(x_s, \tilde{y}_s)$ is generated by sample selection as described in in Section 3.3, we have*

$$\tilde{R}_s(h) = (1 - \rho)R_s(h) + \rho\mathbb{E}_{q_{x_s}(x_s)}[\boldsymbol{\eta}_{\boldsymbol{q}}^T(x_s)\boldsymbol{\ell}(h(x_s))], \quad (6)$$

*where $\boldsymbol{\eta}_{\boldsymbol{q}}(x_s) = [q_{Y_s|X_s}(1|x_s), ..., q_{Y_s|X_s}(K|x_s)]^T$.*

*Proof (sketch).* For Eq. (5), we can prove it using the definition of the transition matrix defined in Section 3.2 and the fact $\tilde{p}_s(x_s, \tilde{y}_s) = \tilde{p}_{\tilde{Y}_s|X_s}(\tilde{y}_s|x_s)p_{x_s}(x_s)$. For Eq. (6), we can prove it using Eq. (2) and the definition of $R_s(h)$. $\quad\square$

**Remark 2.** In Eq. (6), $\mathbb{E}_{q_{x_s}(x_s)}[\boldsymbol{\eta}_{\boldsymbol{q}}^T(x_s)\boldsymbol{\ell}(h(x_s))]$ represents the expected risk on the incorrect m.r.v.. To ensure to obtain useful knowledge from $\tilde{P}_s$, we need to avoid $\tilde{R}_s(h) \approx \mathbb{E}_{q_{x_s}(x_s)}[\boldsymbol{\eta}_{\boldsymbol{q}}^T(x_s)\boldsymbol{\ell}(h(x_s))]$. Specifically, we assume: there is a constant $0 < M_s < \infty$ such that $\mathbb{E}_{q_{x_s}(x_s)}[\boldsymbol{\eta}_{\boldsymbol{q}}^T(x_s)\boldsymbol{\ell}(h(x_s))] \leq R_s(h) + M_s$.

Theorem 1 shows that $\tilde{R}_s(h)$ equals $R_s(h)$ if only two cases happen: 1) $Q = I$ and $\rho = 0$, or 2) some special combinations (e.g., special $p_{x_s}$, $q_{x_s}$, $Q$, $\eta$ and $\ell$) make the second term in Eq. (5) equal zero or make the second term in Eq. (6) equal $\rho R_s(h)$. Case 1) means that source data are clean, which is not real in the wild. Case 2) rarely happens, since it is difficult to find such special combinations when $p_{x_s}$, $q_{x_s}$, $Q$ and $\eta$ are unknown. As a result, $\tilde{R}_s(h)$ has an essential difference with $R_s(h)$. Then, following the proof skills in [1], we derive the upper bound of $R_t(h)$ as below.

**Theorem 2.** *For any $h \in \mathcal{H}$, we have*

$$R_t(h, f_t) \leq \underbrace{\tilde{R}_s(h)}_{(i)\ \textit{noisy-data risk}} + \underbrace{|R_t(h, \tilde{f}_t) - \tilde{R}_s(h, \tilde{f}_t)|}_{(ii)\ \textit{discrepancy between distributions}}$$
$$+ \underbrace{|R_s(h, \tilde{f}_t) - R_s(h)|}_{(iii)\ \textit{domain dissimilarity}}$$
$$+ \underbrace{|\tilde{R}_s(h) - R_s(h)| + |\tilde{R}_s(h, \tilde{f}_t) - R_s(h, \tilde{f}_t)|}_{(iv)\ \textit{noise effects from source } \Delta_s}$$
$$+ \underbrace{|R_t(h, f_t) - R_t(h, \tilde{f}_t)|}_{(v)\ \textit{noise effects from target } \Delta_t}. \quad (7)$$

*Proof (sketch).* For any $h \in \mathcal{H}$, we have

$$R_t(h, f_t)$$
$$= R_t(h, f_t) + \tilde{R}_s(h) - \tilde{R}_s(h) + R_s(h, f_t) - R_s(h, f_t)$$
$$= \tilde{R}_s(h) + R_t(h, f_t) - \tilde{R}_s(h, f_t) + R_s(h, f_t) - R_s(h)$$
$$+ R_s(h) - \tilde{R}_s(h) + \tilde{R}_s(h, f_t) - R_s(h, f_t).$$

Since we do not know $f_t$, we substitute the following equations into the above equation,

$$R_t(h, f_t) = R_t(h, \tilde{f}_t) + R_t(h, f_t) - R_t(h, \tilde{f}_t),$$
$$\tilde{R}_s(h, f_t) = \tilde{R}_s(h, \tilde{f}_t) + \tilde{R}_s(h, f_t) - \tilde{R}_s(h, \tilde{f}_t),$$
$$R_s(h, f_t) = R_s(h, \tilde{f}_t) + R_s(h, f_t) - R_s(h, \tilde{f}_t),$$

which proves this theorem. $\quad\square$

**Remark 3.** To ensure that we can gain useful knowledge from $\tilde{f}_t(x_t)$, we assume: there is a constant $0 < M_t < \infty$ such that $\mathbb{E}_{q_{x_s}(x)}[\ell(h(x), \tilde{f}_t(x))] \leq R_s(h, \tilde{f}_t) + M_t$ and $\mathbb{E}_{q_{x_t}(x)}[\ell(h(x), \tilde{f}_t(x))] \leq R_t(h, f_t) + M_t$. Since we do not have labels in the target domain, we also assume that there exists an $h \in \mathcal{H}$ such that $R_t(h, f_t) + R_s(h)$ is a small value. This assumption follows common assumption of UDA problem [1] and ensures that the adaptation is possible.

It is clear that the upper bound of $R_t(h, f_t)$, shown in Eq. (7), has 5 terms. However, existing UDA methods only focus on minimizing $(i) + (ii)$ [23], [33], [36] or $(i) + (ii) + (iii)$ [13], which ignores terms $(iv)$ and $(v)$ (i.e., $\Delta = \Delta_s + \Delta_t$). Thus, existing UDA methods cannot handle WUDA well.

## 5.2 Two-step approach is a compromise solution

To reduce noise effects, a straightforward solution is two-step approach. For example, in the first step, we can train a classifier with noisy source data using co-teaching [56] and use this classifier to annotate pseudo labels for source data. In the second step, we use ATDA [13] to train a target-domain classifier with pseudo-labeled-source data and pseudo-labeled target data.

Nonetheless, the pseudo-labeled source data are still noisy. Let labels of noisy source data $\tilde{y}_s$ be replaced with pseudo labels $\tilde{y}_s'$ after using co-teaching. Noise effects $\Delta$ will become pseudo-label effects $\Delta_p$ as follows.

$$\Delta_p = \underbrace{|\tilde{R}_s'(h) - R_s(h)| + |\tilde{R}_s'(h, \tilde{f}_t) - R_s(h, \tilde{f}_t)|}_{\textbf{pseudo-labeled-source effects } \Delta_s'} + \Delta_t, \quad (8)$$

where $\tilde{R}_s'(h)$ and $\tilde{R}_s'(h, \tilde{f}_t)$ correspond to $\tilde{R}_s(h)$ and $\tilde{R}_s(h, \tilde{f}_t)$ in $\Delta_s$. It is clear that the difference between $\Delta_p$ and $\Delta$ is $\Delta_s' - \Delta_s$. The left term in $\Delta_s'$ may be less than that in $\Delta_s$ due to a label-noise algorithm (e.g., co-teaching [56]), but the right term in $\Delta_s'$ may be higher than that in $\Delta_s$ since a label-noise algorithm does not consider minimizing it. Thus, it is hard to say whether $\Delta_s' < \Delta_s$ (i.e., $\Delta_p < \Delta$). This means that two-step approach may not really reduce noise effects.

## 6 HOW TO ADDRESS WUDA IN PRINCIPLE

To eliminate noise effects $\Delta$, we aim to select correct data simultaneously from noisy source data and pseudo-labeled target data. In theory, we prove that noise effects will be eliminated if we can select correct data with a high probability. Let $\rho_{01}^s$ represent the probability that incorrect data is selected from noisy source data, and $\rho_{01}^t$ represent the probability that incorrect data is selected from pseudo-labeled target data. Theorem 3 shows that $\Delta \to 0$ if $\rho_{01}^s \to 0$ and $\rho_{01}^t \to 0$ and presents a new upper bound of $R_t(h, f_t)$. Before stating Theorem 3, we first present two m.r.v.s below.

- $(X_s, Y_s, V_s)$ defined on $\mathcal{X} \times \mathcal{Y} \times \mathcal{V}$ with the probability density $p_s^{\text{po}}(x_s, y_s, v_s)$, where $\mathcal{V} = \{0, 1\}$;
- $(X_t, V_t)$ defined on $\mathcal{X} \times \mathcal{V}$ with the probability density $p_t^{\text{po}}(x_t, v_t)$, where $\mathcal{V} = \{0, 1\}$. $p_{V_t}^{\text{po}}(v_t)$ is the marginal density of $p_t^{\text{po}}(x_t, v_t)$.

The $V_s$ has been introduced in Section 3.3. Similar with $V_s$, $V_t$ is also a *perfect-selection random variable*. Data drawn from the distribution of $(X_t, V_t)$ can be regarded as a pool that mixes the correct ($v_t = 1$) and incorrect ($v_t = 0$) pseudo-labeled target data. Namely, $V_t = 1$ means $f_t(x_t) = \tilde{f}_t(x_t)$ and $V_t = 0$ means $f_t(x_t) \neq \tilde{f}_t(x_t)$. It is clear that, higher value of $p_{V_t}^{\text{po}}(V_t = 1)$ means that $\tilde{f}_t$ is more like $f_t$. In following, we use $\rho_{v_t}$ to represent $p_{V_t}^{\text{po}}(v_t = 0)$. Note that both perfect-selection random variables $V_s$ and $V_t$ cannot be observed and we can only observe following m.r.v.s.

- $(X_s, Y_s, U_s)$ defined on $\mathcal{X} \times \mathcal{Y} \times \mathcal{V}$ with the probability density $\tilde{p}_s^{\text{po}}(x_s, y_s, u_s)$;
- $(X_t, U_t)$ defined on $\mathcal{X} \times \mathcal{V}$ with the probability density $\tilde{p}_t^{\text{po}}(x_t, u_t)$. $\tilde{p}_{U_t}^{\text{po}}(u_t)$ is the marginal density of $\tilde{p}_t^{\text{po}}(x_t, u_t)$.

The $U_s$ and $U_t$ are *algorithm-selection random variables*. Data drawn from the distribution of $(X_s, Y_s, U_s)$ can be regarded as a pool that mixes the selected ($u_s = 1$) and unselected ($u_s = 0$) noisy source data. Data drawn from the distribution of $(X_t, U_t)$ can be regarded as a pool that mixes the selected ($u_t = 1$) and unselected ($u_t = 0$) pseudo-labeled target data. We can obtain observations of $(X_s, Y_s, U_s)$ and $(X_t, U_t)$ using an algorithm that is used to select correct data. After executing the algorithm, we can obtain observations $\{x_{si}, \tilde{y}_{si}, u_{si}\}_{i=1}^{n_s}$ and $\{x_{ti}, u_{ti}\}_{i=1}^{n_t}$. Based on $(X_s, Y_s, U_s)$ and $(X_t, U_t)$, we can define the following expected risks.

$$\tilde{R}_s^{\text{po}}(h, u_s) = (1 - \rho_{u_s})^{-1} \mathbb{E}_{\tilde{p}_s^{\text{po}}(x_s, y_s, u_s)}[u_s \ell(h(x_s), y_s)],$$

$$\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) = (1 - \rho_{u_t})^{-1} \mathbb{E}_{\tilde{p}_t^{\text{po}}(x_t, u_t)}[u_t \ell(h(x_t), \tilde{f}_t(x_t))],$$

$$\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s) = (1 - \rho_{u_s})^{-1} \mathbb{E}_{\tilde{p}_s^{\text{po}}(x_s, y_s, u_s)}[u_s \ell(h(x_s), \tilde{f}_t(x_s))].$$

where $\rho_{u_s} = \tilde{p}_{U_s}^{\text{po}}(u_s = 0)$ and $\rho_{u_t} = \tilde{p}_{U_t}^{\text{po}}(u_t = 0)$. Since we can observe $(X_s, Y_s, U_s)$ and $(X_t, U_t)$, the empirical estimators of these three risks can be easily computed. Then, we define following probabilities to describe the relation between perfect-selection random variables and algorithm-selection random variables, where $i, j = 0, 1$.

- $\rho_{ji}^s = \text{Pr}(V_s = j | U_s = i)$ represents the probability of the event: $V_s = j$ given $U_s = i$,
- $\rho_{ji}^t = \text{Pr}(V_t = j | U_t = i)$ represents the probability of the event: $V_t = j$ given $U_t = i$.

**Remark 4.** Based on above definitions, we know that 1) $\rho_{01}^s$ is the probability that incorrect data is selected from noisy source data, and 2) $\rho_{01}^t$ is the probability that incorrect data is selected from pseudo-labeled target data.

Using $\rho_{ji}^s$ and $\rho_{ji}^t$, we can show the relation between probability densities of $(X_s, Y_s | V_s)$ and $(X_s, Y_s | U_s)$, and the relation between probability densities of $(X_t | V_t)$ as follows.

$$\tilde{p}_{X_s, Y_s | U_s}^{\text{po}}(x_s, y_s | i) = \rho_{0i}^s p_{X_s, Y_s | V_s}^{\text{po}}(x_s, y_s | 0) + \rho_{1i}^s p_{X_s, Y_s | V_s}^{\text{po}}(x_s, y_s | 1),$$

$$\tilde{p}_{X_t | U_t}^{\text{po}}(x_t | i) = \rho_{0i}^t p_{X_t | V_t}^{\text{po}}(x_t | 0) + \rho_{1i}^t p_{X_t | V_t}^{\text{po}}(x_t | 1).$$

Since

$$p_{X_s, Y_s | V_s}^{\text{po}}(x_s, y_s | 1) = p_s(x_s, y_s),$$

$$p_{X_s, Y_s | V_s}^{\text{po}}(x_s, y_s | 0) = q_s(x_s, y_s),$$

$$p_{X_t | V_t}^{\text{po}}(x_t | 0) = p_{x_t}(x_t) 1_A(x_t) / P_{x_t}(A) = q_{x_t}(x_t),$$

$$p_{X_t | V_t}^{\text{po}}(x_t | 1) = p_{x_t}(x_t) 1_B(x_t) / P_{x_t}(B) = p'_{x_t}(x_t),$$

we have

$$\tilde{p}_{X_s, Y_s | U_s}^{\text{po}}(x_s, y_s | i) = \rho_{0i}^s q_s(x_s, y_s) + \rho_{1i}^s p_s(x_s, y_s), \quad (9)$$

$$\tilde{p}_{X_t | U_t}^{\text{po}}(x_t | i) = \rho_{0i}^t q_{x_t}(x_t) + \rho_{1i}^t p'_{x_t}(x_t). \quad (10)$$

**Remark 5.** Eq. (9) and Eq. (10) show that, if $\rho_{01}^s \to 0$ and $\rho_{01}^t \to 0$, we have 1) $\tilde{p}_{X_s, Y_s | U_s}^{\text{po}}(x_s, y_s | 1) \to p_s(x_s, y_s)$ and 2) $\tilde{p}_{X_t | U_t}^{\text{po}}(x_t | 1) \to p'_{x_t}(x_t)$. However, we cannot prove the main theorem (Theorem 3) using 1) and 2), since we only take care of risks instead of densities (like 1) and 2)).

Next, we present a lemma to show the relation between $\tilde{R}_s^{\text{po}}(h, u_s)$ and $R_s(h)$.

**Lemma 1.** *Given the m.r.v.* $(X_s, Y_s, U_s)$ *with the probability density* $\tilde{p}_s^{\text{po}}(x_s, y_s, u_s)$ *and Eq. (9), we have*

$$|\tilde{R}_s^{\text{po}}(h, u_s) - R_s(h)|$$
$$\leq \rho_{01}^s \max\{\mathbb{E}_{q_s(x_s, y_s)}[\ell(h(x_s), y_s)], R_s(h)\}. \quad (11)$$

*Proof (sketch).* Based on definition of $\tilde{R}_s^{\text{po}}(h, u_s)$ and the fact $\tilde{p}_s^{\text{po}}(x_s, y_s, u_s) = \tilde{p}_{X_s, Y_s | U_s}^{\text{po}}(x_s, y_s | 1) \tilde{p}_{U_s}^{\text{po}}(1)$, $\tilde{R}_s^{\text{po}}(h, u_s)$ equals

$$\frac{\int_{\mathcal{X}} \sum_{y_s=1}^{K} \ell(h(x_s), y_s) \tilde{p}_{X_s, Y_s | U_s}^{\text{po}}(x_s, y_s | 1) \tilde{p}_{U_s}^{\text{po}}(1) dx_s}{1 - \rho_{u_s}}$$

Then, we can use the definition of $\rho_{u_s}$ and the Eq. (9) to prove this lemma. $\qquad \square$

Similar with Lemma 1, we can obtain

$$|\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s) - R_s(h, \tilde{f}_t)|$$
$$\leq \rho_{01}^s \max\{\mathbb{E}_{q_{x_s}(x_s)}[\ell(h(x_s), \tilde{f}_t(x_s))], R_s(h, \tilde{f}_t)\}. \quad (12)$$

Then, we give another lemma to show relation between $\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t)$ and $R_t(h, \tilde{f}_t)$.

**Lemma 2.** *Given the m.r.v.* $(X_t, U_t)$ *with the probability density* $\tilde{p}_s^{\text{po}}(x_t, u_t)$ *and Eq. (10), if* $\mathbb{E}_{p'_{x_t}(x_t)}[\ell(h(x_t), f_t(x_t))] \leq R_t(h, f_t) + \rho_{01}^s M_t$, *then we have*

$$|\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) - R_t(h, f_t)|$$
$$\leq \rho_{01}^t \max\{\mathbb{E}_{q_{x_t}(x_t)}[\ell(h(x_t), \tilde{f}_t(x_t))], R_t(h, f_t)\} + \rho_{11}^t \rho_{01}^s M_t. \quad (13)$$

*Proof (sketch).* According to definition of $\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t)$, we can unfold it to be

$$\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t)$$
$$= (1 - \rho_{u_t})^{-1} \int_{\mathcal{X}} \ell(h(x_t), \tilde{f}_t(x_t)) \tilde{p}_{X_t | U_t}^{\text{po}}(x_t | 1) \tilde{p}_{U_t}^{\text{po}}(1) dx_t.$$

Then, using the definition of $\rho_{u_s}$, Eq. (9), the definition of $V_t$ ($f_t(x_t) = \tilde{f}_t(x_t)$ when $V_t = 1$) and the assumption that $\mathbb{E}_{p'_{x_t}(x_t)}[\ell(h(x_t), f_t(x_t))] \leq R_t(h, f_t) + \rho_{01}^s M_t$, we have

$$\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t)$$
$$\leq \rho_{01}^t \mathbb{E}_{q_{x_t}(x_t)}[\ell(h(x_t), \tilde{f}_t(x_t))] + \rho_{11}^t(R_t(h, f_t) + \rho_{01}^s M_t).$$

Finally, we can upper bound $|\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) - R_t(h, f_t)|$ using the above inequality, which proves this lemma.    □

**Remark 6.** In Lemma 2, $\mathbb{E}_{p'_{x_t}(x_t)}[\ell(h(x_t), f_t(x_t))] \leq R_t(h, f_t) + \rho_{01}^s M_t$ means that the expected risk restricted in $B$ (i.e., $\mathbb{E}_{p'_{x_t}(x_t)}[\ell(h(x_t), f_t(x_t))]$) can represent the true risk $R_t(h, f_t)$ when $\rho_{01}^s$ is small. If this assumption fails, we cannot gain useful knowledge from $\tilde{f}_t$ even when we can select correct data from pseudo-labeled target data ($\rho_{01}^t = 0$).

Inequalities (11), (12) and (13) show that if we can perfectly avoid annotating incorrect data as "correct" (i.e., $\rho_{01}^s = 0$ and $\rho_{01}^t = 0$), we have $\tilde{R}_s^{\text{po}}(h, u_s) = R_s(h)$, $\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_t) = R_s(h, \tilde{f}_t)$ and $\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) = R_t(h, f_t)$. Nonetheless, $\rho_{01}^s$ and $\rho_{01}^t$ never equal zero, and $\mathbb{E}_{q_s(x_s, y_s)}[\ell(h(x), y)]$, $\mathbb{E}_{q_{x_s}(x_s)}[\ell(h(x_s), \tilde{f}_t(x_s))]$ and $\mathbb{E}_{q_{x_t}(x_t)}[\ell(h(x_t), \tilde{f}_t(x_t))]$ may equal $+\infty$ for some $h \in \mathcal{H}$. Namely, even when $\rho_{01}^s$ and $\rho_{01}^t$ are very small, $\tilde{R}_s^{\text{po}}(h, u_s)$ is probably far away from $R_s(h)$. Thus, without proper assumptions, it is useless to use $(X_s, Y_s, U_s)$ to represent $(X_s, Y_s | V_s = 1)$.

In Theorem 3, we prove that, under assumptions in Remarks 2, 3 and Lemma 2, $\tilde{R}_s^{\text{po}}(h, u_s) \to R_s(h)$, $\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_t) \to R_s(h, \tilde{f}_t)$ and $\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) \to R_t(h, f_t)$ if $\rho_{01}^s \to 0$ and $\rho_{01}^t \to 0$. Moreover, we give a new upper bound of $R_t(h, f_t)$. In the new upper bound, we show that: $\Delta \to 0$ if $\rho_{01}^s \to 0$ and $\rho_{01}^t \to 0$.

**Theorem 3.** *Given two m.r.v.s $(X_s, Y_s, U_s)$ defined on $\mathcal{X} \times \mathcal{Y} \times \mathcal{V}$ and $(X_t, U_t)$ defined on $\mathcal{X} \times \mathcal{V}$, under the assumptions in Remark 2, Remark 3 and Lemma 2, $\forall \epsilon \in (0, 1)$, there are $\delta_s$ and $\delta_t$, if $\rho_{01}^s < \delta_s$ and $\rho_{01}^t < \delta_t$, for any $h \in \mathcal{H}$, we will have*

$$|\tilde{R}_s^{po}(h, \tilde{f}_t, u_s) - R_s(h, \tilde{f}_t)| + |\tilde{R}_s^{po}(h, u_s) - R_s(h)| < 2\epsilon. \quad (14)$$

*Moreover, we will have*

$$R_t(h, f_t) \leq \underbrace{\tilde{R}_s^{po}(h, u_s)}_{(i) \text{ noisy-data risk}} + \underbrace{|\tilde{R}_t^{po}(h, \tilde{f}_t, u_t) - \tilde{R}_s^{po}(h, \tilde{f}_t, u_s)|}_{(ii) \text{ discrepancy between distributions}}$$
$$+ \underbrace{|R_s(h, \tilde{f}_t) - R_s(h)|}_{(iii) \text{ domain dissimilarity}} + \underbrace{2\epsilon}_{(iv) \text{ noise effects } \Delta_s}$$
$$+ \underbrace{2\epsilon}_{(iv) \text{ noise effects } \Delta_t} . \quad (15)$$

*Proof.* We first prove upper bounds of $|\tilde{R}_s^{\text{po}}(h, u_s) - R_s(h)|$, $|\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_t) - R_s(h, \tilde{f}_t)|$ and $|\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) - R_t(h, f_t)|$ under assumptions in Theorem 3. Based on Lemma 1,

$$|\tilde{R}_s^{\text{po}}(h, u_s) - R_s(h)|$$
$$= |\rho_{01}^s \mathbb{E}_{q_s(x_s, y_s)}[\ell(h(x_s), y_s)] - (1 - \rho_{11}^s)R_s(h)|$$
$$\leq |\rho_{01}^s(R_s(h) + M_s) - \rho_{01}^s R_s(h)|$$
$$= \rho_{01}^s M_s. \quad (16)$$

Similar, we have

$$|\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s) - R_s(h, \tilde{f}_t)| \leq \rho_{01}^t M_t, \quad (17)$$

$$|\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) - R_t(h, f_t)| \leq \rho_{01}^t M_t + \rho_{11}^t \rho_{01}^s M_t. \quad (18)$$

Since $M_s$ and $M_t$ are positive constants, it is clear that $\tilde{R}_s^{\text{po}}(h, u_s) \to R_s(h)$, $\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s) \to R_s(h, \tilde{f}_t)$ and $\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) \to R_t(h, f_t)$ when $\rho_{01}^s \to 0$ and $\rho_{01}^t \to 0$.

Specifically, $\forall \epsilon \in (0, 1)$, let $\delta_t = \epsilon/M_t$ and $\delta_s = \epsilon/\max\{M_s, \rho_{11}^t M_t\}$. When $\rho_{01}^s < \delta_s$ and $\rho_{01}^t < \delta_t$, we have

$$|\tilde{R}_s^{\text{po}}(h, u_s) - R_s(h)| + |\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s) - R_s(h, \tilde{f}_t)| < 2\epsilon \quad (19)$$

$$|\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) - R_t(h, f_t)| < 2\epsilon. \quad (20)$$

Hence, we prove the Eq. (14). In following, we give a new upper bound of $R_t(h, f_t)$. Recall Theorem 2, we replace 1) $\tilde{R}_s(h)$ with $\tilde{R}_s^{\text{po}}(h, u_s)$, 2) $\tilde{R}_s(h, \tilde{f}_t)$ with $\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s)$, 3) $R_t(h, \tilde{f}_t)$ with $\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t)$. Then, we have

$$R_t(h, f_t) \leq \tilde{R}_s^{\text{po}}(h, u_s) + |\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) - \tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s)|$$
$$+ |R_s(h, \tilde{f}_t) - R_s(h)| + |\tilde{R}_s^{\text{po}}(h, u_s) - R_s(h)|$$
$$+ |\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s) - R_s(h, \tilde{f}_t)|$$
$$+ |R_t(h, f_t) - \tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t)|. \quad (21)$$

Let $\rho_{01}^s \leq \delta_s$ and $\rho_{01}^t \leq \delta_t$, based on Eqs. (19) and (20), we have

$$R_t(h, f_t) \leq \underbrace{\tilde{R}_s^{\text{po}}(h, u_s)}_{(i) \text{ noisy-data risk}} + \underbrace{|\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t) - \tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s)|}_{(ii) \text{ discrepancy between distributions}}$$
$$+ \underbrace{|R_s(h, \tilde{f}_t) - R_s(h)|}_{(iii) \text{ domain dissimilarity}} + \underbrace{2\epsilon}_{(iv) \text{ noise effects } \Delta_s}$$
$$+ \underbrace{2\epsilon}_{(iv) \text{ noise effects } \Delta_t} .$$

Hence, we prove this theorem.    □

Theorem 3 shows that if selected data have a high probability to be correct ones ($\rho_{01}^s \to 0$ and $\rho_{01}^t \to 0$), then $\Delta_s$ and $\Delta_t$ approach zero, meaning that noise effects are eliminated. This motivates us to find a reliable way to select correct data from noisy source data and pseudo-labeled target data and propose the butterfly to WUDA problem.

**Remark 7.** Note that, since Theorems 2 and 3 hold for any hypothesis and any data distributions, the bounds in both theorems are loose and pessimistic. However, both theorems are proposed to show which factors we should take care of in the WUDA problem and both theorems point out the major difference between WUDA and UDA. From this perspective, both theorems are very important for positioning and understanding the WUDA problem.

# 7 BUTTERFLY: TOWARDS ROBUST ONE-STEP APPROACH

This section presents Butterfly to solve the WUDA problem.

## 7.1 What is the Principle-guided Solution?

Guided by Theorem 3, a robust approach should check high-correctness data out (meaning $\rho_{01}^s \to 0$ and $\rho_{01}^t \to 0$). This checking process will make $(iv)$ and $(v)$, $2\epsilon + 2\epsilon$, become 0. Then, we can obtain gradients of $\tilde{R}_s^{\text{po}}(h, u_s)$, $\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s)$ and $\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t)$ w.r.t. parameters of $h$ and use these gradients to minimize them, which minimizes $(i)$ and $(ii)$ as $(i) +$

$(ii) \leq \tilde{R}_s^{\text{po}}(h, u_s) + \tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s) + \tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t)$. Note that $(iii)$ cannot be directly minimized since we cannot pinpoint clean source data. However, following [13], we can indirectly minimize $(iii)$ via minimizing $\tilde{R}_s^{\text{po}}(h, u_s) + \tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s)$, as $(iii) \leq R_s(h, \tilde{f}_t) + R_s(h) \leq \tilde{R}_s^{\text{po}}(h, u_s) + \tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s) + 2\epsilon$, where the last inequality follows Eq. (14). This means that a robust approach guided by Theorem 3 can minimize all terms in the right side of inequality in Eq. (15).

## 7.2 Dual-checking principle

**Memorization effects of deep networks.** Recently, an interesting observation for deep networks is that they can memorize easy samples first, and gradually adapt to hard samples as increasing training epochs [59]. Namely, although deep networks can fit everything (e.g., mislabeled data) in the end, they *learn patterns first* [59]: this suggests deep networks can gradually memorize the data, moving from regular data to irregular data such as outliers. To utilize this memorization effects, previous studies have shown that we can regard small-loss data as correct ones (also known as the *small-loss trick*). Then we can obtain a good classifier that is trained with the small-loss data [54].

**Co-teaching learning paradigm.** However, if we only use small-loss trick to select correct data (like [54]), we will get accumulated errors caused by sample-selection bias [56]. Therefore, researchers also consider a new deep learning paradigm called *co-teaching*, where we train two deep networks simultaneously, and let them *teach each other* [56]. Based on this novel learning paradigm, we can effectively reduce the negative effects from the accumulated errors caused by sample-selection bias.

**Dual-checking principle.** Motivated by Section 7.1, we propose the *dual-checking principle* (DCP): we need to check high-correctness data out in the source and target domains simultaneously. According to the memorization effects of deep networks, we realize DCP based on deep networks, small-loss trick and the co-teaching learning paradigm (i.e., the Butterfly introduced below).

## 7.3 Principle-guided Butterfly

To realize the robust approach for addressing the WUDA problem, we propose a Butterfly framework, which trains four networks dividing into two branches (Figure 3). By using DCP, Branch-I checks which data is correct in the mixture domain; while Branch-II checks which pseudo-labeled target data is correct. To ensure these checked data highly-correct, we apply the small-loss trick based on memorization effects of deep learning [59]. After cross-propagating these checked data [60], Butterfly can obtain high-quality DIR and TSR simultaneously in an iterative manner. Theoretically, Branch-I minimizes $(i)+(ii)+(iii)+(iv)$; while Branch-II minimizes $(ii) + (v)$. This means that Butterfly can minimize all terms in the right side of inequality in Eq. (15).

## 7.4 Loss function in Butterfly

According to $\tilde{R}_s^{\text{po}}(h, u_s)$, $\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, u_t)$ and $\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, u_s)$ defined in Section 6, four networks trained by Butterfly share the same loss function but with different inputs.

$$\mathcal{L}(\theta, \boldsymbol{u}; F, D) = \frac{1}{\sum_{i=1}^n u_i} \sum_{i=1}^n u_i \ell(F(x_i), \check{y}_i), \quad (22)$$

where $n$ is the batch size (i.e., $n = |D|$), and $F$ represents a network (e.g., $F_1, F_2, F_{t1}$ and $F_{t2}$). $D = \{(x_i, \check{y}_i)\}_{i=1}^n$ is a mini-batch for training a network, where $\{x_i, \check{y}_i\}_{i=1}^n$ could be data in MD or TD (Figure 3), and $\theta$ represents parameters of $F$ and $\boldsymbol{u} = [u_1, ..., u_n]^T$ is an $n$-by-1 vector whose elements equal 0 or 1. For two networks in Branch-I, following [13], we also add a regularizer $|\theta_{f11}^T \theta_{f21}|$ in their loss functions, where $\theta_{f11}$ and $\theta_{f21}$ are weights of the first fully-connect layer of $F_1$ and $F_2$. With this regularizer, $F_1$ and $F_2$ will learn from different features.

**Nature of the loss $\mathcal{L}$.** In the loss function $\mathcal{L}$, we have $n$ samples: $\{(x_i, \check{y}_i)\}_{i=1}^n$. For the $i^{th}$ sample, we will compute its cross-entropy loss (i.e., $\ell(F(x_i), \check{y}_i)$), and we will denote this sample as "selected" if $u_i = 1$. Thus, the nature of $\mathcal{L}$ is actually the average value of cross-entropy loss of these "selected" samples. Note that, we need to set a constrain to prevent $\sum_{i=1}^n u_i = 0$ in $\mathcal{L}$, which means that we should select at least one sample to compute $\mathcal{L}$.

## 7.5 Training procedures of Butterfly

This subsection will first present the checking process in Butterfly (Algorithm 1). Then, the full training procedure of Butterfly (Algorithm 2) will be introduced in detail.

### 7.5.1 Checking process in Butterfly (Algorithm 1)

We first obtain four inputs: 1) networks $F_1$ and $F_2$, and 2) a mini-batch $D$, and 3) learning rate $\eta$ and 4) remember rate $\alpha$ (line 1). Then, we will obtain the best $\boldsymbol{u}_1$ by solving a minimization problem (line 2). $\mathcal{L}$ represents the loss function defined in Eq. (22). $\theta_1$ represents the parameters of the network $F_1$. Similarly we will obtain the best $\boldsymbol{u}_2$ (line 3). $\theta_2$ represents the parameters of the network $F_2$. Next, $\theta_1$ and $\theta_2$ are updated using gradient descent, where the gradients are computed using a given optimizer (lines 4-5). Finally, we substitute the updated $\theta_1$ into $F_1$ and the updated $\theta_2$ into $F_2$ and output $F_1$ and $F_2$ (line 6). Note that, lines 2-3 correspond to the small-loss trick mentioned in Section 7.2, and lines 4-5 corresponds to the co-teaching paradigm in Section 7.2.

**Remark 8.** In line 2 or 3 in Algorithm 1, we need to solve a minimization problem: $\min_{\boldsymbol{u}': \mathbf{1}\boldsymbol{u}' > \alpha|D|} \mathcal{L}(\theta, \boldsymbol{u}'; F, D)$ and return the best $\boldsymbol{u}'$ as $\boldsymbol{u}$ ($\boldsymbol{u}_1$ in line 2 and $\boldsymbol{u}_2$ in line 3). In this paragraph, we will show how to quickly solve this problem using a sorting algorithm. Recall the nature of the loss $\mathcal{L}$, we know $\mathcal{L}$ is the average value of cross-entropy losses of "selected" samples, and $\mathbf{1}\boldsymbol{u}'$ is the number of these "selected" samples. Therefore, this minimization problem is equivalent to "given a fixed $F$ ($F_1$ or $F_2$) and $n$ samples in $D$, how to select at least $k$ samples such that $\mathcal{L}$ is minimized", where $k = \lceil \alpha|D| \rceil$. To solve this problem, we first use a sorting algorithm (top_k function in TensorFlow) to sort these $n$ samples according to their cross-entropy losses $\ell(F_1(x_i), \check{y}_i)$. Then, we select $k$ samples with the smallest cross-entropy losses. Finally, let $u_i$ of these $k$ samples be 1 and $u_i$ of the other samples be 0, and we can get the best $\boldsymbol{u} = [u_1, \ldots, u_n]$. The average value of cross-entropy losses of these $k$ samples

---

**Algorithm 1** Checking($F_1, F_2, D, \eta, \alpha$)

---

**1: Input** networks $F_1$, $F_2$, mini-batch $D$, learning rate $\eta$, remember rate $\alpha$;

**2: Obtain** $\boldsymbol{u}_1 = \arg\min_{\boldsymbol{u}_1': \mathbf{1}\boldsymbol{u}_1' > \alpha|D|} \mathcal{L}(\theta_1, \boldsymbol{u}_1'; F_1, D)$;　　　　　　　　　// Check high-correctness data

**3: Obtain** $\boldsymbol{u}_2 = \arg\min_{\boldsymbol{u}_2': \mathbf{1}\boldsymbol{u}_2' > \alpha|D|} \mathcal{L}(\theta_2, \boldsymbol{u}_2'; F_2, D)$;　　　　　　　　　// Check high-correctness data

**4: Update** $\theta_1 = \theta_1 - \eta\nabla\mathcal{L}(\theta_1, \boldsymbol{u}_2; F_1, D)$;　　　　　　　　　　　　　　// Update $\theta_1$

**5: Update** $\theta_2 = \theta_2 - \eta\nabla\mathcal{L}(\theta_2, \boldsymbol{u}_1; F_2, D)$;　　　　　　　　　　　　　　// Update $\theta_2$

**6: Output** $F_1$ and $F_2$

---

**Algorithm 2** Butterfly Framework: quadruple training for WUDA problem

---

**1: Input** $\tilde{D}_s$, $D_t$, learning rate $\eta$, fixed $\tau$, fixed $\tau_t$, epoch $T_k$ and $T_{max}$, iteration $N_{max}$, # of pseudo-labeled target data $n_{init}$, max of $n_{init}$ $n_{t,max}^l$;

**2: Initial** $F_1, F_2, F_{t1}, F_{t2}, \tilde{D}_t^l = \tilde{D}_s, \tilde{D} = \tilde{D}_s, n_t^l = n_{init}$;

**for** $T = 1, 2, \ldots, T_{max}$ **do**

　　**3: Shuffle** training set $\tilde{D}$;　　　　　　　　　　　　　　　// Noisy dataset

　　**for** $N = 1, \ldots, N_{max}$ **do**

　　　　**4: Fetch** mini-batch $\check{D}$ from $\tilde{D}$;

　　　　**5: Update** Branch-I: $F_1, F_2$ = Checking($F_1, F_2, \check{D}, \eta, R(T)$);　　　// Check data in MD using Algorithm 1

　　　　**6: Fetch** mini-batch $\check{D}_t$ from $\tilde{D}_t^l$;

　　　　**7: Update** Branch-II: $F_{t1}, F_{t2}$ = Checking($F_{t1}, F_{t2}, \check{D}_t, \eta, R_t(T)$);　// Check data in TD using Algorithm 1

　　**end**

　　**8: Obtain** $\tilde{D}_t^l$ = Labeling($F_1, F_2, D_t, n_t^l$);　　　　　　　　　// Label $D_t$, following [13]

　　**9: Obtain** $\tilde{D} = \tilde{D}_s \cup \tilde{D}_t^l$;　　　　　　　　　　　　　　// Update MD

　　**10: Update** $n_t^l = \min\{T/20 * n_t, n_{t,max}^l\}$;

　　**11: Update** $R(T) = 1 - \min\{\frac{T}{T_k}\tau, \tau\}$, $R_t(T) = 1 - \min\{\frac{T}{T_k}\tau_t, \tau_t\}$;

**end**

**12: Output** $F_{t1}$ and $F_{t2}$

---

is the minimized value of $\mathcal{L}(\theta, \boldsymbol{u}'; F, D)$ under the constrain $\mathbf{1}\boldsymbol{u}' > \alpha|D|$. It is clear that this solving process is equivalent to finding small-loss samples.

### 7.5.2 Training procedures of Butterfly (Algorithm 2)

**Update parameters of networks.** First, we initialize training data for two branches ($\tilde{D}$ for Branch-I and $\tilde{D}_t^l$ for Branch-II), four networks ($F_1, F_2, F_{t1}$ and $F_{t2}$) and the number of pseudo labels (line 2). In the first epoch ($T = 1$), following [13], $\tilde{D}_t^l$ is the same with $\tilde{D}_s$ (i.e., we use noisy source data as pseudo-labeled target data), since we cannot annotate pseudo labels for target data when $T = 1$. After mini-batch $\check{D}$ is fetched from $\tilde{D}$ (line 4), $F_1$ and $F_2$ check high-correctness data out and update their parameters (lines 5) using Algorithm 1. Using similar procedures, $F_{t1}$ and $F_{t2}$ also update their parameters using Algorithm 1 (lines 6-7).

**Assign pseudo labels.** In each epoch, after $N_{max}$ mini-batch updating, we randomly select $n_t^l$ unlabeled target data and assign them pseudo labels using the Labeling function [13], $F_1$ and $F_2$ (lines 8). Following [13], the Labeling function in Algorithm 2 (line 8) assigns pseudo labels to unlabeled target data, when predictions of $F_1$ and $F_2$ agree and at least one of them is confident about their predictions (probability above 0.9 or 0.95). Using this function, we can obtain the pseudo-labeled target data $\tilde{D}_t^l$ for training Branch-II in the next epoch. Then, we merge $\tilde{D}_t^l$ and $\tilde{D}_s$ to be $\tilde{D}$ for training Branch-I in the next epoch (line 9).

**Update other parameters.** Finally, we update $n_t^l$, $R(T)$ and $R_t(T)$ in lines 10-11. Note that $R(T)$ and $R_t(T)$ are actually piecewise-defined linear functions:

$$R(T) = \begin{cases} 1 - \tau, & T \geq T_k, \\ 1 - T/T_k \times \tau, & T \leq T_k, \end{cases}$$

$$R_t(T) = \begin{cases} 1 - \tau_t, & T \geq T_k, \\ 1 - T/T_k \times \tau_t, & T \leq T_k. \end{cases}$$

In Algorithm 2, we use $\tau$ to represent the noise rate (i.e., the ratio of data with incorrect labels) in MD and use $\tau_t$ to represent the noise rate in TD. However, in WUDA, we cannot obtain the ground-truth $\tau$ and $\tau_t$. Thus, we regard $\tau$ and $\tau_t$ as hyper-parameters.

### 7.6 Can we realize DCP using other models?

Based on Theorem 3, if we check high-correctness source data and pseudo-labeled target data out, we can reduce the negative effects of noisy source data significantly. Thus, we propose the DCP to check correct data out, which is introduced in Section 7.2. In Butterfly, we realize DCP using deep networks, since the memorization effects of deep networks ensures that we can check correct data out. For non-network models, if they also have memorization effects like deep networks, they can also be used into our approach. We also tried other models. Unfortunately, these models cannot fit the pattern first (like what deep networks did when fitting training data), meaning that, currently, we can only realize our approach using deep networks.

### 7.7 A Generalization Bound for WUDA

In this subsection, we prove a generalization bound for WUDA problem using the loss function Eq. (22) and Theorem 3[1]. Practitioner may safely skip it. First, we introduce the Rademacher complexity of a class of vector-valued functions [61], [62], [63], [64], [65], [66], which measures the degree to

---

1. Please note that this is a generalization bound for WUDA problem rather than Butterfly. In Butterfly, we essentially have four classifiers ($F_1, F_2, F_{t1}, F_{t2}$), which is very difficult to analyze it. We will develop a generalization and estimation error bound for Butterfly in the future.

which a class can fit random noise. Rademacher Complexity of $\mathcal{H}$ is defined as follows.

**Definition 1** (Rademacher Complexity of $\mathcal{H}$). *Given a sample $S = \{(x_i)\}_{i=1}^n$, the empirical Rademacher complexity of the set $\mathcal{H}$ is defined as follows.*

$$\hat{\Re}_S(\mathcal{H}) = \frac{2}{n}\mathbb{E}_\sigma\Big(\sup_{h\in\mathcal{H}}\sum_{i=1}^n\sum_{k=1}^K\sigma_{ik}h_k(x_i)\Big),$$

*where $h_k(\cdot)$ is the $k^{th}$ component of function $h\in\mathcal{H}$ and the $\sigma_{ik}$ are $n\times K$ matrix of independent Rademacher variables [63]. The Rademacher complexity of the set $\mathcal{H}$ is defined as the expectation of $\hat{\Re}_H(\mathcal{H})$ over all samples of size $n$:*

$$\Re_n(\mathcal{H}) = \mathbb{E}_S\Big(\hat{\Re}_S(\mathcal{H})\Big||S| = n\Big).$$

Then, using the Rademacher complexity, we can prove an upper bound of $\tilde{R}_s^{\text{po}}(h, u_s)$ to show the relation between $\tilde{R}_s^{\text{po}}(h, u_s)$ and the loss function Eq. (22). As a common practice [67], [68], we assume that, 1) there are $C_h > 0$ and $C_L > 0$ such that $\sup_{h\in\mathcal{H}}\|h\|_\infty \le C_h$ and $\sup_{\|t\|_\infty \le C_h}\max_y \ell(t, y) \le C_L$, and 2) $\ell(t, y)$ is Lipschitz continuous in $\|t\|_\infty \le C_h$ with a Lipschitz constant $L_\ell$.

**Lemma 3.** *Given a sample $S_s = \{(x_{si}, y_{si}, u_{si})\}_{i=1}^n$ drawn from the probability density $\tilde{p}_s^{po}(x_s, y_s, u_s)$, with the probability of at least $1-\delta$ over samples $S_s$ of size $n$ drawn from $\tilde{p}_s^{po}(x_s, y_s, u_s)$, the following inequality holds.*

$$\tilde{R}_s^{po}(h, \boldsymbol{u}_s) \le \mathcal{L}(\theta, h; \boldsymbol{u}_s, D_s^{xy}) + \frac{\sqrt{2}L_\ell\hat{\Re}_{D_s^x}(\mathcal{H})}{1 - \tau_s}$$
$$+ \frac{3C_L}{1 - \tau_s}\sqrt{\frac{\ln\frac{\delta}{2}}{2n}}, \tag{23}$$

*where $\mathcal{L}$ is defined in Eq. (22), $D_s^{xy} = \{x_{si}, y_{si}\}_{i=1}^n$, $D_s^x = \{x_{si}\}_{i=1}^n$, $\boldsymbol{u}_s = [u_{s1}, \ldots, u_{sn}]^T$ and $\tau_s = \rho_{u_s} = 1 - \sum_{i=1}^n u_{si}/n$.*

*Proof (sketch).* For simplicity, in this proof, we let $\mathcal{L}_{S_s}(\ell, h) = \mathcal{L}(\theta, h; \boldsymbol{u}_s, D_s^{xy})$, $\tilde{R}_s^{\text{po}}(\ell, h) = \tilde{R}_s^{\text{po}}(h, u_s)$, and $\mathbb{E}_{S_s}[\cdot] = \mathbb{E}_{S_s\sim(\tilde{P}_s^{\text{po}})^n}[\cdot]$, where $\tilde{P}_s^{\text{po}}$ is the probability measure corresponding to the density $\tilde{p}_s^{\text{po}}$. We first prove that $\mathcal{L}_{S_s}(\ell, h)$ is an unbiased estimator of $\tilde{R}_s^{\text{po}}(\ell, h)$ based on the definition of $\tilde{R}_s^{\text{po}}(h, u_s)$ in Section 6.

Then, let $\Phi(S_s) = \sup_{\ell\in\mathbb{L}_\mathcal{H}}\left(\tilde{R}_s^{\text{po}}(\ell, h) - \mathcal{L}_{S_s}(\ell, h)\right)$. Changing a point of $S_s$ affects $\Phi(S_s)$ at most $C_L/(n(1 - \tau_s))$. Thus, by McDiarmid's inequality applied to $\Phi(S_s)$, for any $\delta > 0$, with probability of at least $1 - \delta/2$, the following inequality holds.

$$\Phi(S_s) \le \mathbb{E}_{S_s}[\Phi(S_s)] + \frac{C_L}{1 - \tau_s}\sqrt{\frac{\ln(\delta/2)}{2n}}.$$

Then, we have

$$\mathbb{E}_{S_s}[\Phi(S_s)] = \mathbb{E}_{S_s}\Big[\sup_{\ell\in\mathbb{L}_\mathcal{H}}\big(\tilde{R}_s^{\text{po}}(h, u_s) - \mathcal{L}_{S_s}(h)\big)\Big]$$
$$\le \frac{2}{n(1 - \tau_s)}\mathbb{E}_{\sigma, S_s}\Big[\sup_{\ell\in\mathbb{L}_\mathcal{H}}\sum_{i=1}^n\sigma_i u_{si}\ell(h(x_{si}), y_{si})\Big]. \tag{24}$$

Because of the existence of $u_{si}$, Eq. (24) is not the Rademacher complexity of $\mathbb{L}_\mathcal{H}$ (i.e., $\Re(\mathbb{L}_\mathcal{H})$). However, we can prove

that Eq. (24) can be bounded by $\Re(\mathbb{L}_\mathcal{H})/(1 - \tau_s)$ using the property of sup.

Since changing a point of $S_s$ affects $\Re_n(\mathbb{L}_\mathcal{H})$ at most $2C_L/n$, by McDiarmid's inequality, for any $\delta > 0$, with probability of at least $1 - \delta/2$, the following inequality holds.

$$\Re_n(\mathbb{L}_\mathcal{H}) \le \hat{\Re}_{S_s}(\mathbb{L}_\mathcal{H}) + 2C_L\sqrt{\frac{\ln(\delta/2)}{2n}}.$$

Since $\ell$ is Lipschitz continuous, according to [63], we have

$$\hat{\Re}_{S_s}(\mathbb{L}_\mathcal{H}) \le \sqrt{2}L_\ell\hat{\Re}_{D_s^x}(\mathcal{H}),$$

which proves this lemma. $\square$

Finally, we prove the generalization bound for WUDA problem as follows.

**Theorem 4.** *Given a sample $S_s = \{(x_{si}, y_{si}, u_{si})\}_{i=1}^{n_s}$ drawn from the probability density $\tilde{p}_s^{po}(x_s, y_s, u_s)$ and a sample $S_t = \{(x_{ti}, u_{ti})\}_{i=1}^{n_t}$ drawn from the probability density $\tilde{p}_t^{po}(x_t, u_t)$, under the assumptions in Remark 2, Remark 3 and Lemma 2, $\forall\epsilon\in(0, 1)$, there are $\delta_s$ and $\delta_t$, if $\rho_{01}^s < \delta_s$ and $\rho_{01}^t < \delta_t$, then, with the probability of at least $1 - 3\delta$, for any $h\in\mathcal{H}$, the following inequality holds.*

$$R_t(h, f_t) \le 2\Big(\mathcal{L}(\theta, h; \boldsymbol{u}_s, D_s^{xy}) + \mathcal{L}(\theta, h; \boldsymbol{u}_s, D_{\tilde{s}}^{xy})\Big)$$
$$+ \mathcal{L}(\theta, h; \boldsymbol{u}_t, D_{\tilde{t}}^{xy}) + \frac{4\sqrt{2}L_\ell\hat{\Re}_{D_s^x}(\mathcal{H})}{1 - \tau_s}$$
$$+ \frac{\sqrt{2}L_\ell\hat{\Re}_{D_t^x}(\mathcal{H})}{1 - \tau_t} + \frac{12C_L}{1 - \tau_s}\sqrt{\frac{\ln\frac{\delta}{2}}{2n_s}}$$
$$+ \frac{3C_L}{1 - \tau_t}\sqrt{\frac{\ln\frac{\delta}{2}}{2n_t}} + 6\epsilon, \tag{25}$$

*where $\mathcal{L}$ is defined in Eq. (22), $D_s^{xy} = \{x_{si}, y_{si}\}_{i=1}^{n_s}$, $D_{\tilde{s}}^{xy} = \{x_{si}, \tilde{f}_t(x_{si})\}_{i=1}^{n_s}$, $D_{\tilde{t}}^{xy} = \{x_{ti}, \tilde{f}_t(x_{ti})\}_{i=1}^{n_t}$ $D_s^x = \{x_{si}\}_{i=1}^{n_s}$, $D_t^x = \{x_{ti}\}_{i=1}^{n_t}$, $\boldsymbol{u}_s = [u_{s1}, \ldots, u_{sn_s}]^T$, $\tau_s = \rho_{u_s} = 1 - \sum_{i=1}^{n_s} u_{si}/n_s$, $\boldsymbol{u}_t = [u_{t1}, \ldots, u_{tn_t}]^T$ and $\tau_t = \rho_{u_t} = 1 - \sum_{i=1}^{n_t} u_{ti}/n_t$.*

*Proof.* We prove this theorem (i.e., Inequality (25)) according to Inequality (21), where (25) has 7 terms in the right side and (21) have 6 terms in the right side.

1) For last 3 terms in (21), according to (16), (17) and (18), we know the sum of last three terms of (21) is less than or equal to $4\epsilon$.

2) For first 3 terms in (21), we have shown that (in Section 7.1) the sum of the first 3 terms in (21) is less than or equal to $(*)$:

$$2\tilde{R}_s^{\text{po}}(h, u_s) + 2\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, \boldsymbol{u}_s) + \tilde{R}_t^{\text{po}}(h, \tilde{f}_t, \boldsymbol{u}_t) + 2\epsilon.$$

Then, we can prove that (similar with Lemma 3), with probability of at least $1 - \delta$, for any $h\in\mathcal{H}$,

$$\tilde{R}_s^{\text{po}}(h, \tilde{f}_t, \boldsymbol{u}_s) \le \mathcal{L}(\theta, h; \boldsymbol{u}_s, D_{\tilde{s}}^{xy}) + \frac{\sqrt{2}L_\ell\hat{\Re}_{D_s^x}(\mathcal{H})}{1 - \tau_s}$$
$$+ \frac{3C_L}{1 - \tau_s}\sqrt{\frac{\ln\frac{\delta}{2}}{2n_s}}, \tag{26}$$

$$\tilde{R}_t^{\text{po}}(h, \tilde{f}_t, \boldsymbol{u}_t) \leq \mathcal{L}(\theta, h; \boldsymbol{u}_t, D_{\tilde{t}}^{xy}) + \frac{\sqrt{2}L_\ell \hat{\Re}_{D_t^x}(\mathcal{H})}{1 - \tau_s}$$

$$+ \frac{3C_L}{1 - \tau_t} \sqrt{\frac{\ln \frac{\delta}{2}}{2n_t}}. \tag{27}$$

Combining (23), (26), (27) with (∗), based on 1), we prove this theorem. Note that, $6\epsilon$ equals $4\epsilon$ (in 1)) + $2\epsilon$ (in (∗)). □

**Corollary 1** (Generalization Bound for WUDA). *Given a sample $S_s$ and a sample $S_t$ defined in Theorem 4, under the assumptions in Remark 2, Remark 3 and Lemma 2, if $\rho_{01}^s < C_\rho^s / \sqrt{n_s T}$ and $\rho_{01}^t < C_\rho^t / \sqrt{n_t T}$, then, with the probability of at least $1 - 3\delta$, for any $h \in \mathcal{H}$, the following inequality holds.*

$$R_t(h, f_t)$$
$$\leq 2\Big(\mathcal{L}(\theta, h; \boldsymbol{u}_s, D_s^{xy}) + \mathcal{L}(\theta, h; \boldsymbol{u}_s, D_{\tilde{s}}^{xy})\Big)$$
$$+ \mathcal{L}(\theta, h; \boldsymbol{u}_t, D_{\tilde{t}}^{xy}) + \frac{4\sqrt{2}L_\ell \hat{\Re}_{D_s^x}(\mathcal{H})}{1 - \tau_s}$$
$$+ \frac{\sqrt{2}L_\ell \hat{\Re}_{D_t^x}(\mathcal{H})}{1 - \tau_t} + \frac{12C_L}{1 - \tau_s} \sqrt{\frac{\ln \frac{\delta}{2}}{2n_s}} + \frac{3C_L}{1 - \tau_t} \sqrt{\frac{\ln \frac{\delta}{2}}{2n_t}}$$
$$+ \frac{C_\rho^s (2M_s + M_t)}{\sqrt{n_s T}} + \frac{3C_\rho^t M_t}{\sqrt{n_t T}}, \tag{28}$$

*where $\mathcal{L}$, $D_s^{xy}$, $D_{\tilde{s}}^{xy}$, $D_{\tilde{t}}^{xy}$, $D_t^x$, $\boldsymbol{u}_s$, $\tau_s$, $\boldsymbol{u}_t$, $\tau_t$ are defined in Theorem 4, $T$ is the number of training epochs, and $C_\rho^s$ and $C_\rho^t$ are two finite constants.*

**Remark 9.** In Corollary 4, we assume that $\rho_{01}^s$ and $\rho_{01}^t$ will go to zero with the convergence speed of $O(1/\sqrt{n_s T})$ and $O(1/\sqrt{n_t T})$, respectively. In Section 9.8, we verify this assumption through our experiments.

Corollary 4 shows the empirical upper bound of the target risk (i.e., $R_t(h, f_t)$). Based on this bound, we can obtain the estimation error bound of $R_t(h, f_t)$ as follows. First, let

$$\hat{R}_t^{\mathcal{L}}(h, S_s, S_t) = 2\Big(\mathcal{L}(\theta, h; \boldsymbol{u}_s, D_s^{xy}) + \mathcal{L}(\theta, h; \boldsymbol{u}_s, D_{\tilde{s}}^{xy})\Big)$$
$$+ \mathcal{L}(\theta, h; \boldsymbol{u}_t, D_{\tilde{t}}^{xy}), \tag{29}$$

where $D_s^{xy}$, $D_{\tilde{s}}^{xy}$ and $D_{\tilde{t}}^{xy}$ are defined in Theorem 4, and $\tilde{h} = \arg\min_{h \in \mathcal{H}} \hat{R}_t^{\mathcal{L}}(h, S_s, S_t)$ means the empirical minimizer of $\hat{R}_t^{\mathcal{L}}(h, S_s, S_t)$, and $h^* = \arg\min_{h \in \mathcal{H}} R_t(h, f_t)$ means the true risk minimizer of $R_t(h, f_t)$, and $\mathcal{H}' = \{h | \hat{R}_t^{\mathcal{L}}(h, S_s, S_t) \leq \epsilon'\}$. Then, we have

$$R_t(\tilde{h}, f_t) - R_t(h^*, f_t)$$
$$= R_t(\tilde{h}, f_t) - \hat{R}_t^{\mathcal{L}}(\tilde{h}, S_s, S_t) + \hat{R}_t^{\mathcal{L}}(\tilde{h}, S_s, S_t) - R_t(h^*, f_t)$$
$$\quad + \hat{R}_t^{\mathcal{L}}(h^*, S_s, S_t) - \hat{R}_t^{\mathcal{L}}(h^*, S_s, S_t)$$
$$= R_t(\tilde{h}, f_t) - \hat{R}_t^{\mathcal{L}}(\tilde{h}, S_s, S_t) + \hat{R}_t^{\mathcal{L}}(h^*, S_s, S_t) - R_t(h^*, f_t)$$
$$\quad + \hat{R}_t^{\mathcal{L}}(\tilde{h}, S_s, S_t) - \hat{R}_t^{\mathcal{L}}(h^*, S_s, S_t)$$
$$\leq \sup_{h \in \mathcal{H}'} (R_t(h, f_t) - \hat{R}_t^{\mathcal{L}}(h, S_s, S_t)) + \epsilon' + 0, \tag{30}$$

where $\hat{R}_t^{\mathcal{L}}(\tilde{h}, S_s, S_t) \leq \hat{R}_t^{\mathcal{L}}(h^*, S_s, S_t)$ due to the definition of $\tilde{h}$. If all conditions in Theorem 4 are satisfied, with the

probability of at least $1 - 3\delta$, for any $h \in \mathcal{H}$, we have

$$R_t(\tilde{h}, f_t) - R_t(h^*, f_t)$$
$$\leq \frac{4\sqrt{2}L_\ell \hat{\Re}_{D_s^x}(\mathcal{H})}{1 - \tau_s} + \frac{\sqrt{2}L_\ell \hat{\Re}_{D_t^x}(\mathcal{H})}{1 - \tau_t} + \frac{12C_L}{1 - \tau_s} \sqrt{\frac{\ln \frac{\delta}{2}}{2n_s}}$$
$$+ \frac{3C_L}{1 - \tau_t} \sqrt{\frac{\ln \frac{\delta}{2}}{2n_t}} + \frac{C_\rho^s (2M_s + M_t)}{\sqrt{n_s T}} + \frac{3C_\rho^t M_t}{\sqrt{n_t T}} + \epsilon'. \tag{31}$$

Eq. (31) ensures that learning with $\hat{R}_t^{\mathcal{L}}(\tilde{h}, S_s, S_t)$ is consistent: as $n_s, n_t \to \infty$ and $\epsilon' \to 0$, $R_t(\tilde{h}, f_t) \to R_t(h^*, f_t)$. For linear-in-parameter model with a bounded norm, $\hat{\Re}_{D_s^x}(\mathcal{H}) = \mathcal{O}(1/\sqrt{n_s})$ and $\hat{\Re}_{D_t^x}(\mathcal{H}) = \mathcal{O}(1/\sqrt{n_t})$ and thus $R_t(\tilde{h}, f_t) \to R_t(h^*, f_t)$ in $\mathcal{O}(1/\sqrt{n_s} + 1/\sqrt{n_t})$.

## 8 COMPARISON TO RELATED WORKS

In this section, we compare Butterfly with related works and show why related works cannot handle WUDA problem.

**Relations to co-teaching.** As Butterfly is related to co-teaching, we discuss their major differences here. Although co-teaching applies the small-loss trick and the cross-update technique to train deep networks against noisy data, it can only deal with one-domain problem instead cross-domain problem. Besides, we argue that Butterfly is not a simple mixtrue of co-teaching and ATDA for two reasons.

First, network structure of Butterfly is different with that of ATDA and co-teaching: Butterfly maintains four networks; while ATDA maintains three and co-teaching maintains two. We cannot simply combine ADTA and co-teaching to derive Butterfly. Second, we have justified that the sequential mixture of co-teaching and ATDA (i.e., two-step method) cannot eliminate noise effects caused by noisy source data (see Section 5.2). Specifically, two-step methods only take care of part of noise effects but Butterfly takes care of the whole noise effects. Thus, Butterfly is the first method to eliminate noise effects rather than alleviate it.

**Relations to TCL.** Recently, *transferable curriculum learning* (TCL) is a robust UDA method to handle noise [44]. TCL uses small-loss trick to train DANN [23]. However, TCL can only minimize $(i) + (ii) + (iv)$, while Butterfly can minimize all terms in the right side of Eq. (15).

## 9 EXPERIMENTS

We conduct experiments on 32 simulated WUDA tasks and 3 real-world WUDA tasks to verify the efficacy of Butterfly.

### 9.1 Simulated WUDA tasks

We verify the effectiveness of our approach on three benchmark datasets (vision and text), including *MNIST*, *SYN-DIGITS (SYND)*[2] and *human-sentiment* analysis (i.e., *Amazon products reviews* on *book*, *dvd*, *electronics* and *kitchen*)[3]. They are used to construct 14 basic tasks: *MNIST→SYND (M→S)*, *SYND→MNIST (S→M)*, *book→dvd (B→D)*, *book→electronics*

---

2. *Digit* datasets (*MNIST* and *SYN Digit*) can be downloaded from official code of ATDA. The link is https://github.com/ksaito-ut/atda.

3. *Sentiment* datasets (*Amazon products reviews*) can be downloaded from the official code of marginalized Stacked Denoising Autoencoder. The link is https://www.cse.wustl.edu/~mchen/code/mSDA.tar.

(a) *MNIST*                    (b) *SYND*

Fig. 4. Visualization of *MNIST* and *SYND*.



(a) *Bing* provided by [69]     (b) *Caltech256* provided by [70]     (c) *ImageNet* provided by [71]     (d) *SUN* provided by [72]

Fig. 5. Visualization of *Bing*, *Caltech256*, *ImageNet* and *SUN* (taking "horse" as the common class).

($B{\rightarrow}E$), . . . , and *kitchen $\rightarrow$ electronics* ($K{\rightarrow}E$). These tasks are often used for evaluation of UDA methods [13], [23], [25]. Figure 4 shows datasets *MNIST* and *SYND*.

Since all source datasets are clean, we corrupt source data using symmetry flipping [58] and pair flipping [56] with noise rate $\rho$ chosen from $\{0.2, 0.45\}$. Note that, there are other ways to generate the noisy source data, such as asymmetry flipping. However, since the asymmetry flipping can be regarded as the combination of symmetry flipping and pair flipping, we only use symmetry flipping and pair flipping to generate *simulated* WUDA tasks. In *real-world* WUDA tasks, we have more complex noisy source data, where the noisy type in the source domain is unknown.

Therefore, for each basic task, we have four kinds of noisy source data: *Pair*-45% (P45), *Pair*-20% (P20), *Symmetry*-45% (S45), *Symmetry*-20% (S20). Following [54], [56], we can corrupt clean-label datasets manually using the noise transition matrix $Q_S$ and $Q_P$. Namely, we evaluate the performance of each method using 32 simulated WUDA tasks: 8 digit tasks and 24 human-sentiment tasks. Since the human-sentiment task is a binary classification problem, pair flipping is equal to symmetry flipping, meaning that we have 24 human-sentiment tasks.

### 9.2 Real-world WUDA tasks

We also verify the efficacy of our approach on "cross-dataset benchmark" including *Bing*, *Caltech256*, *Imagenet* and *SUN* [48] [4]. In this benchmark, *Bing*, *Caltech256*, *Imagenet* and *SUN* contain common 40 classes. Since *Bing* dataset was formed by

---

4. *Real-world* datasets (*BCIS*) can be downloaded from the website of the project "A Testbed for Cross-Dataset Analysis": https://sites.google.com/site/crossdataset/home/files ("setup DENSE decaf7", 1.3GB, decaf7 features).

collecting images retrieved by Bing image search, it contains rich noisy data, with presence of multiple objects in the same image and caricaturization [48]. We use *Bing* as noisy source data, and *Caltech256*, *Imagenet* and *SUN* as unlabeled target data, which can form three real-world WUDA tasks. Figure 5 shows datasets *Bing*, *Caltech256*, *Imagenet* and *SUN* (taking "horse" as the common class).

### 9.3 Baselines

We realize Butterfly using four networks (B-Net) and compare B-Net with following baselines: 1) ATDA: representative pseudo-labeling-based UDA method [13]; 2) DAN: representative IPM-based UDA method [36]; 3) DANN: representative adversarial-training-based UDA method [23]; 4) *Manifold embedded distribution alignment* (MEDA): a representative non-deep UDA method [73]; 5) TCL: an existing robust UDA method; 6) co-teaching+ATDA (Co+ATDA): a two-step method (see Section 5.2); and 7) co-teaching+TCL (Co+TCL). Since MEDA cannot extract features from images, we only compare with MEDA on human-sentiment tasks, where features are already given.

### 9.4 Network structure and optimizer

We implement all methods on Python 3.6 with a NIVIDIA P100 GPU. We use MomentumSGD for optimization in digit and real-world tasks, and set the momentum as 0.9. We use Adagrad for optimization in human-sentiment tasks because of sparsity of review data [13]. $F_1$, $F_2$, $F_{t1}$ and $F_{t2}$ are 6-layer CNN (3 convolutional and 3 fully-connected layers) for digit tasks; and are 3-layer neural networks (3 fully-connected layers) for human-sentiment tasks; and are 4-layer neural networks (4 fully-connected layers) for real-world tasks. The ReLU active function is used as activation function of these
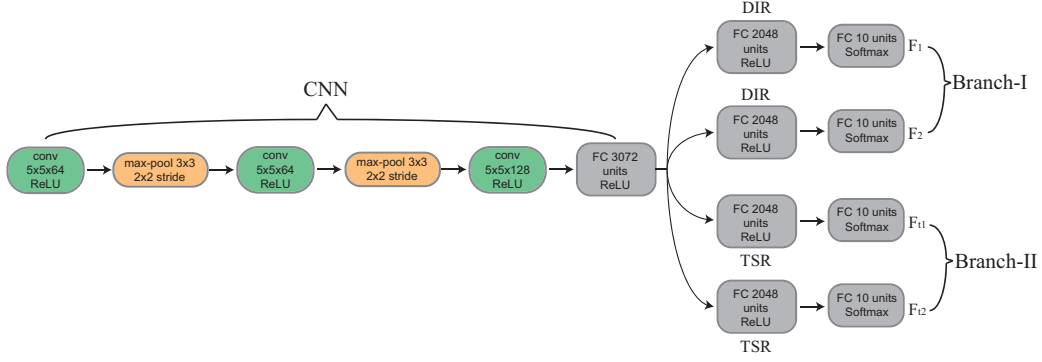
Fig. 6. The architecture of B-Net for digit WUDA tasks *SYND* ↔ *MNIST*. We added BN layer in the last convolution layer in CNN and FC layers in $F_1$ and $F_2$. We also used dropout in the last convolution layer in CNN and FC layers in $F_1$, $F_2$, $F_{t1}$ and $F_{t2}$ (dropout probability is set to $0.5$).



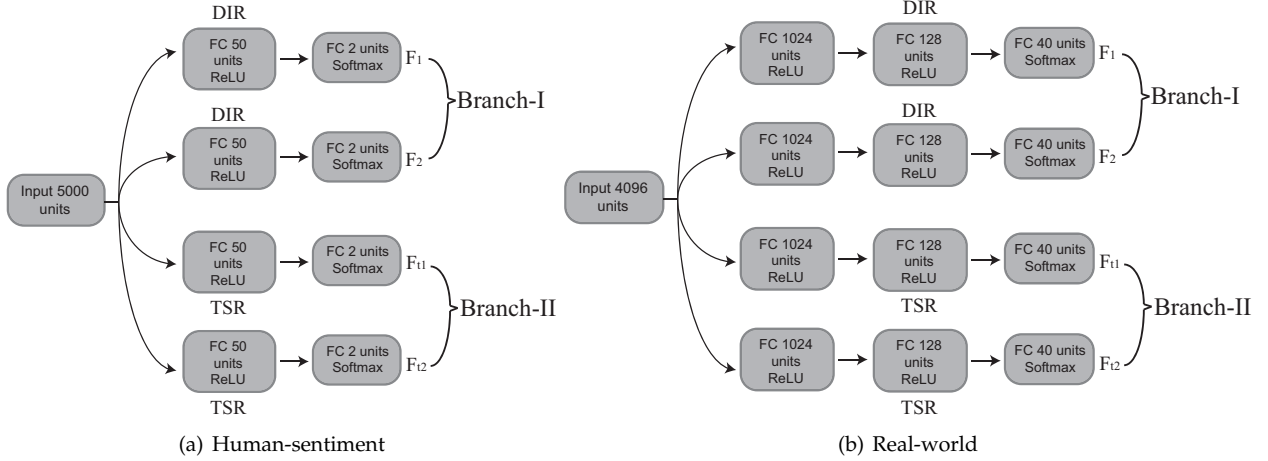(a) Human-sentiment

(b) Real-world

Fig. 7. The architecture of B-Net for (a) human-sentiment WUDA tasks and (b) real-world WUDA tasks. We added BN layer in the first FC layers in $F_1$ and $F_2$. We also used dropout in the first FC layers in $F_1$, $F_2$, $F_{t1}$ and $F_{t2}$ (dropout probability is set to $0.5$).

networks. Besides, dropout and batch normalization are also used. The network topology is shown in Figures 6 and 7. As deep networks are highly nonconvex, even with the same network and optimization method, different initializations can lead to different local optimal. Thus, following [57], we take four networks with the same architecture but different initialization as four classifiers.

## 9.5 Experimental setup

Since this paper deals with the challenging situation where no labeled data are available in the target domain, we follow the common protocol to set hyperparameters that the similar tasks have the same hyperparameters [37]. For example, we set the same hyperparameters for all WUDA tasks regarding digit datasets (there are 8 WUDA tasks regarding digit datasets). The selected hyperparameters are robust to many tasks rather than a specific task. Details can be found below.

For all 35 WUDA tasks, $T_k$ is set to 5, and $T_{max}$ is set to 30, and $\ell(\cdot, \cdot)$ is the cross-entropy loss function. Learning rate is set to $0.01$ for simulated tasks and $0.05$ for real-world WUDA tasks, $\tau_t$ is set to $0.05$ for simulated tasks and $0.02$ for real-world WUDA tasks. Confidence level of labeling function in line 8 of Algorithm 2 is set to $0.95$ for 8 digit tasks, and $0.9$ for 24 human-sentiment tasks and $0.8$ for real-world WUDA tasks. $\tau$ is set to $0.4$ for digit tasks, $0.1$ for human-sentiment tasks, $0.2$ for real-world WUDA tasks. $n_{t,max}^l$ is

set to $15,000$ for digit tasks, $500$ for human-sentiment tasks and $4000$ for real-world WUDA tasks. $N_{max}$ is set to $1000$ for digit tasks, and $200$ for human-sentiment and real-world tasks. Batch size is set to $128$ for digit, real-world WUDA tasks, and $24$ for human-sentiment tasks. Penalty parameter is set to $0.01$ for digit, real-world WUDA tasks, and $0.001$ for human-sentiment tasks.

To fairly compare all methods, they have the same network structure. Namely, ATDA, DAN, DANN, TCL and B-Net adopt the same network structure for each dataset. Note that DANN and TCL use the same structure for their discriminate networks. All experiments are repeated 10 times and we report the average accuracy values and *standard deviation* (STD) of accuracy values of 10 experiments.

## 9.6 Results on simulated WUDA tasks

This subsection presents accuracy on unlabled target data (i.e., target-domain accuracy) in 32 simulated WUDA tasks.

### 9.6.1 Results on digits WUDA tasks

Table 1 reports the target-domain accuracy in 8 digit tasks. As can be seen, average target-domain accuracy of B-Net is higher than those of all baselines. On S20 case (the easiest case), most methods work well. ATDA has a satisfactory performance although it does not consider the noise effects explicitly. Then, when facing harder cases (i.e., P20 and P45),
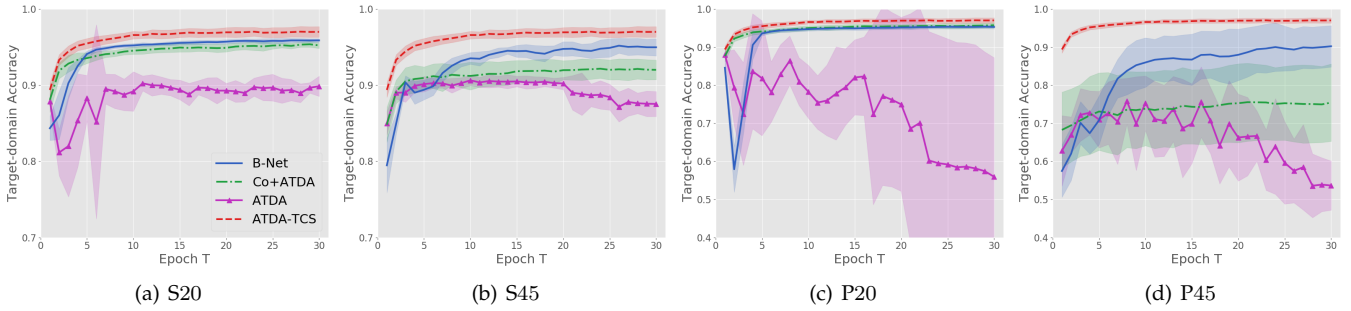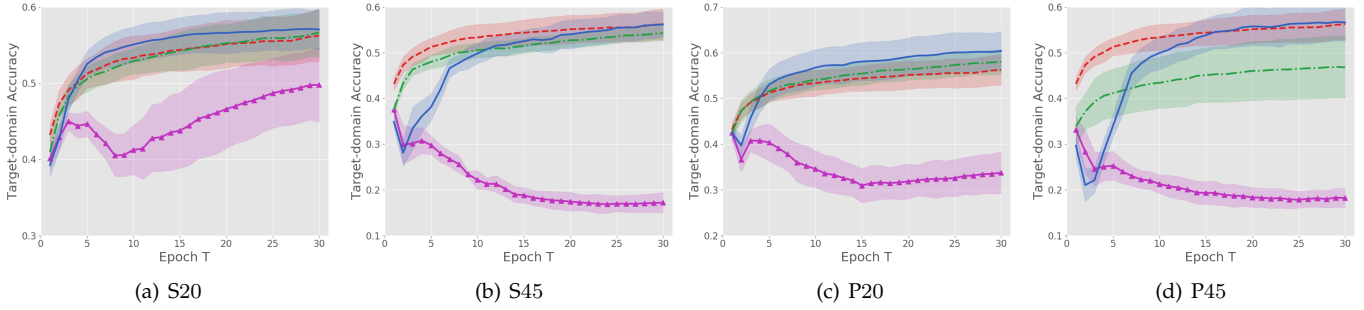
Fig. 8. Target-domain accuracy vs. number of epochs on four *SYND→MNIST* WUDA tasks.



Fig. 9. Target-domain accuracy vs. number of epochs on four *MNIST→SYND* WUDA tasks.

TABLE 1
Target-domain accuracy on $8$ digit WUDA tasks (*SYND↔MNIST*). Bold value represents the highest accuracy in each row.

| Tasks | Type | DAN | DANN | ATDA | TCL | Co+TCL | Co+ATDA | B-Net |
|-------|------|-----|------|------|-----|--------|---------|-------|
| $S{\to}M$ | P20 | 90.17% | 79.06% | 55.95% | 80.81% | 88.56% | **95.37%** | 95.29% |
| | P45 | 67.00% | 55.34% | 53.66% | 55.97% | 73.27% | 75.43% | **90.21%** |
| | S20 | 90.74% | 75.19% | 89.87% | 80.23% | 85.88% | 95.22% | **95.88%** |
| | S45 | 89.31% | 65.87% | 87.53% | 68.54% | 75.69% | 92.03% | **94.97%** |
| $M{\to}S$ | P20 | 40.82% | 58.78% | 33.74% | 58.88% | 59.08% | 58.02% | **60.36%** |
| | P45 | 28.41% | 43.70% | 19.50% | 45.31% | 47.15% | 46.80% | **56.62%** |
| | S20 | 30.62% | 53.52% | 49.80% | 56.74% | 56.91% | 56.64% | **57.05%** |
| | S45 | 28.21% | 43.76% | 17.20% | 49.91% | 51.22% | 54.29% | **56.18%** |
| Average | | 58.16% | 58.01% | 50.91% | 62.05% | 67.22% | 71.73% | **75.82%** |

ATDA fails to transfer useful knowledge from noisy source data to unlabeled target data. When facing the hardest cases (i.e., $M{\to}S$ with P45 and S45), DANN has higher accuracy than DAN and ATDA have. However, when facing the easiest cases (i.e., $S{\to}M$ with P20 and S20), the performance of DANN is worse than that of DAN and ATDA.

Although two-step method Co+ATDA (or Co+TCL) outperforms ATDA (or TCL) in all $8$ tasks, it cannot beat one-step method: B-Net in terms of average target-domain accuracy. This result is an evidence for the claim in Section 5.2. In the task $S{\to}M$ with P20, Co+ATDA outperforms all methods (slightly higher than B-Net), since pseudo-labeled source data are almost correct.

Figures 8 and 9 show the target-domain accuracy vs. number of epochs among ATDA, Co+ATDA and B-Net. Besides, we show the accuracy of ATDA trained with clean source data (ATDA-TCS) as a reference point. When accuracy of one method is close to that of ATDA-TCS (red dash line), this method successfully eliminates noise effects. From our observations, it is clear that B-Net is very close to ATDA-TCS in 7 out of 8 tasks (except for $S{\to}M$ task with P45, Figure 8-(d)), which is an evidence that Butterfly can eliminate noise effects. Since P45 case is the hardest one and we only have finite samples, it is reasonable that B-Net cannot perfectly

eliminate noise effects. An interesting phenomenon is that, B-Net outperforms ATDA-TCS in 2 $M{\to}S$ tasks (Figure 9-(a), (c)). This means that B-Net transfers more useful knowledge (from noisy source data to unlabeled target data) even than ATDA-TCS (from clean source data to unlabeled target data).

### 9.6.2 Results on human sentiment WUDA tasks

Tables 2 and 3 report the target-domain accuracy of each method in $24$ human-sentiment WUDA tasks. For these tasks, B-Net has the highest average target-domain accuracy. It should be noted that two-step method does not always perform better than existing UDA methods, such as for $20\%$-noise situation. The reason is that co-teaching performs poorly when pinpointing clean source data from noisy source data. Another observation is that noise effects is not eliminated like target-domain accuracy in $8$ digit WUDA tasks. The reason mainly includes that 1) these datasets only provide predefined features (i.e., we cannot extract better features from original contents in the training process), and 2) we only have finite samples and the number of samples in these datasets is smaller than those of digit datasets.

### 9.7 Results on real-world WUDA tasks

Table 4 reports the target-domain accuracy in 3 tasks. B-Net enjoys the best performance on all tasks. It should be noted

TABLE 2
Target-domain accuracy on $12$ human-sentiment WUDA tasks with the $20\%$ noise rate. Bold values mean the highest values in each row.

| Tasks | DAN | DANN | ATDA | TCL | MEDA | Co+TCL | Co+ATDA | B-Net |
|-------|-----|------|------|-----|------|--------|---------|-------|
| $B \rightarrow D$ | 68.28% | 68.08% | 70.31% | 71.40% | 74.81% | 67.81% | 66.70% | **71.84%** |
| $B \rightarrow E$ | 63.78% | 63.53% | 72.79% | 65.08% | 65.18% | 60.54% | 68.89% | **75.92%** |
| $B \rightarrow K$ | 65.48% | 64.63% | 71.79% | 66.80% | 68.65% | 61.23% | 66.51% | **76.32%** |
| $D \rightarrow B$ | 64.63% | 64.52% | 70.25% | 67.33% | 67.63% | 65.22% | 68.04% | **70.56%** |
| $D \rightarrow E$ | 65.33% | 65.16% | 69.99% | 66.74% | 69.51% | 64.55% | 67.32% | **73.73%** |
| $D \rightarrow K$ | 65.68% | 66.28% | 74.53% | 68.82% | 72.24% | 67.98% | 72.20% | **77.97%** |
| $E \rightarrow B$ | 60.41% | 60.15% | **63.89%** | 63.13% | 63.36% | 61.18% | 61.08% | 62.22% |
| $E \rightarrow D$ | 62.35% | 61.67% | 62.30% | 62.93% | 66.18% | 60.81% | 59.77% | **63.53%** |
| $E \rightarrow K$ | 72.05% | 71.51% | 74.00% | 75.36% | 75.42% | 72.65% | 70.85% | **78.96%** |
| $K \rightarrow B$ | 59.94% | 59.40% | **63.53%** | 62.77% | 65.13% | 60.71% | 61.22% | 63.36% |
| $K \rightarrow D$ | 61.46% | 61.51% | 64.66% | 64.16% | 66.87% | 64.15% | 64.94% | **66.98%** |
| $K \rightarrow E$ | 70.60% | 72.23% | 74.75% | 74.14% | 75.99% | 68.95% | 69.69% | **76.96%** |
| Average | 65.00% | 64.89% | 69.40% | 67.39% | 69.25% | 64.65% | 66.43% | **71.53%** |

TABLE 3
Target-domain accuracy on $12$ human-sentiment WUDA tasks with the $45\%$ noise rate. Bold values mean the highest values in each row.

| Tasks | DAN | DANN | ATDA | TCL | MEDA | Co+TCL | Co+ATDA | B-Net |
|-------|-----|------|------|-----|------|--------|---------|-------|
| $B \rightarrow D$ | 52.43% | 52.98% | 53.56% | 54.44% | 54.50% | 53.21% | 54.32% | **56.59%** |
| $B \rightarrow E$ | 52.17% | 53.50% | 55.14% | 54.14% | 54.29% | 53.98% | **57.34%** | 55.74% |
| $B \rightarrow K$ | 52.89% | 51.84% | 51.14% | 53.32% | 53.68% | 51.77% | 53.28% | **57.00%** |
| $D \rightarrow B$ | 53.11% | 53.04% | 54.48% | 53.27% | 53.66% | 54.85% | **55.95%** | 55.15% |
| $D \rightarrow E$ | 51.30% | 53.04% | 54.21% | 53.77% | 54.11% | 55.63% | 56.08% | **58.91%** |
| $D \rightarrow K$ | 52.15% | 53.17% | 57.99% | 52.45% | 52.45% | 58.10% | 59.94% | **66.20%** |
| $E \rightarrow B$ | 51.38% | 51.08% | 52.54% | 52.14% | 52.56% | 54.88% | 53.30% | **54.93%** |
| $E \rightarrow D$ | 52.83% | 51.24% | 49.02% | 52.57% | 53.03% | 50.03% | 49.62% | **52.88%** |
| $E \rightarrow K$ | 54.21% | 53.58% | 51.66% | 55.04% | 55.42% | 56.15% | 52.10% | **56.12%** |
| $K \rightarrow B$ | 50.44% | 51.77% | **51.96%** | 51.50% | 51.52% | 53.81% | 52.59% | 51.39% |
| $K \rightarrow D$ | 52.20% | 51.45% | 52.86% | 53.19% | 53.38% | 55.69% | 54.52% | **53.53%** |
| $K \rightarrow E$ | **54.72%** | 53.33% | 52.11% | 53.46% | 53.81% | 51.26% | 52.62% | 53.71% |
| Average | 52.49% | 52.50% | 53.65% | 53.27% | 53.54% | 54.11% | 54.31% | **56.01%** |

TABLE 4
Target-domain accuracy on $3$ real-world WUDA tasks. The source domain is the *Bing* dataset that contains noisy information from the Internet. Bold value represents the highest accuracy in each row.

| Target | DAN | DANN | ATDA | TCL | Co+TCL | Co+ATDA | B-Net |
|--------|-----|------|------|-----|--------|---------|-------|
| *Caltech256* | 77.83% | 78.00% | 80.84% | 79.35% | 79.27% | 79.89% | **81.71%** |
| *Imagenet* | 70.29% | 72.16% | 74.89% | 72.53% | 72.33% | 74.73% | **75.00%** |
| *SUN* | 24.56% | 26.80% | 26.26% | 28.80% | 29.15% | 26.31% | **30.54%** |
| Average | 57.56% | 58.99% | 60.66% | 60.23% | 60.25% | 60.31% | **62.42%** |

that, in *Bing→Caltech256* and *Bing→ImageNet* tasks, ATDA is slightly worse than B-Net. However, in *Bing→SUN* task, ATDA is much worse than B-Net. The reason is that the DIR between *Bing* and *SUN* are more affected by noisy source data. This is also observed when comparing DANN and TCL. Compared to Co+ATDA, ATDA is slightly better than Co+ATDA. This abnormal phenomenon can be explained using $\Delta$ (see Section 5.2), after using co-teaching to assign pseudo labels to noisy source data, the second term in $\Delta_s$ may increase, which results in that $\Delta$ increases, i.e., noise effects actually increase. This phenomenon is an evidence that a two-step method may not really reduce noise effects.

### 9.8 Can we check correct data out?

This subsection verifies that $\rho_{01}^s$ and $\rho_{01}^t$ will go to zero with the convergence speed of $O(1/\sqrt{n_s T})$ and $O(1/\sqrt{n_t T})$, respectively. Figure 10 shows the values of $\rho_{01}^s$ and $\rho_{01}^t$. It can be seen that $\rho_{01}^s$ and $\rho_{01}^t$ will go to zero when increasing the training epochs. $\rho_{01}^s$ is always lower than $\rho_{01}^t$ because that $n_s$ is much larger than $n_t$, indicating that we can check more

correct data out when more samples are available. Figure 10-(b) shows that we can always find two finite $C_\rho^s$ such that $\rho_{01}^s$ goes to the zero with the convergence speed of $O(1/\sqrt{n_s T})$. So do $C_\rho^t$ and $\rho_{01}^t$ in Figure 10-(c).

### 9.9 Ablation study

Finally, we conduct thorough experiments to show the contribution of individual components in B-Net. We report average target-domain accuracy on $32$ simulated WUDA tasks (8 digit and 24 human-sentiment WUDA tasks) and 3 real-world WUDA tasks. We consider following baselines:

- Tri-C-Net: <u>tri</u>ply <u>c</u>heck data in SD, MD and TD. Compared to B-Net, Tri-C-Net has another branch (denoted by Branch-III) to check data in SD. Namely, Tri-C-Net has three branches (i.e., six networks). Parameters of CNN of the Branch-III are the same with that of Branch-I and Branch-II.
- B w/o C: train <u>B</u>-Net by Algorithm 2, <u>without</u> adding $|\theta_{f11}^T \theta_{f21}|$ into the loss function of B-Net.

(a) Values of $\rho_{01}^s$ and $\rho_{01}^t$     (b) Convergence speed of $\rho_{01}^s$     (c) Convergence speed of $\rho_{01}^t$
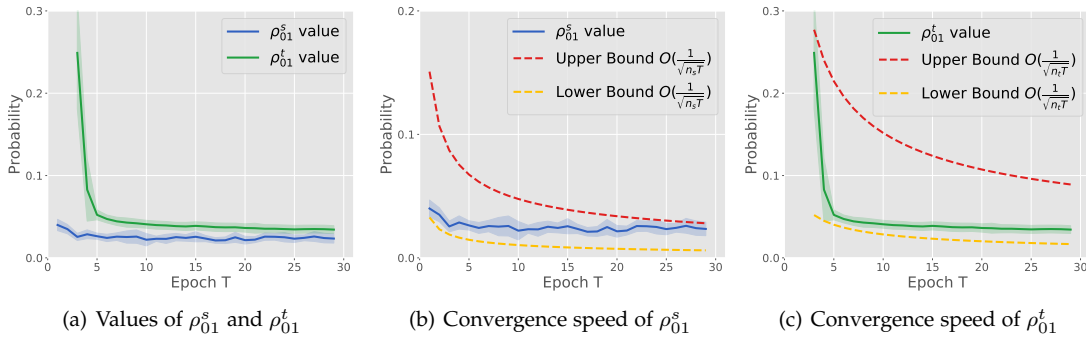
Fig. 10. The values of $\rho_{01}^s$ and $\rho_{01}^t$ on the $S{\rightarrow}M$ task (P20), where $n_s = 494,000$ is much larger than $n_t = 10,000$. Since we do not have pseudo-labeled target data in the first epoch, we illustrate $\rho_{01}^t$ from the second epoch.

TABLE 5
Results of ablation study. Average target-domain accuracy on 8 simulated digit WUDA tasks (*Digit*), 24 simulated human-sentiment WUDA tasks (*Sentiment*) and 3 real-world WUDA tasks (*Real-world*). Bold value represents the highest accuracy in each row.

| Datasets | Tri-C-Net | B w/o C | DCP-D | DCP-M | B-Net-S | B-Net-T | B-Net-ST | B-Net-M | B-Net |
|---|---|---|---|---|---|---|---|---|---|
| *Digit* | 59.80% | 74.52% | 59.19% | 70.85% | 71.93% | 52.00% | 72.27% | 73.89% | **75.82%** |
| *Sentiment* | 61.25% | 63.57% | 61.37% | 63.39% | 61.49% | 61.12% | 61.73% | 62.21% | **63.77%** |
| *Real-world* | 61.50% | 62.27% | 59.82% | 62.34% | 61.91% | 60.87% | 62.24% | 62.17% | **62.42%** |

- DCP-D: realize <u>DCP</u> via <u>D</u>ecoupling [57] to check data in MD and TD.
- DCP-M: realize <u>DCP</u> via <u>M</u>entorNet [54] to check data in MD and TD.
- B-Net-S: train <u>B-Net</u> where the check is turned on for <u>S</u>ource data in MD.
- B-Net-T: train <u>B-Net</u> where the check is turned on for <u>T</u>arget data in TD.
- B-Net-ST: train <u>B-Net</u> where the checks are turned on for <u>S</u>ource data in MD and <u>T</u>arget data in TD.
- B-Net-M: train <u>B-Net</u> where the check is turned on for all data in <u>M</u>D.

Note that in the full B-Net, the checks are turned on for all data in MD and TD. Comparing B-Net with Tri-C-Net shows whether two branches (i.e., four networks) are the optimal design. Comparing B-Net with B w/o C reveals if the constraint $|\theta_{f11}^T \theta_{f21}|$ takes effects. Comparing B-Net with DCP-D and $\theta_{f21}$M shows whether realizing DCP via co-teaching is the optimal way. Comparing B-Net with B-Net-S, B-Net-T, B-Net-ST and B-Net-M reveals if DCP is necessary.

Table 5 reports average target-domain accuracy of above baselines and B-Net. As can be seen, 1) maintaining 4 networks (like B-Net) is better than maintaining 6 networks (like Tri-C-Net) since B-Net outperforms Tri-C-Net in terms of average target-domain accuracy; 2) B-Net benefits from adding the constraint to the loss function $\mathcal{L}$; 3) realizing DCP by co-teaching is better than using Decoupling or MentorNet; and 4) DCP is necessary since accuracy of B-Net is higher than those of B-Net-S, B-Net-T, B-Net-ST and B-Net-M.

## 10 CONCLUSIONS

This paper opens a new problem called *wildly unsupervised domain adaptation* (WUDA). However, existing UDA methods cannot handle WUDA well. To address this problem, we propose a robust one-step approach called *Butterfly*. Butterfly maintains four deep networks simultaneously: Two take care of all adaptations; while the other two can focus on classification in target domain. We compare Butterfly with existing UDA methods on 32 simulated and 3 real-world WUDA tasks. Empirical results demonstrate that Butterfly can robustly transfer knowledge from noisy source data to unlabeled target data. In the future, we will extend our Butterfly framework to address open-set WUDA, where label space of target domain is larger than that of source domain.

## REFERENCES

[1] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan, "A theory of learning from different domains," *Machine Learning*, vol. 79, no. 1-2, pp. 151–175, 2010.

[2] Y. Ganin and V. S. Lempitsky, "Unsupervised domain adaptation by backpropagation," in *ICML*, 2015, pp. 1180–1189.

[3] M. Xiao and Y. Guo, "Feature space independent semi-supervised domain adaptation via kernel matching," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 1, pp. 54–66, 2015.

[4] K. Zhang, M. Gong, and B. Schölkopf, "Multi-source domain adaptation: A causal view," in *AAAI*, 2015, pp. 3150–3157.

[5] K. Zhang, B. Schölkopf, K. Muandet, and Z. Wang, "Domain adaptation under target and conditional shift," in *ICML*, 2013, pp. 819–827.

[6] P. Stojanov, M. Gong, J. G. Carbonell, and K. Zhang, "Data-driven approach to multiple-source domain adaptation," in *AISTATS*, vol. 89, 2019, pp. 3487–3496.

[7] E. Tzeng, J. Hoffman, T. Darrell, and K. Saenko, "Simultaneous deep transfer across domains and tasks," in *ICCV*, 2015, pp. 4068–4076.

[8] Y. Guo and M. Xiao, "Cross language text classification via subspace co-regularized multi-view learning," in *ICML*, 2012.

[9] W. Li, L. Duan, D. Xu, and I. W. Tsang, "Learning with augmented features for supervised and semi-supervised heterogeneous domain adaptation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 6, pp. 1134–1148, 2014.

[10] L. Duan, D. Xu, and I. Tsang, "Learning with augmented features for heterogeneous domain adaptation," in *ICML*, Edinburgh, UK, 2012, pp. 711–718.

[11] M. Xiao and Y. Guo, "Feature space independent semi-supervised domain adaptation via kernel matching," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 1, pp. 54–66, 2015.

[12] Y. Yan, W. Li, M. K. Ng, M. Tan, H. Wu, H. Min, and Q. Wu, "Learning discriminative correlation subspace for heterogeneous domain adaptation." in *IJCAI*, 2017, pp. 3252–3258.

[13] K. Saito, Y. Ushiku, and T. Harada, "Asymmetric tri-training for unsupervised domain adaptation," in *ICML*, 2017, pp. 2988–2997.

[14] M. Gong, K. Zhang, T. Liu, D. Tao, and C. Glymour, "Domain adaptation with conditional transferable components," in *ICML*, 2016, pp. 2839–2848.

[15] M. Long, Y. Cao, Z. Cao, J. Wang, and M. I. Jordan, "Transferable representation learning with deep adaptation networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. Early Access, pp. 1–14, 2018.

[16] R. Gopalan, R. Li, and R. Chellappa, "Unsupervised adaptation across domain shifts by generating intermediate data representations," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 11, pp. 2288–2302, 2014.

[17] W. Zhang, D. Xu, W. Ouyang, and W. Li, "Self-paced collaborative and adversarial network for unsupervised domain adaptation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. Early Access, pp. 1–15, 2019.

[18] F. Liu, J. Lu, and G. Zhang, "Unsupervised heterogeneous domain adaptation via shared fuzzy relations," *IEEE Transactions on Fuzzy Systems*, vol. 26, no. 6, pp. 3555–3568, 2018.

[19] M. Ghifary, W. B. Kleijn, M. Zhang, D. Balduzzi, and W. Li, "Deep reconstruction-classification networks for unsupervised domain adaptation," in *ECCV*, 2016, pp. 597–613.

[20] Z. Deng, Y. Luo, and J. Zhu, "Cluster alignment with a teacher for unsupervised domain adaptation," in *ICCV*, 2019, pp. 9944–9953.

[21] S. Motiian, Q. Jones, S. M. Iranmanesh, and G. Doretto, "Few-shot adversarial domain adaptation," in *NeurIPS*, 2017, pp. 6673–6683.

[22] W. Li, L. Chen, D. Xu, and L. Van Gool, "Visual recognition in rgb images and videos by learning from rgb-d data," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 8, pp. 2030–2036, 2017.

[23] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. S. Lempitsky, "Domain-adversarial training of neural networks," *Journal of Machine Learning Research*, vol. 17, pp. 59:1–59:35, 2016.

[24] B. Gong, Y. Shi, F. Sha, and K. Grauman, "Geodesic flow kernel for unsupervised domain adaptation," in *CVPR*, 2012, pp. 2066–2073.

[25] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada, "Maximum classifier discrepancy for unsupervised domain adaptation," in *CVPR*, 2018, pp. 3723–3732.

[26] R. Shu, H. H. Bui, H. Narui, and S. Ermon, "A DIRT-T approach to unsupervised domain adaptation," in *ICLR*, 2018.

[27] G. Kang, L. Jiang, Y. Yang, and A. G. Hauptmann, "Contrastive adaptation network for unsupervised domain adaptation," in *CVPR*, 2019, pp. 4893–4902.

[28] X. Ma, T. Zhang, and C. Xu, "GCAN: graph convolutional adversarial network for unsupervised domain adaptation," in *CVPR*, 2019, pp. 8266–8276.

[29] Y. Ziser and R. Reichart, "Task refinement learning for improved accuracy and stability of unsupervised domain adaptation," in *ACL*, 2019, pp. 5895–5906.

[30] J. Xu, S. Ramos, D. Vazquez, and A. M. Lopez, "Domain adaptation of deformable part-based models," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 12, pp. 2367–2380, 2014.

[31] R. Gong, W. Li, Y. Chen, and L. V. Gool, "Dlow: Domain flow for adaptation and generalization," in *CVPR*, 2019, pp. 2477–2486.

[32] H. Li, W. Li, H. Cao, S. Wang, F. Huang, and A. C. Kot, "Unsupervised domain adaptation for face anti-spoofing," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 7, pp. 1794–1809, 2018.

[33] M. Ghifary, D. Balduzzi, W. B. Kleijn, and M. Zhang, "Scatter component analysis : A unified framework for domain adaptation and domain generalization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 7, pp. 1414–1430, 2017.

[34] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. J. Smola, "A kernel two-sample test," *Journal of Machine Learning Research*, vol. 13, pp. 723–773, 2012.

[35] J. Lee and M. Raginsky, "Minimax statistical learning with wasserstein distances," in *NeurIPS*, 2018, pp. 2692–2701.

[36] M. Long, Y. Cao, J. Wang, and M. I. Jordan, "Learning transferable features with deep adaptation networks," in *ICML*, 2015, pp. 97–105.

[37] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Deep transfer learning with joint adaptation networks," in *ICML*, 2017, pp. 2208–2217.

[38] X. Yu, T. Liu, M. Gong, K. Zhang, K. Batmanghelich, and D. Tao, "Transfer learning with label noise," *CoRR*, vol. abs/1707.09724, 2017.

[39] M. Gong, K. Zhang, B. Huang, C. Glymour, D. Tao, and K. Batmanghelich, "Causal generative domain adaptation networks," *CoRR*, vol. abs/1804.04333, 2018.

[40] J. Hoffman, E. Tzeng, T. Park, J. Zhu, P. Isola, K. Saenko, A. A. Efros, and T. Darrell, "Cycada: Cycle-consistent adversarial domain adaptation," in *ICML*, 2018, pp. 1994–2003.

[41] Y. Li, X. Tian, M. Gong, Y. Liu, T. Liu, K. Zhang, and D. Tao, "Deep domain generalization via conditional invariant adversarial networks," in *ECCV*, 2018, pp. 647–663.

[42] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *CVPR*, 2017, pp. 2962–2971.

[43] W. Zhang, W. Ouyang, W. Li, and D. Xu, "Collaborative and adversarial network for unsupervised domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3801–3809.

[44] Y. Shu, Z. Cao, M. Long, and J. Wang, "Transferable curriculum for weakly-supervised domain adaptation," in *AAAI*, 2019, pp. 4951–4958.

[45] M. Tan, I. W. Tsang, and L. Wang, "Towards ultrahigh dimensional feature selection for big data," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1371–1429, 2014.

[46] K. Lee, X. He, L. Zhang, and L. Yang, "Cleannet: Transfer learning for scalable image classifier training with label noise," in *CVPR*, 2018, pp. 5447–5456.

[47] F. Schroff, A. Criminisi, and A. Zisserman, "Harvesting image databases from the web," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 4, pp. 754–766, 2011.

[48] T. Tommasi and T. Tuytelaars, "A testbed for cross-dataset analysis," in *ECCV TASK-CV Workshops*, 2014, pp. 18–31.

[49] X. Yu, T. Liu, M. Gong, K. Zhang, K. Batmanghelich, and D. Tao, "Transfer learning with label noise," *arXiv preprint arXiv:1707.09724*, 2017.

[50] F. Liu, W. Xu, J. Lu, G. Zhang, A. Gretton, and D. J. Sutherland, "Learning deep kernels for non-parametric two-sample tests," in *ICML*, 2020, pp. 6316–6326.

[51] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu, "Transfer joint matching for unsupervised domain adaptation," in *CVPR*, 2014, pp. 1410–1417.

[52] J. Liang, D. Hu, and J. Feng, "Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation," in *ICML*, 2020.

[53] T. Liu and D. Tao, "Classification with noisy labels by importance reweighting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 3, pp. 447–461, 2016.

[54] L. Jiang, Z. Zhou, T. Leung, L. Li, and F. Li, "Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels," in *ICML*, 2018, pp. 2309–2318.

[55] T. Xiao, T. Xia, Y. Yang, C. Huang, and X. Wang, "Learning from massive noisy labeled data for image classification," in *CVPR*, 2015, pp. 2691–2699.

[56] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. W. Tsang, and M. Sugiyama, "Co-teaching: Robust training of deep neural networks with extremely noisy labels," in *NeurIPS*, 2018, pp. 8527–8537.

[57] E. Malach and S. Shalev-Shwartz, "Decoupling "when to update" from "how to update"," in *NeurIPS*, 2017, pp. 961–971.

[58] G. Patrini, A. Rozza, A. K. Menon, R. Nock, and L. Qu, "Making deep neural networks robust to label noise: A loss correction approach," in *CVPR*, 2017, pp. 2233–2241.

[59] D. Arpit, S. Jastrzebski, N. Ballas, D. Krueger, E. Bengio, M. Kanwal, T. Maharaj, A. Fischer, A. Courville, and Y. Bengio, "A closer look at memorization in deep networks," in *ICML*, 2017.

[60] Y. Bengio, "Evolving culture versus local minima," in *Growing Adaptive Machines*, 2014, pp. 109–138.

[61] P. L. Bartlett and S. Mendelson, "Rademacher and gaussian complexities: Risk bounds and structural results," *Journal of Machine Learning Research*, vol. 3, pp. 463–482, 2002.

[62] Y. Mansour, M. Mohri, and A. Rostamizadeh, "Domain adaptation: Learning bounds and algorithms," in *COLT*, 2009, pp. 3:1–3:11.

[63] A. Maurer, "A vector-contraction inequality for rademacher complexities," in *ALT*, 2016, pp. 3–17.

[64] J. Li, Y. Liu, R. Yin, H. Zhang, L. Ding, and W. Wang, "Multi-class learning: From theory to algorithm," in *NeurIPS*, 2018, pp. 1593–1602.

[65] J. Li, Y. Liu, R. Yin, and W. Wang, "Multi-class learning using unlabeled samples: Theory and algorithm," in *IJCAI*, 2019, pp. 2880–2886.

[66] Y. Zhang, T. Liu, M. Long, and M. I. Jordan, "Bridging theory and algorithm for domain adaptation," in *ICML*, 2019, pp. 7404–7413.

[67] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of machine learning*. MIT press, 2018.

[68] R. Kiryo, G. Niu, M. C. du Plessis, and M. Sugiyama, "Positive-unlabeled learning with non-negative risk estimator," in *NeurIPS*, 2017.

[69] A. Bergamo and L. Torresani, "Exploiting weakly-labeled web images to improve object classification: a domain adaptation approach," in *NeurIPS*, 2010, pp. 181–189.

[70] G. Griffin, A. Holub, and P. Perona, "Caltech-256 object category dataset," California Institute of Technology, Tech. Rep., 2007.

[71] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F. Li, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009, pp. 248–255.

[72] J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba, "SUN database: Large-scale scene recognition from abbey to zoo," in *CVPR*, 2010, pp. 3485–3492.

[73] J. Wang, W. Feng, Y. Chen, H. Yu, M. Huang, and P. S. Yu, "Visual domain adaptation with manifold embedded distribution alignment," in *ACM MM*, S. Boll, K. M. Lee, J. Luo, W. Zhu, H. Byun, C. W. Chen, R. Lienhart, and T. Mei, Eds., 2018, pp. 402–410.

**Feng Liu** is a Doctoral candidate in Centre for Artificial intelligence, Faculty of Engineering and Information Technology, University of Technology Sydney, Australia. He received an M.Sc. degree in probability and statistics and a B.Sc. degree in pure mathematics from the School of Mathematics and Statistics, Lanzhou University, China, in 2015 and 2013, respectively. His research interests include domain adaptation and two-sample test. He has served as a senior program committee member for ECAI and program committee members for NeurIPS, ICML, IJCAI, CIKM, ECAI, FUZZ-IEEE and ISKE. He also served as rev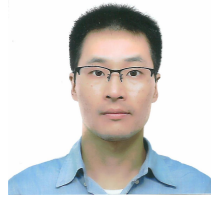iewers for TPAMI, TNNLS, TFS and TCYB. He has received the UTS-FEIT HDR Research Excellence Award (2019), Best Student Paper Award of FUZZ-IEEE (2019) and UTS Research Publication Award (2018).

**Jie Lu** (F'18) is a Distinguished Professor and the Director of the Centre for Artificial Intelligence at the University of Technology Sydney, Australia. She received the Ph.D. degree from Curtin University of Technology, Australia, in 2000.

Her main research expertise is in fuzzy transfer learning, decision support systems, concept drift, and recommender systems. She has published six research books and 400 papers in Artificial Intelligence, IEEE transactions on Fuzzy Systems and other refereed journals and conference proceedings. She has won over 20 Australian Research Council (ARC) discovery grants and other research grants for over $7 million. She serves as Editor-In-Chief for Knowledge-Based Systems (Elsevier) and Editor-In-Chief for International Journal on Computational Intelligence Systems (Atlantis), has delivered 20 keynote speeches at international conferences, and has chaired 10 international conferences. She is a Fellow of IEEE and Fellow of IFSA.

**Bo Han** is currently an Assistant Professor of Computer Science at Hong Kong Baptist University and a Visiting Scientist at RIKEN Center for Advanced Intelligence Project (RIKEN AIP), hosted by Masashi Sugiyama. He was a Postdoc Fellow at RIKEN AIP (2019-2020), advised by Masashi Sugiyama. He received his Ph.D. degree in Computer Science from University of Technology Sydney (2015-2019), advised by Ivor W. Tsang and Ling Chen. During 2018-2019, he was a Research Intern with the AI Residency Program at RIKEN AIP, working on robust deep learning projects with Masashi Sugiyama, Gang Niu and Mingyuan Zhou. His current research interests lie in machine learning, deep learning and artificial intelligence. His long-term goal is to develop trustworthy intelligent systems, which can learn from a massive volume of complex (e.g., weakly-supervised, adversarial, and private) data (e.g, single-/multi-label, ranking, domain, similarity, graph and demonstration) automatically. He has served as program committes of NeurIPS, ICML, ICLR, AISTATS, UAI, AAAI, IJCAI, ACML and ICDM. He received the National Scholarship (2013), UTS International Research Scholarship (2014) and UTS Research Publication Award (2017 and 2018).

**Gang Niu** is a research scientist at RIKEN Center for Advanced Intelligence Project. He received the PhD degree in computer science from Tokyo Institute of Technology in 2013. His research interests include mainly weakly-supervised learning and its applications. He has published 10 NeurIPS (including 1 oral and 1 spotlight) and 10 ICML papers and also served as an area chair for ICML 2019, NeurIPS 2019 and ICML 2020.

**Guangquan Zhang** is an Associate Professor and Director of the Decision Systems and e-Service Intelligent (DeSI) Research Laboratory at the University of Technology Sydney, Australia. He received the Ph.D. degree in applied mathematics from Curtin University of Technology, Australia, in 2001.

His research interests include fuzzy machine learning, fuzzy optimization, and machine learning. He has authored five monographs, five textbooks, and 460 papers including 220 refereed international journal papers. Dr. Zhang has won seven Australian Research Council (ARC) Discovery Projects grants and many other research grants. He was awarded an ARC QEII fellowship in 2005. He has served as a member of the editorial boards of several international journals, as a guest editor of eight special issues for IEEE transactions and other international journals, and co-chaired several international conferences and workshops in the area of fuzzy decision-making and knowledge engineering.

**Masashi Sugiyama** is Director of RIKEN Center for Advanced Intelligence Project and Professor at the University of Tokyo. He received the PhD degree in computer science from Tokyo Institute of Technology in 2001. His research interests include theories and algorithms of machine learning. He was awarded the Japan Society for the Promotion of Science Award and the Japan Academy Medal in 2017.

## APPENDIX A
## PROOFS

This section presents the completed proofs for theoretical results obtained in this paper. Since we have provided completed proofs regarding Theorems 3 and 4, we do not repeat them here.

### A.1 Proof of Theorem 1

*Proof.* We will fist prove Eq. (5) (Case 1) and then prove Eq. (6) (Case 2).

**Case** 1. According to definition of $\tilde{R}_s(h)$, we have

$$
\begin{aligned}
\tilde{R}_s(h) &= \mathbb{E}_{\tilde{p}_s(x_s, \tilde{y}_s)}[\ell(h(x_s), \tilde{y}_s)] \\
&= \int_{\mathcal{X}} \sum_{\tilde{y}_s=1}^{K} \ell(h(x_s), \tilde{y}_s) \tilde{p}_s(x_s, \tilde{y}_s) dx_s \\
&= \int_{\mathcal{X}} \sum_{\tilde{y}_s=1}^{K} \ell(h(x_s), \tilde{y}_s) \tilde{p}_{\tilde{Y}_s|X_s}(\tilde{y}_s|x_s) p_{x_s}(x_s) dx_s \\
&= \int_{\mathcal{X}} \tilde{\boldsymbol{\eta}}^T(x_s) \boldsymbol{\ell}(h(x_s)) p_{x_s}(x_s) dx_s,
\end{aligned}
\tag{32}
$$

where $\boldsymbol{\ell}(h(x_s)) = [\ell(h(x_s), 1), ..., \ell(h(x_s), K)]^T$ and $\tilde{\boldsymbol{\eta}}(x_s) = [\tilde{p}_{\tilde{Y}_s|X_s}(1|x_s), \ldots, \tilde{p}_{\tilde{Y}_s|X_s}(K|x_s)]^T$. According to definition of the transition matrix $Q$, we know that

$$
\tilde{\boldsymbol{\eta}}^T(x_s) = \boldsymbol{\eta}^T(x_s) Q,
\tag{33}
$$

where $\boldsymbol{\eta}(x_s) = [p_{Y_s|X_s}(1|x_s), \ldots, p_{Y_s|X_s}(K|x_s)]^T$. Substituting Eq. (33) into Eq. (32), we have

$$
\begin{aligned}
\tilde{R}_s(h) &= \int_{\mathcal{X}} \boldsymbol{\eta}^T(x_s) Q \boldsymbol{\ell}(h(x_s)) p_{x_s}(x_s) dx_s \\
&= \int_{\mathcal{X}} \boldsymbol{\eta}^T(x_s) I \boldsymbol{\ell}(h(x_s)) p_{x_s}(x_s) dx_s + \int_{\mathcal{X}} \boldsymbol{\eta}^T(x_s) (Q-I) \boldsymbol{\ell}(h(x_s)) p_{x_s}(x_s) dx_s \\
&= R_s(h) + \mathbb{E}_{p_{x_s}(x_s)}[\boldsymbol{\eta}^T(x_s)(Q-I)\boldsymbol{\ell}(h(x_s))].
\end{aligned}
$$

Hence, Case 1 is proved.

**Case** 2. According to definition of $\tilde{R}_s(h)$ and Eq. (2), we have

$$
\begin{aligned}
\tilde{R}_s(h) &= \mathbb{E}_{\tilde{p}_s(x_s, \tilde{y}_s)}[\ell(h(x_s), \tilde{y}_s)] \\
&= \int_{\mathcal{X}} \sum_{\tilde{y}_s=1}^{K} \ell(h(x_s), \tilde{y}_s) \tilde{p}_s(x_s, \tilde{y}_s) dx_s \\
&= \int_{\mathcal{X}} \sum_{y_s=1}^{K} \ell(h(x_s), y_s) \big((1-\rho)p_s(x_s, y_s) + \rho q_s(x_s, y_s)\big) dx_s \\
&= (1-\rho) \int_{\mathcal{X}} \sum_{y_s=1}^{K} \ell(h(x_s), y_s) p_s(x_s, y_s) dx_s + \rho \int_{\mathcal{X}} \sum_{y_s=1}^{K} \ell(h(x_s), y_s) q_s(x_s, y_s) dx_s \\
&= (1-\rho) R_s(h) + \rho \int_{\mathcal{X}} \sum_{y_s=1}^{K} \ell(h(x_s), y_s) q_{Y_s|X_s}(y_s|x_s) q_{x_s}(x_s) dx_s.
\end{aligned}
\tag{34}
$$

Let $\boldsymbol{\eta}_q(x_s) = [q_{Y_s|X_s}(1|x_s), ..., q_{Y_s|X_s}(K|x_s)]^T$, we have

$$
\tilde{R}_s(h) = (1-\rho) R_s(h) + \rho \mathbb{E}_{q_{x_s}(x_s)}[\boldsymbol{\eta}_q^T(x_s)\boldsymbol{\ell}(h(x_s))].
$$

Hence, Case 2 is proved. $\qquad\square$

### A.2 Proof of Theorem 2

*Proof.* For any $h \in \mathcal{H}$, we have

$$
\begin{aligned}
R_t(h, f_t) &= R_t(h, f_t) + \tilde{R}_s(h) - \tilde{R}_s(h) + R_s(h, f_t) - R_s(h, f_t) \\
&= \tilde{R}_s(h) + R_t(h, f_t) - \tilde{R}_s(h, f_t) + R_s(h, f_t) - R_s(h) + R_s(h) - \tilde{R}_s(h) + \tilde{R}_s(h, f_t) - R_s(h, f_t).
\end{aligned}
\tag{35}
$$

Since we do not know $f_t$, we substitute following equations into Eq. (35),

$$
\begin{aligned}
R_t(h, f_t) &= R_t(h, \tilde{f}_t) + R_t(h, f_t) - R_t(h, \tilde{f}_t), \\
\tilde{R}_s(h, f_t) &= \tilde{R}_s(h, \tilde{f}_t) + \tilde{R}_s(h, f_t) - \tilde{R}_s(h, \tilde{f}_t), \\
R_s(h, f_t) &= R_s(h, \tilde{f}_t) + R_s(h, f_t) - R_s(h, \tilde{f}_t).
\end{aligned}
$$

Then, we have

$$
\begin{aligned}
R_t(h, f_t) &= \tilde{R}_s(h) + R_t(h, \tilde{f}_t) - \tilde{R}_s(h, \tilde{f}_t) + R_s(h, \tilde{f}_t) - R_s(h) \\
&\quad + R_s(h) - \tilde{R}_s(h) + \tilde{R}_s(h, \tilde{f}_t) - R_s(h, \tilde{f}_t) + R_t(h, f_t) - R_t(h, \tilde{f}_t) \\
&\leq \tilde{R}_s(h) + |R_t(h, \tilde{f}_t) - \tilde{R}_s(h, \tilde{f}_t)| + |R_s(h, \tilde{f}_t) - R_s(h)| \\
&\quad + |\tilde{R}_s(h) - R_s(h)| + |\tilde{R}_s(h, \tilde{f}_t) - R_s(h, \tilde{f}_t)| + |R_t(h, f_t) - R_t(h, \tilde{f}_t)|.
\end{aligned}
$$

Hence, this theorem is proved.                                                                                                                                                  □

## A.3  Proof of Lemma 1

According to definition of $\tilde{R}_s^{\mathrm{po}}(h, u_s)$ in Section 6, we have

$$
\begin{aligned}
\tilde{R}_s^{\mathrm{po}}(h, u_s) &= \frac{\int_{\mathcal{X}} \sum_{u_s=0}^{1} \sum_{y_s=1}^{K} u_s \ell(h(x_s), y_s) \tilde{p}_s^{\mathrm{po}}(x_s, y_s, u_s) dx_s}{1 - \rho_{u_s}} \\
&= \frac{\int_{\mathcal{X}} \sum_{y_s=1}^{K} \ell(h(x_s), y_s) \tilde{p}_{X_s, Y_s|U_s}^{\mathrm{po}}(x_s, y_s|1) \tilde{p}_{U_s}^{\mathrm{po}}(1) dx_s}{1 - \rho_{u_s}} \\
&\overset{(a)}{=} \frac{1 - \rho_{u_s}}{1 - \rho_{u_s}} \int_{\mathcal{X}} \sum_{y_s=1}^{K} \ell(h(x_s), y_s) \big( \rho_{01}^s q_s(x_s, y_s) + \rho_{11}^s p_s(x_s, y_s) \big) dx_s \\
&= \rho_{01}^s \mathbb{E}_{q_s(x_s, y_s)} [\ell(h(x_s), y_s)] + \rho_{11}^s R_s(h),
\end{aligned}
$$

where $(a)$ is based on the definition of $\rho_{u_s}$ and Eq. (9). Thus, we have

$$
\begin{aligned}
|\tilde{R}_s^{\mathrm{po}}(h, u_s) - R_s(h)| &= |\rho_{01}^s \mathbb{E}_{q_s(x_s, y_s)} [\ell(h(x_s), y_s)] - (1 - \rho_{11}^s) R_s(h)| \\
&\leq \rho_{01}^s \max \{ \mathbb{E}_{q_s(x_s, y_s)} [\ell(h(x_s), y_s)], R_s(h) \}.
\end{aligned}
$$

This lemma is proved.

## A.4  Proof of Lemma 2

According to definition of $\tilde{R}_t^{\mathrm{po}}(h, \tilde{f}_t, u_t)$ in Section 6, we have

$$
\begin{aligned}
\tilde{R}_t^{\mathrm{po}}(h, \tilde{f}_t, u_t) &= (1 - \rho_{u_t})^{-1} \int_{\mathcal{X}} \sum_{u_t=0}^{1} u_t \ell(h(x_t), \tilde{f}_t(x_t)) \tilde{p}_t^{\mathrm{po}}(x_t, u_t) dx_t \\
&= (1 - \rho_{u_t})^{-1} \int_{\mathcal{X}} \ell(h(x_t), \tilde{f}_t(x_t)) \tilde{p}_{X_t|U_t}^{\mathrm{po}}(x_t|1) \tilde{p}_{U_t}^{\mathrm{po}}(1) dx_t \\
&\overset{(a)}{=} \frac{1 - \rho_{u_t}}{1 - \rho_{u_t}} \int_{\mathcal{X}} \ell(h(x_s), \tilde{f}_t(x_t)) \big( \rho_{01}^t q_{x_t}(x_t) + \rho_{11}^t p_{X_t|V_t}^{\mathrm{po}}(x_t|1) \big) dx_t \\
&= \rho_{01}^t \mathbb{E}_{q_{x_t}(x_t)} [\ell(h(x_t), \tilde{f}_t(x_t))] + \rho_{11}^t \int_{\mathcal{X}} \ell(h(x_t), \tilde{f}_t(x_t)) p_{X_t|V_t}^{\mathrm{po}}(x_t|V_t = 1) dx_t \\
&\overset{(b)}{=} \rho_{01}^t \mathbb{E}_{q_{x_t}(x_t)} [\ell(h(x_t), \tilde{f}_t(x_t))] + \rho_{11}^t \int_{\mathcal{X}} \ell(h(x_t), f_t(x_t)) p_{X_t|V_t}^{\mathrm{po}}(x_t|V_t = 1) dx_t \\
&= \rho_{01}^t \mathbb{E}_{q_{x_t}(x_t)} [\ell(h(x_t), \tilde{f}_t(x_t))] + \rho_{11}^t \int_{\mathcal{X}} \ell(h(x_t), f_t(x_t)) p'_{x_t}(x_t) dx_t \\
&= \rho_{01}^t \mathbb{E}_{q_{x_t}(x_t)} [\ell(h(x_t), \tilde{f}_t(x_t))] + \rho_{11}^t \mathbb{E}_{p'_{x_t}(x_t)} [\ell(h(x_t), f_t(x_t))],
\end{aligned}
\tag{36}
$$

where $(a)$ is based on the definition of $\rho_{u_s}$ and Eq. (9) and $(b)$ is based on the definition of $V_t$ ($f_t(x_t) = \tilde{f}_t(x_t)$ when $V_t = 1$). Since $\mathbb{E}_{p'_{x_t}(x_t)} [\ell(h(x_t), f_t(x_t))] \leq R_t(h, f_t) + \rho_{01}^s M_t$, we have

$$
\tilde{R}_t^{\mathrm{po}}(h, \tilde{f}_t, u_t) \leq \rho_{01}^t \mathbb{E}_{q_{x_t}(x_t)} [\ell(h(x_t), \tilde{f}_t(x_t))] + \rho_{11}^t (R_t(h, f_t) + \rho_{01}^s M_t).
\tag{37}
$$

Thus, we have

$$
\begin{aligned}
|\tilde{R}_t^{\mathrm{po}}(h, \tilde{f}_t, u_t) - R_t(h, f_t)| &= |\rho_{01}^t \mathbb{E}_{q_{x_t}(x_t)} [\ell(h(x_t), \tilde{f}_t(x_t))] + \rho_{11}^t \mathbb{E}_{p'_{x_t}(x_t)} [\ell(h(x_t), f_t(x_t))] - R_t(h, f_t)| \\
&\leq |\rho_{01}^t \mathbb{E}_{q_{x_t}(x_t)} [\ell(h(x_t), \tilde{f}_t(x_t))] + \rho_{11}^t (R_t(h, f_t) + \rho_{01}^s M_t) - R_t(h, f_t)| \\
&= |\rho_{01}^t (\mathbb{E}_{q_{x_t}(x_t)} [\ell(h(x_t), \tilde{f}_t(x_t))] - R_t(h, f_t)) + \rho_{11}^t \rho_{01}^s M_t| \\
&\leq \rho_{01}^t \max \{ \mathbb{E}_{q_{x_t}(x_t)} [\ell(h(x_t), \tilde{f}_t(x_t))], R_t(h, f_t) \} + \rho_{11}^t \rho_{01}^s M_t.
\end{aligned}
$$

This lemma is proved.

## A.5 Proof of Lemma 3

For simplicity, in this proof, we let $\mathcal{L}_{S_s}(\ell, h) = \mathcal{L}(\theta, h; \boldsymbol{u}_s, D_s^{xy})$, $\tilde{R}_s^{\mathrm{po}}(\ell, h) = \tilde{R}_s^{\mathrm{po}}(h, u_s)$, and $\mathbb{E}_{S_s}[\cdot] = \mathbb{E}_{S_s \sim (\tilde{P}_s^{\mathrm{po}})^n}[\cdot]$, where $\tilde{P}_s^{\mathrm{po}}$ is the probability measure corresponding to the density $\tilde{p}_s^{\mathrm{po}}$. We first show that $\mathcal{L}_{S_s}(\ell, h)$ is an unbiased estimator of $\tilde{R}_s^{\mathrm{po}}(\ell, h)$ based on the definition of $\tilde{R}_s^{\mathrm{po}}(h, u_s)$ in Section 6. Since $S_s = \{(x_{si}, y_{si}, u_{si})\}_{i=1}^n$ are i.i.d samples from $\tilde{P}_s^{\mathrm{po}}$, $\mathbb{E}_{S_s}[\mathcal{L}_{S_s}(\ell, h)]$ can be expressed as follows.

$$\mathbb{E}_{S_s}\left[\frac{1}{\sum_{i=1}^n u_{si}} \sum_{i=1}^n u_{si} \ell(h(x_{si}), y_{si})\right]$$

$$= \int_{\mathcal{X}} \frac{1}{\sum_{i=1}^n u_{si}} \sum_{i=1}^n \sum_{u_{si}=0}^1 \sum_{y_{si}=1}^K u_{si} \ell(h(x_{si}), y_{si}) d\tilde{P}_s^{\mathrm{po}}$$

$$= \frac{1}{n} \int_{\mathcal{X}} \frac{n}{\sum_{i=1}^n U_i} \sum_{i=1}^n \sum_{u_{si}=0}^1 \sum_{y_{si}=1}^K u_{si} \ell(h(x_{si}), y_{si}) d\tilde{P}_s^{\mathrm{po}}$$

$$= \frac{1}{n} \int_{\mathcal{X}} (1 - \rho_{u_s})^{-1} \sum_{i=1}^n \sum_{u_{si}=0}^1 \sum_{y_{si}=1}^K u_{si} \ell(h(x_{si}), y_{si}) d\tilde{P}_s^{\mathrm{po}}$$

$$= (1 - \rho_{u_s})^{-1} \frac{1}{n} \sum_{i=1}^n \int_{\mathcal{X}} \sum_{u_{si}=0}^1 \sum_{y_{si}=1}^K u_{si} \ell(h(x_{si}), y_{si}) d\tilde{P}_s^{\mathrm{po}}$$

$$= (1 - \rho_{u_s})^{-1} \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\tilde{p}_s^{\mathrm{po}}(x_s, y_s, u_s)}[u_s \ell(h(x_s), y_s)]$$

$$= \tilde{R}_s^{\mathrm{po}}(h, u_s) = \tilde{R}_s^{\mathrm{po}}(\ell, h), \tag{38}$$

which means that $\mathcal{L}_{S_s}(\ell, h)$ is an unbiased estimator of $\tilde{R}_s^{\mathrm{po}}(\ell, h)$. Then, let $\Phi(S_s) = \sup_{\ell \in \mathbb{L}_{\mathcal{H}}}\left(\tilde{R}_s^{\mathrm{po}}(\ell, h) - \mathcal{L}_{S_s}(\ell, h)\right)$. Changing a point of $S_s$ affects $\Phi(S_s)$ at most $C_L/(n(1 - \tau_s))$. Thus, by McDiarmid's inequality applied to $\Phi(S_s)$, for any $\delta > 0$, with probability of at least $1 - \delta/2$, the following inequality holds.

$$\Phi(S_s) \le \mathbb{E}_{S_s}[\Phi(S_s)] + \frac{C_L}{1 - \tau_s}\sqrt{\frac{\ln(\delta/2)}{2n}}. \tag{39}$$

Then, we have

$$\mathbb{E}_{S_s}[\Phi(S_s)] = \mathbb{E}_{S_s}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}}\left(\tilde{R}_s^{\mathrm{po}}(h, u_s) - \mathcal{L}_{S_s}(h)\right)\right]$$

$$= \mathbb{E}_{S_s}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}}\left(\mathbb{E}_{S_s'}[\mathcal{L}_{S_s'}(\ell, h)] - \mathcal{L}_{S_s}(h)\right)\right] \tag{40}$$

$$= \mathbb{E}_{S_s}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}}\left(\mathbb{E}_{S_s'}[\mathcal{L}_{S_s'}(\ell, h) - \mathcal{L}_{S_s}(h)]\right)\right]$$

$$\le \mathbb{E}_{S_s, S_s'}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}}\left(\mathcal{L}_{S_s'}(\ell, h) - \mathcal{L}_{S_s}(h)\right)\right] \tag{41}$$

$$= \frac{1}{n}\mathbb{E}_{S_s, S_s'}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}} \sum_{i=1}^n \left(\frac{u_{si}' \ell(h(x_{si}'), y_{si}') - u_{si}\ell(h(x_{si}), y_{si}))}{1 - \tau_s}\right)\right]$$

$$= \frac{1}{n}\mathbb{E}_{\sigma, S_s, S_s'}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}} \sum_{i=1}^n \sigma_i\left(\frac{u_{si}' \ell(h(x_{si}'), y_{si}') - u_{si}\ell(h(x_{si}), y_{si}))}{1 - \tau_s}\right)\right]$$

$$\le \frac{1}{n(1 - \tau_s)}\mathbb{E}_{\sigma, S_s'}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}} \sum_{i=1}^n \sigma_i u_{si}' \ell(h(x_{si}'), y_{si}')\right] + \frac{1}{n(1 - \tau_s)}\mathbb{E}_{\sigma, S_s}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}} \sum_{i=1}^n -\sigma_i u_{si}\ell(h(x_{si}), y_{si})\right] \tag{42}$$

$$= \frac{2}{n(1 - \tau_s)}\mathbb{E}_{\sigma, S_s}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}} \sum_{i=1}^n \sigma_i u_{si}\ell(h(x_{si}), y_{si})\right], \tag{43}$$

where Eq. (40) is based on Eq. (38), Inequalities (41) and (42) are based on Jensen's Inequality. Because of existence of $u_{si}$, Eq. (43) is not the Rademacher complexity of $\mathbb{L}_{\mathcal{H}}$ (i.e., $\Re(\mathbb{L}_{\mathcal{H}})$). However, in following, we prove that Eq. (43) can be bounded by $\Re(\mathbb{L}_{\mathcal{H}})/(1 - \tau_s)$.

$$\mathbb{E}_{\sigma, S_s}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}} \sum_{i=1}^n \sigma_i u_{si}\ell(h(x_{si}), y_{si})\right]$$

$$= \mathbb{E}_{\sigma, S_s}\left[\sup_{\ell \in \mathbb{L}_{\mathcal{H}}}\left(\sigma_1 u_{s1}\ell(h(x_{s1}), y_{s1}) + \sum_{i=2}^n \sigma_i u_{si}\ell(h(x_{si}), y_{si})\right)\right]$$

$$= \frac{1}{2}\mathbb{E}_{\sigma, S_s}\left[\sup_{\ell, \ell' \in \mathbb{L}_{\mathcal{H}}}\left(u_{s1}\ell(h(x_{s1}), y_{s1}) + \sum_{i=2}^n \sigma_i u_{si}\ell(h(x_{si}), y_{si}) + (-u_{s1})\ell'(h(x_{s1}), y_{s1}) + \sum_{i=2}^n \sigma_i u_{si}\ell'(h(x_{si}), y_{si})\right)\right]$$

$$= \frac{1}{2}\mathbb{E}_{\sigma,S_s}\Big[\sup_{\ell,\ell'\in\mathbb{L}_{\mathcal{H}}}\Big(u_{s1}(\ell(h(x_{s1}),y_{s1})-\ell'(h(x_{s1}),y_{s1}))+\sum_{i=2}^{n}\sigma_i u_{si}\ell(h(x_{si}),y_{si})+\sum_{i=2}^{n}\sigma_i u_{si}\ell'(h(x_{si}),y_{si})\Big)\Big]$$

$$\leq \frac{1}{2}\mathbb{E}_{\sigma,S_s}\Big[\sup_{\ell,\ell'\in\mathbb{L}_{\mathcal{H}}}\Big(\ell(h(x_{s1}),y_{s1})-\ell'(h(x_{s1}),y_{s1})+\sum_{i=2}^{n}\sigma_i u_{si}\ell(h(x_{si}),y_{si})+\sum_{i=2}^{n}\sigma_i u_{si}\ell'(h(x_{si}),y_{si})\Big)\Big] \qquad (44)$$

$$= \mathbb{E}_{\sigma,S_s}\Big[\sup_{\ell\in\mathbb{L}_{\mathcal{H}}}\Big(\sigma_1\ell(h(x_{s1}),y_{s1})+\sum_{i=2}^{n}\sigma_i u_{si}\ell(h(x_{si}),y_{si})\Big)\Big],$$

where Inequality (44) is based on the fact that there are always $\ell,\ell'\in\mathbb{L}_{\mathcal{H}}$ such that $\ell(h(x_{s1}),y_{s1})-\ell'(h(x_{s1}),y_{s1}) > 0$. Repeat above procedures $n-1$ times, we have

$$\frac{2}{n}\mathbb{E}_{\sigma,S_s}\Big[\sup_{\ell\in\mathbb{L}_{\mathcal{H}}}\sum_{i=1}^{n}\sigma_i u_{si}\ell(h(x_{si}),y_{si})\Big] \leq \frac{2}{n}\mathbb{E}_{\sigma,S_s}\Big[\sup_{\ell\in\mathbb{L}_{\mathcal{H}}}\sum_{i=1}^{n}\sigma_i\ell(h(x_{si}),y_{si})\Big] := \Re_n(\mathbb{L}_{\mathcal{H}}). \qquad (45)$$

Changing a point of $S_s$ affects $\Re_n(\mathbb{L}_{\mathcal{H}})$ at most $2C_L/n$. Thus, by McDiarmid's inequality, for any $\delta > 0$, with probability of at least $1-\delta/2$, the following inequality holds.

$$\Re_n(\mathbb{L}_{\mathcal{H}}) \leq \hat{\Re}_{S_s}(\mathbb{L}_{\mathcal{H}}) + 2C_L\sqrt{\frac{\ln(\delta/2)}{2n}}. \qquad (46)$$

Since $\ell$ is Lipschitz continuous, according to [63], we have

$$\hat{\Re}_{S_s}(\mathbb{L}_{\mathcal{H}}) \leq \sqrt{2}L_\ell\hat{\Re}_{D_s^x}(\mathcal{H}). \qquad (47)$$

Combining (39), (43), (45), (46) and (47), we prove this lemma.

## A.6 Proof of Corollary 1

We prove this corollary (i.e., Inequality (28)) according to Inequality (21), where (28) has 8 terms in the right side and (21) have 6 terms in the right side.

1) For last 3 terms in (21), since $\rho_{01}^s < C_\rho^s/\sqrt{n_sT}$ and $\rho_{01}^t < C_\rho^t/\sqrt{n_tT}$, according to (16), (17) and (18), we know the sum of last three terms of (21) is less than or equal to $C_\rho^s(M_s+M_t)/\sqrt{n_sT} + 2C_\rho^t M_t/\sqrt{n_tT}$ (i.e., the last 2 terms in (28)).

2) For first 3 terms in (21), we have shown that (in Section 7.1) the sum of the first 3 terms in (21) is less than or equal to $(*)$:

$$2\tilde{R}_s^{\text{po}}(h,u_s) + 2\tilde{R}_s^{\text{po}}(h,\tilde{f}_t,\boldsymbol{u}_s) + \tilde{R}_t^{\text{po}}(h,\tilde{f}_t,\boldsymbol{u}_t) + \frac{C_\rho^s M_s}{\sqrt{n_sT}} + \frac{C_\rho^t M_t}{\sqrt{n_tT}}.$$

Then, we can prove that (similar with Lemma 3), with probability of at least $1-\delta$, for any $h\in\mathcal{H}$,

$$\tilde{R}_s^{\text{po}}(h,\tilde{f}_t,\boldsymbol{u}_s) \leq \mathcal{L}(\theta,h;\boldsymbol{u}_s,D_{\tilde{s}}^{xy}) + \frac{\sqrt{2}L_\ell\hat{\Re}_{D_{\tilde{s}}^x}(\mathcal{H})}{1-\tau_s} + \frac{3C_L}{1-\tau_s}\sqrt{\frac{\ln\frac{\delta}{2}}{2n_s}}, \qquad (48)$$

$$\tilde{R}_t^{\text{po}}(h,\tilde{f}_t,\boldsymbol{u}_t) \leq \mathcal{L}(\theta,h;\boldsymbol{u}_t,D_{\tilde{t}}^{xy}) + \frac{\sqrt{2}L_\ell\hat{\Re}_{D_{\tilde{t}}^x}(\mathcal{H})}{1-\tau_s} + \frac{3C_L}{1-\tau_t}\sqrt{\frac{\ln\frac{\delta}{2}}{2n_t}}. \qquad (49)$$

Combining (23), (48), (49) with $(*)$, we get the first 6 terms in (25). Hence we obtain all 6 terms in (25) and prove this corollary.

# APPENDIX B
## ADDITIONAL EXPERIMENTAL RESULTS

In this section, we present the standard deviation of target-domain accuracy of all methods on WUDA tasks.

TABLE 6
The standard deviation of target-domain accuracy on $8$ digit WUDA tasks ($SYND\leftrightarrow MNIST$). Bold value represents the highest accuracy in each row.

| Tasks | Type | DAN | DANN | ATDA | TCL | Co+TCL | Co+ATDA | B-Net |
|-------|------|-----|------|------|-----|--------|---------|-------|
| $S\rightarrow M$ | P20 | 0.23% | 1.12% | 31.26% | 3.88% | 3.26% | 0.66% | 0.50% |
| | P45 | 6.43% | 6.88% | 6.45% | 7.08% | 6.45% | 4.02% | 5.43% |
| | S20 | 1.18% | 1.29% | 1.32% | 1.17% | 1.32% | 0.38% | 0.23% |
| | S45 | 1.38% | 1.59% | 1.64% | 1.62% | 1.64% | 1.29% | 1.13% |
| $M\rightarrow S$ | P20 | 4.30% | 4.59% | 4.62% | 4.54% | 4.62% | 2.73% | 4.31% |
| | P45 | 2.01% | 2.05% | 2.06% | 1.87% | 2.06% | 6.81% | 4.06% |
| | S20 | 4.82% | 4.84% | 4.88% | 4.70% | 4.88% | 3.20% | 2.66% |
| | S45 | 2.02% | 2.25% | 2.25% | 2.22% | 2.25% | 1.68% | 2.79% |
| Average | | 2.80% | 3.08% | 6.81% | 3.39% | 3.31% | 2.60% | 2.64% |

TABLE 7
The standard deviation of target-domain accuracy on $12$ human-sentiment WUDA tasks with the $20\%$ noise rate. Bold values mean the highest values in each row.

| Tasks | DAN | DANN | ATDA | TCL | MEDA | Co+TCL | Co+ATDA | B-Net |
|---|---|---|---|---|---|---|---|---|
| $B{\to}D$ | 1.48% | 1.37% | 1.45% | 1.41% | 1.40% | 1.38% | 1.47% | 1.47% |
| $B{\to}E$ | 1.82% | 1.67% | 1.81% | 1.77% | 1.76% | 1.77% | 1.68% | 1.81% |
| $B{\to}K$ | 1.34% | 1.33% | 1.33% | 0.97% | 1.31% | 1.21% | 1.24% | 1.33% |
| $D{\to}B$ | 1.84% | 1.50% | 1.83% | 1.78% | 1.83% | 1.68% | 1.79% | 1.63% |
| $D{\to}E$ | 1.78% | 1.72% | 1.75% | 1.77% | 1.74% | 1.66% | 1.72% | 1.77% |
| $D{\to}K$ | 2.02% | 1.98% | 2.00% | 1.81% | 1.97% | 1.96% | 1.88% | 1.90% |
| $E{\to}B$ | 1.54% | 1.42% | 1.52% | 1.22% | 1.53% | 1.45% | 1.51% | 1.53% |
| $E{\to}D$ | 1.72% | 1.65% | 1.71% | 1.48% | 1.67% | 1.79% | 1.53% | 1.70% |
| $E{\to}K$ | 1.29% | 1.12% | 1.27% | 1.22% | 1.27% | 1.15% | 1.14% | 1.28% |
| $K{\to}B$ | 1.86% | 1.74% | 1.84% | 1.72% | 1.82% | 1.72% | 1.63% | 1.85% |
| $K{\to}D$ | 0.44% | 0.11% | 0.42% | 0.27% | 0.39% | 0.31% | 0.21% | 0.43% |
| $K{\to}E$ | 1.00% | 0.68% | 0.98% | 0.64% | 0.98% | 0.96% | 0.79% | 0.99% |
| Average | 1.51% | 1.36% | 1.49% | 1.34% | 1.47% | 1.42% | 1.38% | 1.48% |

TABLE 8
The standard deviation of target-domain accuracy on $12$ human-sentiment WUDA tasks with the $45\%$ noise rate. Bold values mean the highest values in each row.

| Tasks | DAN | DANN | ATDA | TCL | MEDA | Co+TCL | Co+ATDA | B-Net |
|---|---|---|---|---|---|---|---|---|
| $B{\to}D$ | 1.11% | 0.83% | 0.92% | 1.11% | 1.11% | 1.07% | 0.97% | 0.88% |
| $B{\to}E$ | 2.92% | 2.86% | 2.37% | 2.57% | 2.90% | 2.90% | 2.85% | 2.69% |
| $B{\to}K$ | 2.12% | 1.95% | 1.89% | 1.90% | 2.11% | 2.03% | 1.76% | 1.91% |
| $D{\to}B$ | 1.81% | 1.71% | 1.26% | 1.28% | 1.81% | 1.70% | 1.54% | 1.52% |
| $D{\to}E$ | 1.71% | 1.52% | 1.14% | 1.71% | 1.70% | 1.62% | 1.43% | 1.55% |
| $D{\to}K$ | 1.91% | 1.62% | 1.86% | 1.65% | 1.90% | 1.85% | 1.51% | 1.74% |
| $E{\to}B$ | 1.37% | 1.02% | 1.16% | 1.12% | 1.36% | 1.26% | 0.90% | 1.24% |
| $E{\to}D$ | 1.53% | 1.23% | 1.35% | 1.27% | 1.51% | 1.43% | 1.32% | 1.23% |
| $E{\to}K$ | 1.29% | 0.71% | 0.75% | 0.85% | 1.28% | 1.18% | 0.89% | 1.05% |
| $K{\to}B$ | 2.26% | 1.92% | 1.58% | 2.08% | 2.24% | 2.16% | 2.01% | 2.06% |
| $K{\to}D$ | 2.86% | 2.23% | 2.58% | 2.41% | 2.85% | 2.70% | 2.35% | 2.62% |
| $K{\to}E$ | 1.89% | 1.46% | 1.25% | 1.62% | 1.86% | 1.85% | 1.38% | 1.67% |
| Average | 1.90% | 1.59% | 1.51% | 1.63% | 1.88% | 1.81% | 1.58% | 1.68% |

TABLE 9
The standard deviation of target-domain accuracy on $3$ real-world WUDA tasks. The source domain is the *Bing* dataset that contains noisy information from the Internet. Bold value represents the highest accuracy in each row.

| Target | DAN | DANN | ATDA | TCL | Co+TCL | Co+ATDA | B-Net |
|---|---|---|---|---|---|---|---|
| *Caltech256* | 0.65% | 0.52% | 0.60% | 0.48% | 0.61% | 0.34% | 0.36% |
| *Imagenet* | 0.32% | 0.24% | 0.29% | 0.21% | 0.26% | 0.34% | 0.51% |
| *SUN* | 1.61% | 1.51% | 1.61% | 1.55% | 1.59% | 1.87% | 1.46% |
| Average | 0.86% | 0.75% | 0.83% | 0.75% | 0.82% | 0.85% | 0.78% |