

# Open-Set Heterogeneous Domain Adaptation: Theoretical Analysis and Algorithm

Thai-Hoang Pham<sup>1,2</sup>, Yuanlong Wang<sup>1,2</sup>, Changchang Yin<sup>1,2</sup>, Xueru Zhang<sup>1</sup>, Ping Zhang<sup>1,2</sup>

<sup>1</sup>Department of Computer Science and Engineering, The Ohio State University, USA

<sup>2</sup>Department of Biomedical Informatics, The Ohio State University, USA  
{pham.375,wang.16050,yin.731,zhang.12807,zhang.10631}@osu.edu

## Abstract

Domain adaptation (DA) tackles the issue of distribution shift by learning a model from a source domain that generalizes to a target domain. However, most existing DA methods are designed for scenarios where the source and target domain data lie within the same feature space, which limits their applicability in real-world situations. Recently, heterogeneous DA (HeDA) methods have been introduced to address the challenges posed by heterogeneous feature space between source and target domains. Despite their successes, current HeDA techniques fall short when there is a mismatch in both feature and label spaces. To address this, this paper explores a new DA scenario called open-set HeDA (OSHeDA). In OSHeDA, the model must not only handle heterogeneity in feature space but also identify samples belonging to novel classes. To tackle this challenge, we first develop a novel theoretical framework that constructs learning bounds for prediction error on target domain. Guided by this framework, we propose a new DA method called Representation Learning for OSHeDA (RL-OSHeDA). This method is designed to simultaneously transfer knowledge between heterogeneous data sources and identify novel classes. Experiments across text, image, and clinical data demonstrate the effectiveness of our algorithm. Model implementation is available at <https://github.com/pth1993/OSHeDA>.

## 1 Introduction

Machine learning (ML) techniques have achieved unprecedented success over the past decades in numerous areas (LeCun, Bengio, and Hinton 2015). However, ML systems are often built on the assumption that training and testing data are independent and identically distributed, which is commonly violated in real-world applications where the environment changes during model deployment. Existing works have shown that the performance of ML models often deteriorates due to distribution shifts between training and testing data (Ben-David et al. 2010; Quiñero-Candela et al. 2022). To learn a model robust under distribution shifts, domain adaptation (DA) (Ben-David et al. 2010; Mansour, Mohri, and Rostamizadeh 2009) has been proposed to transfer knowledge from a source domain that possesses abundant labeled data to a different but relevant target domain.

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

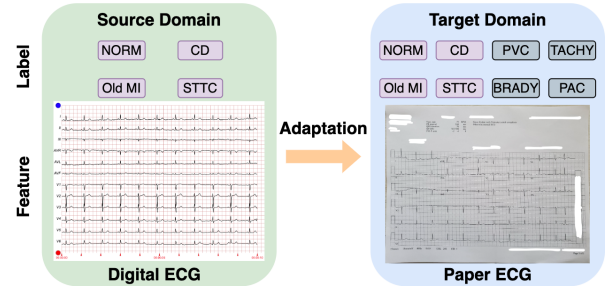


Figure 1: A motivating example about OSHeDA in the context of screening diseases using electrocardiogram (ECG) data. While digital ECGs comprise the majority of labeled data for training ML models for disease screening, physical or paper ECGs remain prevalent worldwide. Thus, the transfer of knowledge from digital ECG datasets is essential to support the training of ML models that analyze paper ECGs. Moreover, ML systems must effectively manage rare abnormalities (indicated with gray boxes), which may not be available in training data, to prevent misdiagnosis.

Existing DA methods, however, typically assume a homogeneous scenario where the source and target domains have the same feature and label spaces. Consequently, they may fail when the source and target domain data lie in different spaces. For example, heterogeneous feature space is common in biomedical domains in which medical terms undergoes continuous evolution, leading to the retirement of outdated terms (e.g., ICD-9 coding system) and the introduction of novel ones (e.g., ICD-10 coding system) (Grief et al. 2016). In such cases, acquiring training data that seamlessly aligns with target domain’s feature space can be impractical or excessively costly. Heterogeneous domain adaptation (HeDA) methods have emerged to handle the heterogeneity observed in distinct feature spaces, which often vary significantly between domains (Li et al. 2020; Zhao et al. 2022).

Despite the significant successes achieved by these HeDA methods, they face a major limitation: current HeDA techniques can only address heterogeneity in feature space and are inadequate when there is a mismatch in both feature and label spaces. This limitation restricts the practical application of HeDA methods in many real-world scenarios because neglecting label mismatch, such as new classes emerging in

the target domain, can lead to negative transfer effects from the source to the target domains (Liu et al. 2019).

To overcome this limitation, this study explores a new DA scenario called open-set heterogeneous domain adaptation (OSHeDA). In OSHeDA, ML methods must not only manage heterogeneity in feature space between source and target domains but also identify samples belonging to novel classes in the target domain. Figure 1 illustrates a real example from clinical applications for this novel learning scenario. In this instance, the adaptation process aims to transfer knowledge from digital electrocardiogram (ECG) to paper ECG formats (*heterogeneous*). Moreover, ML models for ECG-based diagnosis must also detect rare abnormalities that were not included in the training data (*open-set*).

To address the challenge of feature and label mismatch in OSHeDA, we first develop a novel theoretical analysis that constructs learning bounds for the prediction error of ML models on the target domain. Guided by this theoretical analysis, we then design a novel representation learning method named **Representation Learning for Open-Set Heterogeneous Domain Adaptation (RL-OSHeDA)**. This method is proposed to transfer knowledge between heterogeneous data sources and identify novel class simultaneously. Unlike existing HeDA methods, RL-OSHeDA transfer knowledge from source to target domains by aligning representations between source and target domains for known classes while also enforcing the representations of novel class in target domains to move apart from the known classes of the source and target domains. To effectively identify samples from novel class within unlabeled data, RL-OSHeDA optimizes a non-negative risk estimator for open-set and employs pseudo labeling to enrich the labeled data.

In summary, the contributions of our work are as follows:

- We conduct a theoretical analysis to establish learning bounds in the OSHeDA scenario. This analysis emphasizes the importance of minimizing the distance between source and target domains for known classes, while maximizing the separation from unknown classes. Moreover, we investigate the impact of pseudo-label and the non-negative risk estimator for open-set in OSHeDA.
- Motivated by the theoretical results, we propose a novel algorithm (RL-OSHeDA) based on representation learning to transfer knowledge from source to target domains.
- We conduct experiments on real data from clinical, computer vision, and natural language processing domains to validate the effectiveness of our method for OSHeDA.

## 2 Related Works

In this section, we summarize existing research from related areas including heterogeneous domain adaptation, open-set domain adaptation, and open-set semi-supervised learning.

**Heterogeneous Domain Adaptation (HeDA).** HeDA aims to transfer knowledge across domains with distinct feature spaces and data distributions. Depending on whether unlabeled target data are used in the adaptation process, HeDA approaches are categorized into three types: supervised, semi-supervised, and unsupervised methods. Supervised HeDA methods utilize ample labeled data from both

source and target domains for adaptation (Hoffman et al. 2013; Li et al. 2013; Hoffman et al. 2014). In contrast, semi-supervised HeDA methods require only a small number of labeled target domain data and utilize unlabeled instances from the target domain to facilitate transfer (Yao et al. 2020; Li et al. 2020; Fang et al. 2022; Zhao et al. 2022; Yao et al. 2019). Finally, unsupervised HeDA methods operate without any labeled target data, relying solely on unlabeled instances and labeled source data to align cross-domain feature representations (Shen and Guo 2018; Li et al. 2018; Zou et al. 2018). However, successful transfer in unsupervised settings depends on specific assumptions about domain relationships (Liu, Zhang, and Lu 2020).

**Open-Set Domain Adaptation (OSDA).** OSDA represents a realistic and challenging scenario in DA where the target domain includes instances whose classes are not observed in the source domain, alongside a shift in feature distribution between the two domains. In contrast to the OSHeDA, OSDA assumes a homogeneous feature space between the source and target domains. Existing approaches for OSDA can be categorized into two main groups: adversarial learning and self-supervised learning. Adversarial learning methods employ adversarial networks to detect unknown samples and align the distributions of known samples between the source and target domains (Saito et al. 2018; Luo et al. 2020). On the other hand, self-supervised learning methods utilize techniques like data augmentation to distinguish between known and unknown instances in the target domain (Bucci, Loghmani, and Tommasi 2020; Li et al. 2021).

**Open-Set Semi-supervised Learning (OS-SSL).** OS-SSL is a SSL scenario that addresses novel classes within unlabeled data during training. Unlike OSDA, OS-SSL requires only a small amount of labeled data. However, it assumes that both labeled and unlabeled data of known classes are drawn from the same distribution, and this setting does not account for novel classes during inference. Methods designed for OS-SSL can be broadly categorized into two types based on how they detect novel classes: criterion-based approaches and detector-based approaches. Criterion-based approaches use heuristic rules to identify novel classes (Chen et al. 2020; Huang, Yang, and Gong 2022; Du et al. 2023; He et al. 2022). In contrast, detector-based approaches employ parameterized detectors to filter outliers (Yu et al. 2020; Huang et al. 2021; Wang et al. 2023; Saito, Kim, and Saenko 2021).

## 3 Problem Formulation

**Notations.** Let  $\mathcal{X}^d$  and  $\mathcal{Y}^d$  denote the feature and label spaces of a domain  $d$  associated with a distribution  $P_d(X_d, Y_d) : \mathcal{X}^d \times \mathcal{Y}^d \rightarrow [0, 1]$  and labeling function  $h_d : \mathcal{X}^d \rightarrow \Delta(\mathcal{Y}^d)$  where  $X_d$  and  $Y_d$  are random variables that take values in  $\mathcal{X}^d$  and  $\mathcal{Y}^d$ , and  $\Delta(\mathcal{Y}^d)$  is a probability simplex over  $\mathcal{Y}^d$ . Consider a model  $h : \mathcal{X}^d \rightarrow \Delta(\mathcal{Y}^d)$ , then the *expected error* of  $h$  under domain  $d$  for some loss function  $L : \Delta(\mathcal{Y}^d) \times \mathcal{Y}^d \rightarrow \mathbb{R}_+$  (e.g., 0-1, cross-entropy loss) can be defined as  $E(P_d, h) = \mathbb{E}_{P_d} [L(h(X_d), Y_d)]$ .

**Open-Set Heterogeneous Domain Adaptation (OSHeDA) Setup.** In DA, we consider  $d \in \{s, t\}$  where  $s$  and  $t$  denote

the source and target domains, respectively. Different from conventional DA setup where feature and label spaces remain the same between source and target domains, in OSHeDA, we have  $\mathcal{X}^s \neq \mathcal{X}^t$  (*heterogeneous*) and  $\mathcal{Y}^s \subset \mathcal{Y}^t$  (*open-set*). Because  $\mathcal{Y}^s \subset \mathcal{Y}^t$ , we use  $Y$  to denote the random variable of label in both source and target domains, and we have  $P_s(Y \in \mathcal{Y}^t \setminus \mathcal{Y}^s) = 0$ . Moreover, classes in the sets  $\mathcal{Y}^t \setminus \mathcal{Y}^s$  are referred to as unknown in our setting. Given sets of samples  $D_s = \{x_i^s, y_i^s\}_{i=1}^{n_s} \stackrel{i.i.d.}{\sim} P_s(X_s, Y)$  (*source dataset*),  $D_{t_l} = \{x_i^t, y_i^t\}_{i=1}^{n_{t_l}} \stackrel{i.i.d.}{\sim} P_t(X_t, Y|Y \in \mathcal{Y}^s)$  (*labeled target dataset*), and  $D_{t_u} = \{x_i^t\}_{i=1}^{n_{t_u}} \stackrel{i.i.d.}{\sim} P_t(X_t, Y)$  (*unlabeled target dataset*), where  $n_s, n_{t_l}, n_{t_u}$  are size of datasets and  $n_{t_l} \ll n_s, n_{t_u}$ , the goal of OSHeDA is to learn a model  $h : \mathcal{X}^t \rightarrow \Delta(\mathcal{Y}^t)$  from  $D_s, D_{t_l}, D_{t_u}$  such that the expected error on the target domain  $E(P_t, h)$  is small.

**Representation learning.** Representation learning is a common approach for transferring knowledge from a source to a target domain in DA (Zhao et al. 2019; Ganin et al. 2016; Albuquerque et al. 2019; Pham, Zhang, and Zhang 2023), and we will leverage this method in OSHeDA. Specifically, it maps the input spaces  $\mathcal{X}^s$  and  $\mathcal{X}^t$  of the source and target domains to a shared representation space  $\mathcal{Z}$  using two representation mappings:  $f_s : \mathcal{X}^s \rightarrow \mathcal{Z}$  and  $f_t : \mathcal{X}^t \rightarrow \mathcal{Z}$ . A shared classifier  $h : \mathcal{Z} \rightarrow \Delta(\mathcal{Y}^t)$  can then be employed to make predictions from this representation space. Notably,  $h$  can be utilized for both domains because  $\mathcal{Y}^s \subset \mathcal{Y}^t$ .

## 4 Theoretical Analysis

In our analysis, we consider Jensen–Shannon (JS) divergence ( $\mathcal{D}_{JS}$ ) as the statistical distance between two domains. While different distances (Ben-David et al. 2010) were used in domain adaptation literature, we adopt JS divergence because it is aligned with the training objective of adversarial learning (Goodfellow et al. 2014), a technique used in many representation learning-based domain adaptation works (Zhao et al. 2019; Ganin et al. 2016; Pham, Zhang, and Zhang 2023). Next, we present the main theorems, with detailed proofs provided in Appendix A.

### 4.1 Learning bounds for OSHeDA (infinite case)

To simplify notations used in our following analysis, we denote  $P_{t,k}(\cdot) = P_t(\cdot|Y \in \mathcal{Y}^s)$  and  $P_{t,u}(\cdot) = P_t(\cdot|Y \notin \mathcal{Y}^s)$  as the distributions of target domain conditioned on known and unknown classes, respectively. We also introduce two distributions  $P_s^u$  and  $P_t^u$  induced from  $P_s$  and  $P_t$  by the two mappings  $f_s^u$  and  $f_t^u$  such that  $f_s^u(X^s, Y) = (X^s, unk)$  and  $f_t^u(X^t, Y) = (X^t, unk)$  where *unk* denotes unknown class. In addition, we adopt an assumption commonly used in DA literature (Nguyen et al. 2021; Mansour, Mohri, and Rostamizadeh 2009; Cortes and Mohri 2014) as follows.

**Assumption 1 (Bounded loss)** Assume loss function  $L$  defined on input space  $\mathcal{X}$  and output space  $\mathcal{Y}$  is upper bounded by a constant  $C$ , i.e.,  $\forall x \in \mathcal{X}, y \in \mathcal{Y}, h \in \mathcal{H}$ , we have  $L(h(x), y) \leq C$ .

We note that this assumption is indeed reasonable rather than stringent. For example, while Assumption 1 does not hold for the cross-entropy loss typically utilized in classification,

we can adjust this loss to ensure that it satisfies Assumption 1 (Pham, Zhang, and Zhang 2024). Based on this assumption, we then provide an upper bound for prediction error on the target domain in OSHeDA as follows.

**Theorem 1** Given a loss function  $L$  satisfying Assumption 1, then for any  $h \in \mathcal{H}, f_s \in \mathcal{F}_s, f_t \in \mathcal{F}_t$ , we have:

$$\begin{aligned} E(P_t, h \circ f_t) &\leq \underbrace{\lambda E(P_s, h \circ f_s)}_{\text{source error}} \\ &\quad + \underbrace{E(P_t^u, h \circ f_t) - \lambda E(P_s^u, h \circ f_s)}_{\text{open-set difference}} \\ &\quad + \sqrt{2}\lambda C \left( (\mathcal{D}_{JS}(P_s(Z) \parallel P_{t,k}(Z)))^{\frac{1}{2}} \right. \\ &\quad \left. + \underbrace{(\mathcal{D}_{JS}(P_s(Z, Y) \parallel P_{t,k}(Z, Y)))^{\frac{1}{2}}}_{\text{domain distance}} \right) \end{aligned}$$

where  $\lambda = P_t(Y \in \mathcal{Y}^s)$ ,  $\mathcal{H}, \mathcal{F}_s, \mathcal{F}_t$  are hypothesis classes for  $h, f_s, f_t$ , and  $P_s(Z)$  and  $P_{t,k}(Z)$  are the distributions induced from  $P_s(X_s)$  and  $P_{t,k}(X_t)$  by  $f_s$  and  $f_t$ , respectively.

**Remark 1** The upper bound in Theorem 1 shed a light on achieving good accuracy on target domain. Specifically, to minimize  $E(P_t, h \circ f_t)$ , the model need to optimize three terms: (i) the source error  $E(P_s, h \circ f_s)$ , (ii) the open-set difference  $E(P_t^u, h \circ f_t) - \lambda E(P_s^u, h \circ f_s)$ , and (iii) the distances of marginal and joint distributions between source domain and target domain conditioned on known labels  $\mathcal{D}_{JS}(P_s(Z) \parallel P_{t,k}(Z))$  and  $\mathcal{D}_{JS}(P_s(Z, Y) \parallel P_{t,k}(Z, Y))$ .

We want to emphasize that minimizing the distance of the joint distribution between the source and target domains,  $\mathcal{D}_{JS}(P_s(Z, Y) \parallel P_{t,k}(Z, Y))$ , requires knowledge of the label distribution in the target domain  $P_{t,k}(Y)$ . Therefore, access to labeled target data during training is essential to avoid negative transfer. Note that the concept of open-set difference is not exclusive to OSHeDA. This term also appears in existing works for OSDA (Fang et al. 2020) and positive-unlabeled learning (Kiryo et al. 2017) which are special cases of our setting. Thus, this demonstrates the consistency between our work and the existing literature. Next, we present a lower bound for OSHeDA.

**Proposition 1** Given a loss function  $L$  satisfying Assumption 1, then for any  $h \in \mathcal{H}, f_s \in \mathcal{F}_s, f_t \in \mathcal{F}_t$ , we have:

$$\begin{aligned} E(P_t, h \circ f_t) &\geq \lambda E(P_{t,k}, h \circ f_t) + (1 - \lambda) E(P_s^u, h \circ f_s) \\ &\quad - \sqrt{2}(1 - \lambda)C (\mathcal{D}_{JS}(P_s(Z) \parallel P_{t,u}(Z)))^{\frac{1}{2}} \end{aligned}$$

where  $P_{t,u}(Z)$  is distribution induced from  $P_{t,u}(X_t)$  by  $f_t$ .

**Remark 2** Theorem 1 shows the necessity of reducing  $E(P_s, h \circ f_s)$  to achieve high accuracy on target domain. However, it may unavoidably increase  $E(P_s^u, h \circ f_s)$ . This observation, combined with Proposition 1, suggests that to avoid the large lower bound for the target error  $E(P_t, h \circ f_t)$ , we should increase the distance of the marginal distribution between the source domain and the unknown data in target domain,  $\mathcal{D}_{JS}(P_s(Z) \parallel P_{t,u}(Z))$ . In other words, we should segregate the representations of known classes from those of unknown class.

## 4.2 Learning bound for OSHeDA (finite case)

The learning bounds discussed in Section 4.1 are only applicable for the setting when we have access to unlimited data from source and target domains. In such cases, minimizing JS divergence of data distribution between these domains is equivalent to achieving invariant representations through adversarial learning (Goodfellow et al. 2014). However, we only work with finite data in practice. Thus, we present the following result, which provides a guarantee for using adversarial learning to optimize JS divergence from finite data.

**Proposition 2 (Adapted from Biau et al. (2020))** *The error in minimizing JS divergence of data distributions between source and target domains in representation space, using finite data, is up to  $\mathcal{O}\left(\frac{1}{\sqrt{n_s}} + \frac{1}{\sqrt{n_t}}\right)$ .*

where  $n_s$  and  $n_t$  are the size of source and target datasets.

**Remark 3** *Proposition 2 states that the performance of minimizing JS divergence from finite data is proportional to the dataset size. Note that in OSHeDA, we only have access to limited label data from target domain which then results in significant error in estimating JS divergence only from labeled source and target data. In essence, this result underscores the need for the development of an effective approach to utilize unlabeled target data for estimating the JS divergence, which involves techniques like pseudo-labeling.*

Therefore, we apply pseudo-labeling on unlabeled data to enrich labeled target data. Let  $g$  be pseudo-label model and denote  $N(P_{t,k}, g) = \mathbb{E}[\mathcal{D}_{JS}(P_{t,k}(g(Z)) \| P_{t,k}(Y|Z))]$  as the noise of  $g$  with respect to the target domain conditioned on known labels. Then, the impact of pseudo-labeled data can be illustrated in a new bound for OSHeDA as follows.

**Theorem 2** *Given a loss function  $L$  satisfying Assumption 1, for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$ , the following holds for all  $h \in \mathcal{H}$ ,  $f_s \in \mathcal{F}_s$ ,  $f_t \in \mathcal{F}_t$ :*

$$\begin{aligned} \mathbb{E}(P_t, h \circ f_t) &\leq \lambda \widehat{\mathbb{E}}(P_s, h \circ f_s) + \widehat{\mathbb{E}}(P_t^u, h \circ f_t) \\ &- \lambda \widehat{\mathbb{E}}(P_s^u, h \circ f_s) + \sqrt{2}\lambda C \left( (\mathcal{D}_{JS}(P_s(Z) \| P_{t,k}(Z)))^{\frac{1}{2}} \right. \\ &+ \left. (\mathcal{D}_{JS}(P_s(Z, Y) \| P_{t,k}(Z, g(Z))))^{\frac{1}{2}} + (N(P_{t,k}, g))^{\frac{1}{2}} \right) \\ &+ \mathcal{O} \left( \lambda C \sqrt{\frac{d_s \log n_s + d_s \log |\mathcal{Y}^t| + \log \frac{1}{\delta}}{n_s}} \right) \\ &+ C \sqrt{\frac{d_t \log n_t + d_t \log |\mathcal{Y}^t| + \log \frac{1}{\delta}}{n_t}} \end{aligned}$$

where  $\widehat{\mathbb{E}}(P_s, h \circ f_s)$ ,  $\widehat{\mathbb{E}}(P_t^u, h \circ f_t)$ ,  $\widehat{\mathbb{E}}(P_s^u, h \circ f_s)$  are empirical errors calculated on samples from distributions  $P_s$ ,  $P_t^u$ ,  $P_s^u$ ,  $n_t = n_{t_l} + n_{t_u}$ , and  $d_s$ ,  $d_t$  are Natarajan dimension (Natarajan 1989) of hypothesis classes  $\mathcal{H} \circ \mathcal{F}_s$ ,  $\mathcal{H} \circ \mathcal{F}_t$ .

Theorem 2 shows that the error in the target domain depends on the quality of the pseudo-label model  $g$ , with higher-quality  $g$  being more effective at reducing noise. Additionally, the bound emphasizes the importance of aligning the joint distributions between the source and target domains

in OSHeDA. This makes OSHeDA more challenging compared to homogeneous DA (HoDA), where source and target data lie on the same space. In contrast, HoDA methods can attain good performance under certain conditions by solely aligning the marginal distributions of representations between source and target domains. We will illustrate this contrast through the bound for HoDA in Section 4.3.

## 4.3 Learning bound for HoDA

Before constructing the learning bound for HoDA, we introduce an assumption about the representation  $Z$  as follows.

**Assumption 2 (Sufficient representation)** *Let  $I_s(\cdot, \cdot)$  be the mutual information between two random variables in the source domain. We assume  $I_s(Z, Y) = I_s(X_s, Y)$ . In particular,  $I_s(Z, Y) = \mathcal{D}_{KL}(P_s(Z, Y) \| P_s(Z) \otimes P_s(Y))$  and  $I_s(X_s, Y) = \mathcal{D}_{KL}(P_s(X_s, Y) \| P_s(X_s) \otimes P_s(Y))$  where  $\mathcal{D}_{KL}$  is KL divergence between two distributions.*

Note that Assumption 2 is reasonable because we have access to labeled data of source domain and the dimension of  $\mathcal{Y}$  is often smaller than that of  $\mathcal{Z}$ . Based on this, we establish the learning bound in HoDA under the covariate shift below.

**Proposition 3** *Suppose Assumptions 1 and 2 hold and the distribution shift between source and target domains is covariate shift (i.e.,  $P_s(X) \neq P_t(X)$ ,  $P_s(Y|X) = P_t(Y|X)$ ), then for any  $h \in \mathcal{H}$  and  $f \in \mathcal{F}$ , we have:*

$$\mathbb{E}(P_t, h \circ f) \leq \mathbb{E}(P_s, h \circ f) + \sqrt{2}C (\mathcal{D}_{JS}(P_s(Z) \| P_t(Z)))^{\frac{1}{2}}$$

In HoDA, due to the homogeneity of the input space, we can utilize a single representation mapping  $f$  for both the source and target domains. Note that the bound in Proposition 3 depends solely on the distance of the marginal distributions between the source and target domains,  $\mathcal{D}_{JS}(P_s(Z) \| P_t(Z))$ , which can be effectively minimized even without access to labeled data in the target domain. Clearly, covariate shift assumption is only reasonable in HoDA, where the source and target data share the same feature and label spaces.

## 5 Methodology

Motivated by theoretical results presented in Section 4, we introduce RL-OSHeDA, a representation learning method specifically designed for OSHeDA. Our method aims to simultaneously optimize both the upper bound in Theorem 2 and the lower bound in Proposition 1. RL-OSHeDA features two distinct representation mappings,  $f_s$  and  $f_t$ , which map heterogeneous source and target feature spaces to a shared representation space, along with a classifier  $h$  that makes predictions based on these representations. Figure 2 presents the overall architecture of RL-OSHeDA, while pseudo code describing training process can be found in Appendix B.2.

### 5.1 Objective function

To improve predictive performance in OSHeDA, our method targets the following: (i) minimizing prediction errors on both source and labeled target data, (ii) minimizing the distances of marginal and label-conditioned representation distributions for known classes between source and target

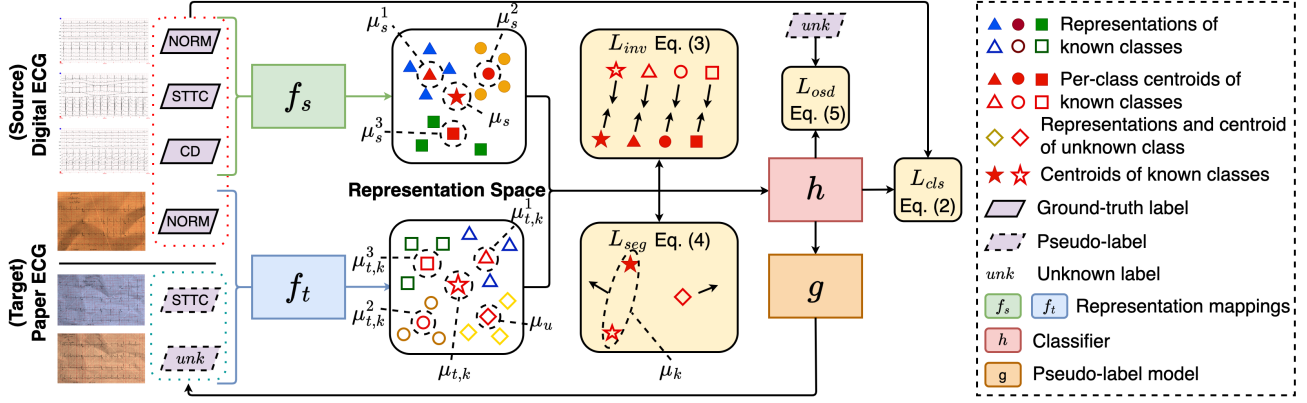


Figure 2: Overall architecture of RL-OSHeDA is illustrated with a motivating example from ECG-based diagnosis application. We leverage 2-stage learning process to update model parameters. In stage 1, model parameters are updated by optimizing  $L_{cls}$ . In stage 2, model parameters are updated by optimizing  $L_{cls}$ ,  $L_{inv}$ ,  $L_{seg}$ , and  $L_{osc}$  with the help from pseudo-label model  $g$ .

data, (iii) maximizing the distances of marginal representation distributions between known and unknown classes, and (iv) minimizing the open-set difference. Specifically, RL-OSHeDA optimizes the following objective function:

$$L = L_{cls} + L_{inv} - L_{seg} + L_{osc} \quad (1)$$

where  $L_{cls}$  is the classification error computed from source and labeled target datasets  $D_s$  and  $D_{t_l}$ , defined as follows:

$$L_{cls} = \frac{\lambda}{n_s} \sum_{i=1}^{n_s} \text{CE}(h(f_s(x_i^s)), y_i^s) + \frac{1}{n_{t_l}} \sum_{i=1}^{n_{t_l}} \text{CE}(h(f_t(x_i^t)), y_i^t) \quad (2)$$

Here CE is the cross-entropy loss.

$L_{inv}$  denotes the distances of marginal and label-conditioned representation distributions for known classes between source and target datasets. Note that, we minimize the distance of the label-conditioned representation distribution  $P(Z|Y)$ , rather than the joint distribution  $P(Z, Y)$ , as noted by Pham, Zhang, and Zhang (2023). As shown in Proposition 2,  $L_{inv}$  can be defined based on JS divergence and minimized through adversarial learning. However, the number of discriminators required for this approach scales linearly with the number of classes, leading to instability in training when the dataset has a large number of classes. To address this issue, we implement  $L_{inv}$  using maximum mean discrepancy (MMD) defined as follows:

$$L_{inv} = \|\mu_s - \mu_{t,k}\|_2^2 + \sum_{m=1}^{|\mathcal{Y}^s|} \|\mu_s^m - \mu_{t,k}^m\|_2^2 \quad (3)$$

where  $\mu_s$  (resp.  $\mu_{t,k}$ ) is centroid of representations from source data (resp. target data belonging to known classes), and  $\mu_s^m$  (resp.  $\mu_{t,k}^m$ ) is centroid of representations from source data (resp. target data) belonging to known class  $m$ . Note that  $\mu_{t,k}$  and  $\mu_{t,k}^m$  are computed using both instances with ground-truth labels from labeled target data and those with high-quality pseudo-labels (see Section 5.2) from unlabeled target data to provide a more accurate estimation.

$L_{seg}$  is the distances between marginal representation distributions of known and unknown classes. Similarly, we im-

plement  $L_{seg}$  with MMD as follows:

$$L_{seg} = \|\mu_k - \mu_u\|_2^2 \quad (4)$$

where  $\mu_k$  (resp.  $\mu_u$ ) are centroids of representations from both source and target datasets belonging to ground-truth and pseudo known (resp. unknown) classes.

$L_{osc}$  represents the open-set difference, as detailed in Theorem 2. The optimal value for the open-set difference is 0. However, due to the flexibility of deep neural networks, this term can become excessively negative during training and adversely affect model performance. To address this issue, we implement  $L_{osc}$  as a non-negative risk estimator:

$$L_{osc} = \max\left(0, \frac{1}{n_t} \sum_{i=1}^{n_t} \text{CE}(h(f_t(x_i^t)), \text{unk}) - \frac{\lambda}{n_s} \sum_{i=1}^{n_s} \text{CE}(h(f_s(x_i^s)), \text{unk})\right) \quad (5)$$

where  $n_t = n_{t_l} + n_{t_u}$  is the size of the target dataset.

## 5.2 Pseudo-labeling using 2-stage learning

The accuracy of  $L_{inv}$  and  $L_{seg}$  highly depends on the quality of pseudo-labels. Traditionally, the pseudo-label model  $g$  is derived by modifying the classifier  $h$  (e.g., using hard labels calculated from  $h$ 's outputs as pseudo-labels), which creates a coupling between  $g$  and  $h$ . Specifically,  $g$  is defined as  $a \circ h$ , where  $a$  is an operator applied to the output of  $h$  (e.g.,  $a := \arg \max$ ). When the distributions of the source and target domains are well-aligned, this coupling is harmless, as the optimal solution for  $g$  also aligns with that for  $h$ . However, at the beginning of the training process, when the distributions of the source and target domains are not aligned,  $g$  and  $h$  have completely different objective functions, resulting in a trade-off between them. To address this issue, we propose a 2-stage learning approach as follows:

- **Stage 1** ( $epoch < T$ ): Update  $f_s, f_t, h$  using  $L_{cls}$ .
- **Stage 2** ( $epoch \geq T$ ): Update  $f_s, f_t, h$  using  $L$ .

where  $T$  is a threshold indicating when to switch from stage 1 to stage 2. In stage 1, optimizing  $L_{cls}$  partially aligns the source and target domains, thereby reducing the trade-off between  $g$  and  $h$  during the optimization of  $L$  in stage 2. Additionally, rather than simply using the hard labels with

Table 1: Prediction performances ( $HOS$ ,  $OS^*$ ,  $UNK$ ) of RL-OSHeDA and baselines for OSHeDA scenario on 7 datasets. We report average results over 10 random seeds for each dataset.

	CIFAR10 & ILSVRC2012			ImageCLEF-DA			Multilingual Reuters Collection			NUSWIDE & ImageNet		
	$HOS$	$OS^*$	$UNK$	$HOS$	$OS^*$	$UNK$	$HOS$	$OS^*$	$UNK$	$HOS$	$OS^*$	$UNK$
DS3L	61.49±0.74	59.04±1.00	64.40±1.06	58.62±2.04	52.87±2.70	66.74±2.77	59.35±0.94	52.92±1.27	67.57±1.24	67.61±1.65	66.17±2.22	69.20±2.30
KPG	57.27±0.50	54.73±0.00	60.30±1.11	40.68±0.94	34.60±0.00	50.79±2.91	11.27±0.07	8.59±0.00	17.04±0.96	55.18±1.17	52.60±0.00	58.10±2.45
OPDA	53.30±0.77	48.22±1.00	60.26±1.11	53.17±1.83	45.28±2.13	65.76±2.76	55.85±0.96	48.47±1.23	65.94±1.23	71.06±1.44	66.60±1.98	76.38±2.09
PL	42.75±0.52	37.12±0.49	52.31±1.10	39.20±1.62	31.93±1.66	54.34±2.91	42.85±0.81	34.56±1.00	57.86±1.29	42.43±0.26	34.05±0.00	61.15±2.10
SCT	59.61±0.75	57.35±1.00	62.33±1.08	58.76±2.05	53.09±2.71	66.71±2.76	61.17±0.94	<b>54.96±1.30</b>	69.00±1.21	70.42±1.49	68.00±2.20	73.10±1.99
SSAN	60.38±0.73	59.01±1.00	62.01±1.08	58.61±2.05	53.18±2.74	66.14±2.74	58.25±0.93	51.99±1.25	66.26±1.24	67.98±1.49	66.25±2.04	69.85±2.21
STN	61.59±0.72	58.80±0.98	64.87±1.05	56.25±2.06	49.80±2.69	65.84±2.76	59.21±0.96	52.91±1.31	67.24±1.23	67.75±1.23	64.80±1.42	71.08±2.16
SL	60.74±0.74	58.29±1.00	63.67±1.08	58.59±2.05	52.84±2.70	66.71±2.76	58.53±0.96	52.14±1.32	66.74±1.21	69.41±1.64	66.63±2.26	72.57±2.23
RL-OSHeDA	<b>72.33±0.70</b>	<b>67.88±0.98</b>	<b>77.81±0.97</b>	<b>63.98±2.04</b>	<b>56.63±2.72</b>	<b>74.80±2.51</b>	<b>65.39±0.91</b>	54.47±1.21	<b>81.97±0.96</b>	<b>80.01±1.30</b>	<b>74.65±2.01</b>	<b>86.35±0.81</b>
	Office & Caltech256			Wikipedia			PTB-XL			Average over datasets		
	$HOS$	$OS^*$	$UNK$	$HOS$	$OS^*$	$UNK$	$HOS$	$OS^*$	$UNK$	$HOS$	$OS^*$	$UNK$
DS3L	72.06±2.48	67.41±3.68	78.15±2.89	56.00±2.01	50.72±2.36	66.24±2.80	30.30±1.19	34.95±0.58	26.74±1.82	57.92±1.58	54.87±1.97	62.72±2.12
KPG	34.46±0.99	29.18±0.00	45.34±3.58	24.82±0.42	16.40±0.00	52.52±3.18	N/A	N/A	N/A	37.28±0.68	32.68±0.00	47.35±2.36
OPDA	65.23±2.58	57.70±3.43	76.23±3.03	52.66±1.92	45.94±1.81	65.42±3.12	31.47±1.22	36.35±0.63	27.74±1.86	54.67±1.53	49.79±1.74	62.53±2.17
PL	48.92±1.36	40.38±1.13	68.41±3.40	41.87±1.65	35.14±1.48	58.40±3.10	26.18±1.37	36.43±0.55	20.43±1.66	40.60±1.08	35.66±0.90	53.27±2.22
SCT	75.72±2.14	71.05±3.20	81.79±2.54	58.41±2.01	52.86±2.39	68.42±2.74	26.23±1.65	<b>46.48±1.71</b>	18.27±1.60	59.89±1.58	57.10±1.93	64.49±1.99
SSAN	72.95±2.36	67.99±3.37	79.67±2.88	58.37±1.76	52.76±1.95	68.36±2.80	25.16±1.47	40.40±0.65	18.27±1.54	57.39±1.54	55.94±1.86	61.51±2.07
STN	72.26±2.28	66.46±3.30	79.84±2.75	57.75±1.91	51.40±2.13	69.00±2.97	27.08±0.96	22.63±0.26	33.72±1.65	57.41±1.45	52.40±1.73	64.51±2.08
SL	72.14±2.54	67.72±3.75	77.89±2.96	57.10±1.97	51.60±2.19	67.04±2.76	25.74±1.55	44.50±0.76	18.11±1.52	57.46±1.64	56.24±2.00	61.82±2.08
RL-OSHeDA	<b>78.18±2.05</b>	<b>73.04±2.91</b>	<b>85.25±2.50</b>	<b>63.10±1.89</b>	<b>57.26±2.45</b>	<b>73.04±2.37</b>	<b>47.48±1.25</b>	44.30±1.39	<b>51.16±1.86</b>	<b>67.21±1.45</b>	<b>61.18±1.95</b>	<b>75.77±1.71</b>

the largest logits from  $h$  as the output of  $g$ , we propose generating pseudo-labels as follows:

- First, select pseudo-labels as  $g(x^t) = a'(h(f_t(x^t)))$  where  $a'$  is arg max operator applied to the logits of the known classes only.
- Then, select  $1 - \lambda$  fraction of instances with the smallest maximum logits and assign pseudo-labels  $unk$  to them.

The motivation behind this design of the pseudo-label model  $g$  is that, at the beginning of stage 2, there is no supervision signal for training the parameters of  $h$  related to unknown class. Therefore, relying solely on logits to determine the unknown class is unreliable. Note that this strategy is only used to generate pseudo-labels during the training of stage 2. Once training is complete and  $h$ 's parameters for unknown class are well-trained by optimizing  $L_{osd}$ , we can simply use arg max across all classes to generate predictions.

## 6 Experiments

Next, we empirically evaluate the performance of our methods across clinical, computer vision, and natural language processing applications. We focus on the OSHeDA scenario, characterized by heterogeneity in the feature space between the source and target domains, with the label space of the target domain encompassing both known and unknown classes.

### 6.1 Experimental setup

**Datasets.** We conduct our experiments on 7 datasets including CIFAR10 (Krizhevsky 2009) & ILSVRC2012 (Russakovsky et al. 2015); Wikipedia (Rasiwasia et al. 2010); Multilingual Reuters Collection (Amini, Usunier, and Goutte 2009); NUSWIDE (Chua et al. 2009) & ImageNet (Deng et al. 2009); Office (Saenko et al. 2010) & Caltech256 (Griffin et al. 2007); ImageCLEF-DA (Griffin et al. 2007); PTB-XL (Wagner et al. 2020). These datasets results in 56 DA tasks. Detailed descriptions and statistics of these datasets are provided in Appendix C.1.

**Baselines.** We compare our method with several representative methods from **HeDA** (SSAN (Li et al. 2020), STN (Yao et al. 2019), SCT (Zhao et al. 2022), KPG (Gu et al.

2022)), **OSDA** (OPDA (Saito et al. 2018)), and **OS-SSL** (DS3L (Guo et al. 2020)) literature. For the HeDA methods, they are trained on both source and target data. In contrast, OSDA and OS-SSL methods are trained only on target data as they cannot handle heterogeneous feature spaces. During inference, HeDA and OS-SSL methods classify instances as  $unk$  using the same method as our pseudo-label model  $g$  (see Section 5.2). Additionally, we explore supervised learning (SL) and pseudo-labeling (PL) methods trained on target data. Among all baselines, only KPG is designed to handle OSHeDA by combining Gromov-Wasserstein distance and partial optimal transport (Xu et al. 2020). Since  $\lambda$  is a required input for most methods in our experiments, we utilize techniques from positive-unlabeled learning (Zeiberg, Jain, and Radivojac 2020) to estimate  $\lambda$ . Detailed architectures of our model and the baselines are in Appendix B.1.

**Evaluation method.** We utilize  $HOS$ , the harmonic mean of  $OS^*$  and  $UNK$  (Bucci, Loghmani, and Tommasi 2020).  $OS^*$  is the class-wise averaged accuracy of known classes, while  $UNK$  measures the accuracy for the unknown class.  $HOS$  is particularly suitable for OSHeDA because it emphasizes the ability to both correctly classify known classes and detect out-of-distribution instances simultaneously. In particular, this metric increases when the performance in both known and unknown classifications is high.

### 6.2 Experimental results

**OSHeDA benchmark.** The prediction performance ( $HOS$ ) of RL-OSHeDA and the baselines is summarized in Table 1. RL-OSHeDA consistently outperforms all baselines across all datasets, demonstrating its effectiveness in simultaneously addressing heterogeneity in the feature space and open-set in the label space during training. Among the baselines, KPG is specifically designed for OSHeDA by using optimal transport. Then, SVM trained on transported source and labeled target data is used to make prediction. However, this method underperforms in our evaluation due to its difficulty in correctly transporting from source to target data. Moreover, this method is not applicable for complex data structures, such as those found in PTB-XL dataset. Other



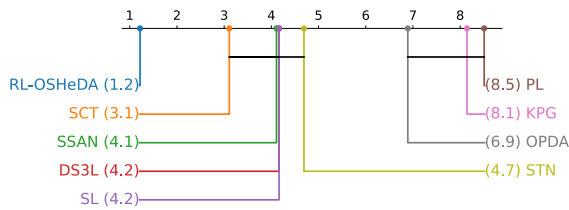


Figure 3: Critical Difference diagram for all methods calculated from 56 DA tasks. RL-OSHeDA is the highest ranked method on  $HOS$  metric, and its performance is significantly better than baselines (as indicated by the lack of connections between RL-OSHeDA and baselines in the diagram).

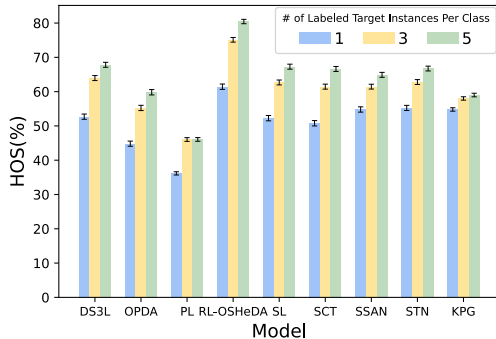


Figure 4: Performances w.r.t. different number of labeled target instances per class on CIFAR10 & ILSVRC2012 dataset.

baselines achieve better prediction performances, but their  $HOS$  remains suboptimal due to their inability to handle novel classes or heterogeneous source data during training.

To further validate the superiority of RL-OSHeDA across all DA tasks, we conduct significance testing, including the Friedman test followed by the Nemenyi test (Demšar 2006). The results (see Figure 3) show that our method significantly outperforms the baselines, with a P-value smaller than 0.05. Among all the baselines, SCT, SSAN, STN, SL, and DS3L exhibit better prediction performances than KPG, OPDA, and PL. Note that all methods, except OPDA and KPG, utilize our approach to detect the unknown class based on logits of known classes. This result suggests that while this approach can partially address the open-set issue, it cannot fully resolve it. For OPDA, although it is designed to handle open-set issue, its inability to leverage heterogeneous source data limits its performance to adapting with only a small labeled target dataset, resulting in suboptimal performance.

**Ablation study.** We conduct an ablation study to better understand the contribution of each component in the objective function of our method. As shown in Table 2, removing any component deteriorates model performance. This finding highlights the importance of achieving a good pseudo-label model using 2-stage learning approach as well as aligning the data distribution of known classes between source and target domains while simultaneously detecting and segregating unknown class from known ones for OSHeDA.

**Impact of labeled target data.** We vary the number of instances per class in the labeled target data to investigate their impact on the DA process. Specifically, we conduct experi-

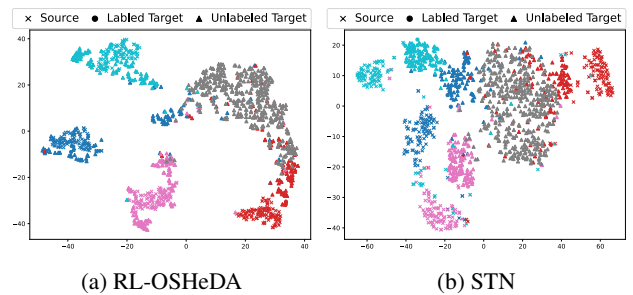


Figure 5: Visualization of representation spaces learned by RL-OSHeDA and STN for NUSWIDE & ImageNet dataset. Different colors represent different classes, with the unknown class denoted in grey.

Table 2: Ablation study for RL-OSHeDA on Multilingual Reuters Collection dataset. **Align** refers to using  $L_{inv}$ ; **Segregate** refers to using  $L_{seg}$ ; **OSD** refers to using  $L_{osd}$ ; **2-stage** refers to using 2-stage learning approach.

Align	Segregate	OSD	2-stage	$HOS$	$OS^*$	$UNK$
✓	✓	✓	✓	<b>65.39</b>	<b>54.47</b>	<b>81.97</b>
✓	✓	✓	✗	59.40	49.47	74.42
✓	✗	✓	✓	61.92	52.63	75.37
✗	✓	✓	✓	58.23	51.13	68.01
✓	✓	✗	✓	59.96	53.10	68.97
✗	✗	✗	✗	58.33	51.86	66.68

ments on CIFAR10 & ILSVRC2012 dataset with 1, 3, and 5 instances per class in the labeled target data and visualize the result in Figure 4. Generally, we observe that increasing the number of labeled target instances facilitates better alignment and enhances the performance of all methods. This result demonstrates the importance of labeled target data for DA methods in OSHeDA.

**Visualization of representation space.** We perform a qualitative analysis to examine the learned representations of RL-OSHeDA and STN for NUSWIDE & ImageNet dataset. Specifically, we use t-SNE (Van der Maaten and Hinton 2008) to project these representations into a 2-dimensional space. As shown in Figure 5, our method effectively aligns representations of the known classes between source and target domains while simultaneously segregating the representations of the unknown class (grey color). This results in improved  $HOS$  scores compared to STN.

## 7 Conclusion

This paper studied a novel domain adaptation scenario called open-set heterogeneous domain adaptation (OSHeDA). We first conducted a theoretical analysis to establish learning bounds in OSHeDA. Based on these theorems, we proposed a representation learning method that aligns the data distribution of known classes between source and target domains while simultaneously detecting and segregating unknown class from known ones. The resulting models trained with the proposed method generalize well to target domains. Experiments on real datasets across diverse domains, including healthcare, natural language processing, and computer vision, demonstrate the effectiveness of our proposed method.

## Acknowledgements

This work was funded in part by the National Science Foundation under award number IIS-2145625 and by the National Institutes of Health under award number R01AI188576. We would like to express our sincere gratitude to Dr. Wei-Lun Chao from The Ohio State University for his valuable suggestions, which greatly improved this paper.

## References

- Albuquerque, I.; Monteiro, J.; Darvishi, M.; Falk, T. H.; and Mitliagkas, I. 2019. Generalizing to unseen domains via distribution matching. *arXiv preprint arXiv:1911.00804*.
- Amini, M. R.; Usunier, N.; and Goutte, C. 2009. Learning from multiple partially observed views—an application to multilingual text categorization. *Advances in neural information processing systems*, 22.
- Ben-David, S.; Blitzer, J.; Crammer, K.; Kulesza, A.; Pereira, F.; and Vaughan, J. W. 2010. A theory of learning from different domains. *Machine learning*, 79: 151–175.
- Biau, G.; Cadre, B.; Sangnier, M.; and Tanielian, U. 2020. Some theoretical properties of GANS. *The Annals of Statistics*, 48(3): 1539 – 1566.
- Bucci, S.; Loghmani, M. R.; and Tommasi, T. 2020. On the effectiveness of image rotation for open set domain adaptation. In *European conference on computer vision*, 422–438. Springer.
- Chen, Y.; Zhu, X.; Li, W.; and Gong, S. 2020. Semi-supervised learning under class distribution mismatch. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 3569–3576.
- Chua, T.-S.; Tang, J.; Hong, R.; Li, H.; Luo, Z.; and Zheng, Y. 2009. Nus-wide: a real-world web image database from national university of singapore. In *Proceedings of the ACM international conference on image and video retrieval*, 1–9.
- Cortes, C.; and Mohri, M. 2014. Domain adaptation and sample bias correction theory and algorithm for regression. *Theoretical Computer Science*, 519: 103–126.
- Demšar, J. 2006. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine learning research*, 7: 1–30.
- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.
- Donahue, J.; Jia, Y.; Vinyals, O.; Hoffman, J.; Zhang, N.; Tzeng, E.; and Darrell, T. 2014. Decaf: A deep convolutional activation feature for generic visual recognition. In *International conference on machine learning*, 647–655. PMLR.
- Du, P.; Zhao, S.; Sheng, Z.; Li, C.; and Chen, H. 2023. Semi-Supervised Learning via Weight-aware Distillation under Class Distribution Mismatch. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 16410–16420.
- Fang, Z.; Lu, J.; Liu, F.; Xuan, J.; and Zhang, G. 2020. Open set domain adaptation: Theoretical bound and algorithm. *IEEE transactions on neural networks and learning systems*, 32(10): 4309–4322.
- Fang, Z.; Lu, J.; Liu, F.; and Zhang, G. 2022. Semi-supervised heterogeneous domain adaptation: Theory and algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1): 1087–1105.
- Ganin, Y.; Ustinova, E.; Ajakan, H.; Germain, P.; Larochelle, H.; Laviolette, F.; Marchand, M.; and Lempitsky, V. 2016. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1): 2096–2030.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2014. Generative adversarial nets. *Advances in neural information processing systems*, 27.
- Grief, S. N.; Patel, J.; Kochendorfer, K. M.; Green, L. A.; Lussier, Y. A.; Li, J.; Burton, M.; and Boyd, A. D. 2016. Simulation of ICD-9 to ICD-10-CM transition for family medicine: simple or convoluted? *The Journal of the American Board of Family Medicine*, 29(1): 29–36.
- Griffin, G.; Holub, A.; Perona, P.; et al. 2007. Caltech-256 object category dataset. Technical report, Technical Report 7694, California Institute of Technology Pasadena.
- Gu, X.; Yang, Y.; Zeng, W.; Sun, J.; and Xu, Z. 2022. Keypoint-guided optimal transport with applications in heterogeneous domain adaptation. *Advances in Neural Information Processing Systems*, 35: 14972–14985.
- Guo, L.-Z.; Zhang, Z.-Y.; Jiang, Y.; Li, Y.-F.; and Zhou, Z.-H. 2020. Safe deep semi-supervised learning for unseen-class unlabeled data. In *International conference on machine learning*, 3897–3906. PMLR.
- He, R.; Han, Z.; Lu, X.; and Yin, Y. 2022. Safe-student for safe deep semi-supervised learning with unseen-class unlabeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14585–14594.
- Hoffman, J.; Rodner, E.; Donahue, J.; Darrell, T.; and Saenko, K. 2013. Efficient learning of domain-invariant image representations. *arXiv preprint arXiv:1301.3224*.
- Hoffman, J.; Rodner, E.; Donahue, J.; Kulis, B.; and Saenko, K. 2014. Asymmetric and category invariant feature transformations for domain adaptation. *International journal of computer vision*, 109: 28–41.
- Huang, J.; Fang, C.; Chen, W.; Chai, Z.; Wei, X.; Wei, P.; Lin, L.; and Li, G. 2021. Trash to treasure: Harvesting odd data with cross-modal matching for open-set semi-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 8310–8319.
- Huang, Z.; Yang, J.; and Gong, C. 2022. They are not completely useless: Towards recycling transferable unlabeled data for class-mismatched semi-supervised learning. *IEEE Transactions on Multimedia*, 25: 1844–1857.
- Kiryo, R.; Niu, G.; Du Plessis, M. C.; and Sugiyama, M. 2017. Positive-unlabeled learning with non-negative risk estimator. *Advances in neural information processing systems*, 30.
- Kolesnikov, A.; Beyer, L.; Zhai, X.; Puigcerver, J.; Yung, J.; Gelly, S.; and Houlsby, N. 2019. Big transfer (BiT): General visual representation learning. *arXiv preprint arXiv:1912.11370*.



- Koltchinskii, V.; and Panchenko, D. 2000. Rademacher processes and bounding the risk of function learning. In *High dimensional probability II*, 443–457. Springer.
- Krizhevsky, A. 2009. Learning Multiple Layers of Features from Tiny Images. *Master's thesis, University of Tront*.
- LeCun, Y.; Bengio, Y.; and Hinton, G. 2015. Deep learning. *nature*, 521(7553): 436–444.
- Li, G.; Kang, G.; Zhu, Y.; Wei, Y.; and Yang, Y. 2021. Domain consensus clustering for universal domain adaptation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9757–9766.
- Li, J.; Lu, K.; Huang, Z.; Zhu, L.; and Shen, H. T. 2018. Heterogeneous domain adaptation through progressive alignment. *IEEE transactions on neural networks and learning systems*, 30(5): 1381–1391.
- Li, S.; Xie, B.; Wu, J.; Zhao, Y.; Liu, C. H.; and Ding, Z. 2020. Simultaneous semantic alignment network for heterogeneous domain adaptation. In *Proceedings of the 28th ACM international conference on multimedia*, 3866–3874.
- Li, W.; Duan, L.; Xu, D.; and Tsang, I. W. 2013. Learning with augmented features for supervised and semi-supervised heterogeneous domain adaptation. *IEEE Transactions on Pattern analysis and machine intelligence*, 36(6): 1134–1148.
- Liang, P. 2016. CS229T/STAT231: Statistical learning theory (Winter 2016).
- Liu, F.; Zhang, G.; and Lu, J. 2020. Heterogeneous domain adaptation: An unsupervised approach. *IEEE transactions on neural networks and learning systems*, 31(12): 5588–5602.
- Liu, H.; Cao, Z.; Long, M.; Wang, J.; and Yang, Q. 2019. Separate to adapt: Open set domain adaptation via progressive separation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2927–2936.
- Luo, Y.; Wang, Z.; Huang, Z.; and Baktashmotlagh, M. 2020. Progressive graph learning for open-set domain adaptation. In *International Conference on Machine Learning*, 6468–6478. PMLR.
- Mansour, Y.; Mohri, M.; and Rostamizadeh, A. 2009. Domain adaptation: Learning bounds and algorithms. *arXiv preprint arXiv:0902.3430*.
- Natarajan, B. K. 1989. On learning sets and functions. *Machine Learning*, 4: 67–97.
- Nguyen, A. T.; Tran, T.; Gal, Y.; Torr, P.; and Baydin, A. G. 2021. KL Guided Domain Adaptation. In *International Conference on Learning Representations*.
- Pham, T.-H.; Zhang, X.; and Zhang, P. 2023. Fairness and Accuracy under Domain Generalization. In *The Eleventh International Conference on Learning Representations*.
- Pham, T.-H.; Zhang, X.; and Zhang, P. 2024. Non-stationary Domain Generalization: Theory and Algorithm. In *The 40th Conference on Uncertainty in Artificial Intelligence*.
- Quiñonero-Candela, J.; Sugiyama, M.; Schwaighofer, A.; and Lawrence, N. D. 2022. *Dataset shift in machine learning*. Mit Press.
- Rasiwasia, N.; Costa Pereira, J.; Coviello, E.; Doyle, G.; Lanckriet, G. R.; Levy, R.; and Vasconcelos, N. 2010. A new approach to cross-modal multimedia retrieval. In *Proceedings of the 18th ACM international conference on Multimedia*, 251–260.
- Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115: 211–252.
- Saenko, K.; Kulis, B.; Fritz, M.; and Darrell, T. 2010. Adapting visual category models to new domains. In *Computer Vision—ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5–11, 2010, Proceedings, Part IV 11*, 213–226. Springer.
- Saito, K.; Kim, D.; and Saenko, K. 2021. Openmatch: Open-set semi-supervised learning with open-set consistency regularization. *Advances in Neural Information Processing Systems*, 34: 25956–25967.
- Saito, K.; Yamamoto, S.; Ushiku, Y.; and Harada, T. 2018. Open set domain adaptation by backpropagation. In *Proceedings of the European conference on computer vision (ECCV)*, 153–168.
- Shalev-Shwartz, S.; and Ben-David, S. 2014. *Understanding machine learning: From theory to algorithms*. Cambridge university press.
- Shen, C.; and Guo, Y. 2018. Unsupervised heterogeneous domain adaptation with sparse feature transformation. In *Asian conference on machine learning*, 375–390. PMLR.
- Simonyan, K.; and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Van der Maaten, L.; and Hinton, G. 2008. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11).
- Wagner, P.; Strodthoff, N.; Bousseljot, R.-D.; Kreiseler, D.; Lunze, F. I.; Samek, W.; and Schaeffter, T. 2020. PTB-XL, a large publicly available electrocardiography dataset. *Scientific data*, 7(1): 1–15.
- Wang, Y.; Qiao, P.; Liu, C.; Song, G.; Zheng, X.; and Chen, J. 2023. Out-of-distributed semantic pruning for robust semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 23849–23858.
- Xu, R.; Liu, P.; Zhang, Y.; Cai, F.; Wang, J.; Liang, S.; Ying, H.; and Yin, J. 2020. Joint Partial Optimal Transport for Open Set Domain Adaptation. In *IJCAI*, 2540–2546.
- Yao, Y.; Zhang, Y.; Li, X.; and Ye, Y. 2019. Heterogeneous domain adaptation via soft transfer network. In *Proceedings of the 27th ACM MM*, 1578–1586.
- Yao, Y.; Zhang, Y.; Li, X.; and Ye, Y. 2020. Discriminative distribution alignment: A unified framework for heterogeneous domain adaptation. *Pattern Recognition*, 101: 107165.
- Yu, Q.; Ikami, D.; Irie, G.; and Aizawa, K. 2020. Multi-task curriculum framework for open-set semi-supervised learning. In *Computer Vision—ECCV 2020: 16th European Con-*

ference, Glasgow, UK, August 23–28, 2020, *Proceedings, Part XII 16*, 438–454. Springer.

Zaheer, M.; Guruganesh, G.; Dubey, K. A.; Ainslie, J.; Alberti, C.; Ontanon, S.; Pham, P.; Ravula, A.; Wang, Q.; Yang, L.; et al. 2020. Big bird: Transformers for longer sequences. *Advances in neural information processing systems*, 33: 17283–17297.

Zeiberg, D.; Jain, S.; and Radivojac, P. 2020. Fast non-parametric estimation of class proportions in the positive-unlabeled classification setting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 6729–6736.

Zhao, H.; Des Combes, R. T.; Zhang, K.; and Gordon, G. 2019. On learning invariant representations for domain adaptation. In *International conference on machine learning*, 7523–7532. PMLR.

Zhao, Y.; Li, S.; Zhang, R.; Liu, C. H.; Cao, W.; Wang, X.; and Tian, S. 2022. Semantic correlation transfer for heterogeneous domain adaptation. *IEEE Transactions on Neural Networks and Learning Systems*.

Zou, Y.; Yu, Z.; Kumar, B.; and Wang, J. 2018. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *Proceedings of the European conference on computer vision (ECCV)*, 289–305.

## A Proofs

### A.1 Additional Lemmas

**Lemma 1** *Given two domains  $s$  and  $t$  associated with two distributions  $P_s(X, Y)$  and  $P_t(X, Y)$ , respectively, then for any classifier  $h : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ , the expected error of  $h$  in domain  $t$  can be upper bounded:*

$$\mathbb{E}(P_t, h) \leq \mathbb{E}(P_s, h) + \sqrt{2}C \times \mathcal{D}_{JS}(P_t(X, Y) \parallel P_s(X, Y))^{1/2}$$

where  $\mathcal{D}_{JS}(\cdot \parallel \cdot)$  is JS-divergence between two distributions.

**Proof of Lemma 1** Let  $\mathcal{D}_{KL}(\cdot \parallel \cdot)$  be KL-divergence,  $p_s, p_t$  are probability density functions associated with  $P_s, P_t$ , and  $U = (X, Y)$  and  $L(U) = L(h(X), Y)$ . We first prove  $\int_{\mathcal{E}} |p_t(u) - p_s(u)| du = \frac{1}{2} \int |p_t(u) - p_s(u)| du$  where  $\mathcal{E}$  is the event that  $p_t(u) \geq p_s(u)$  (\*) as follows:

$$\begin{aligned} \int_{\mathcal{E}} |p_t(u) - p_s(u)| du &= \int_{\mathcal{E}} (p_t(u) - p_s(u)) du \\ &= \int_{\mathcal{E} \cup \bar{\mathcal{E}}} (p_t(u) - p_s(u)) du - \int_{\bar{\mathcal{E}}} (p_t(u) - p_s(u)) du \\ &\stackrel{(1)}{=} \int_{\bar{\mathcal{E}}} (p_s(u) - p_t(u)) du \\ &= \int_{\bar{\mathcal{E}}} |p_t(u) - p_s(u)| du \\ &= \frac{1}{2} \int |p_t(u) - p_s(u)| du \end{aligned}$$

where  $\bar{\mathcal{E}}$  is the complement of  $\mathcal{E}$ . We have  $\stackrel{(1)}{=}$  because  $\int_{\mathcal{E} \cup \bar{\mathcal{E}}} (p_t(u) - p_s(u)) du = \int_{\mathcal{U}} (p_t(u) - p_s(u)) du = 0$ . Then, we have:

$$\begin{aligned} \mathbb{E}(P_t, h) &= \mathbb{E}_{P_t}[L(U)] \\ &= \int_{\mathcal{U}} L(u)p_t(u)du \\ &= \int_{\mathcal{U}} L(u)p_s(u)du + \int_{\mathcal{U}} L(u)(p_t(u) - p_s(u)) du \\ &= \mathbb{E}_{P_s}[L(U)] + \int_{\mathcal{U}} L(u)(p_t(u) - p_s(u)) du \\ &= \mathbb{E}(P_s, h) + \int_{\mathcal{E}} L(u)(p_t(u) - p_s(u)) du + \int_{\bar{\mathcal{E}}} L(u)(p_t(u) - p_s(u)) du \\ &\stackrel{(2)}{\leq} \mathbb{E}(P_s, h) + \int_{\mathcal{E}} L(u)(p_t(u) - p_s(u)) du \\ &\stackrel{(3)}{\leq} \mathbb{E}(P_s, h) + C \int_{\mathcal{E}} (p_t(u) - p_s(u)) du \\ &= \mathbb{E}(P_s, h) + C \int_{\mathcal{E}} |p_t(u) - p_s(u)| du \\ &\stackrel{(4)}{=} \mathbb{E}(P_s, h) + \frac{C}{2} \int |p_t(u) - p_s(u)| du \\ &\stackrel{(5)}{\leq} \mathbb{E}(P_s, h) + \frac{C}{2} \sqrt{2 \min(\mathcal{D}_{KL}(P_s(U) \parallel P_t(U)), \mathcal{D}_{KL}(P_t(U) \parallel P_s(U)))} \\ &\leq \mathbb{E}(P_s, h) + \frac{C}{\sqrt{2}} \sqrt{\mathcal{D}_{KL}(P_t(U) \parallel P_s(U))} \tag{6} \end{aligned}$$

We have  $\stackrel{(2)}{\leq}$  because  $\int_{\bar{\mathcal{E}}} L(u)(p_t(u) - p_s(u)) du \leq 0$ ;  $\stackrel{(3)}{\leq}$  because  $L(u)$  is non-negative function and is bounded by  $C$ ;  $\stackrel{(4)}{=}$  by using (\*);  $\stackrel{(5)}{\leq}$  by using Pinsker's inequality between total variation norm and KL-divergence.

Let  $P_{s,t}(U) = \frac{1}{2}(P_t(U) + P_s(U))$ . Apply Eq.(6) for two distributions  $P_t$  and  $P_{s,t}$ , we have:

$$\mathbb{E}(P_t, h) \leq \mathbb{E}(P_{s,t}, h) + \frac{C}{\sqrt{2}} \sqrt{\mathcal{D}_{KL}(P_t(U) \parallel P_{s,t}(U))} \quad (7)$$

Apply Eq.(6) again for two distributions  $P_{s,t}$  and  $P_s$ , we have:

$$\mathbb{E}(P_{s,t}, h) \leq \mathbb{E}(P_s, h) + \frac{C}{\sqrt{2}} \sqrt{\mathcal{D}_{KL}(P_s(U) \parallel P_{s,t}(U))} \quad (8)$$

Adding Eq. (7) to Eq. (8) and subtracting  $\mathbb{E}(P_{s,t}, h)$ , we have:

$$\begin{aligned} \mathbb{E}(P_t, h) &\leq \mathbb{E}(P_s, h) + \frac{C}{\sqrt{2}} \left( \sqrt{\mathcal{D}_{KL}(P_t(U) \parallel P_{s,t}(U))} + \sqrt{\mathcal{D}_{KL}(P_s(U) \parallel P_{s,t}(U))} \right) \\ &\stackrel{(6)}{\leq} \mathbb{E}(P_s, h) + \frac{C}{\sqrt{2}} \sqrt{2(\mathcal{D}_{KL}(P_t(U) \parallel P_{s,t}(U)) + \mathcal{D}_{KL}(P_s(U) \parallel P_{s,t}(U)))} \\ &= \mathbb{E}(P_s, h) + \frac{C}{\sqrt{2}} \sqrt{4\mathcal{D}_{JS}(P_s(U) \parallel P_t(U))} \\ &= \mathbb{E}(P_s, h) + \sqrt{2}C \sqrt{\mathcal{D}_{JS}(P_s(U) \parallel P_t(U))} \end{aligned}$$

We have  $\stackrel{(6)}{\leq}$  by using Cauchy–Schwarz inequality.

**Lemma 2** Given two domains  $s$  and  $t$  associated with two distributions  $P_s(X, Y)$  and  $P_t(X, Y)$ , respectively, then JS-divergence  $\mathcal{D}_{JS}(P_s(X, Y) \parallel P_t(X, Y))$  can be decomposed as follows:

$$\begin{aligned} \mathcal{D}_{JS}(P_s(X, Y) \parallel P_t(X, Y)) &\leq \mathcal{D}_{JS}(P_s(Y) \parallel P_t(Y)) + \mathbb{E}_{P_s}[\mathcal{D}_{JS}(P_s(X|Y) \parallel P_t(X|Y))] \\ &\quad + \mathbb{E}_{P_t}[\mathcal{D}_{JS}(P_s(X|Y) \parallel P_t(X|Y))] \end{aligned}$$

**Proof of Lemma 2** First, we show the decomposition formulation for KL-divergence as follow.

$$\begin{aligned} &\mathcal{D}_{KL}(P_s(X, Y) \parallel P_t(X, Y)) \\ &= \mathbb{E}_{P_s}[\log p_s(X, Y) - \log p_t(X, Y)] \\ &= \mathbb{E}_{P_s}[\log p_s(Y) + \log p_s(X|Y)] - \mathbb{E}_{P_s}[\log p_t(Y) + \log p_t(X|Y)] \\ &= \mathbb{E}_{P_s}[\log p_s(Y) - \log p_t(Y)] + \mathbb{E}_{P_s}[\log p_s(X|Y) - \log p_t(X|Y)] \\ &= \mathbb{E}_{P_s}[\log p_s(Y) - \log p_t(Y)] + \mathbb{E}_{y \sim P_s(Y)}[\mathbb{E}_{x \sim P_s(X|Y)}[\log p_s(X|Y) - \log p_t(X|Y)]] \\ &= \mathcal{D}_{KL}(P_s(Y) \parallel P_t(Y)) + \mathbb{E}_{P_s}[\mathcal{D}_{KL}(P_s(X|Y) \parallel P_t(X|Y))] \end{aligned} \quad (9)$$

For JS-divergence, we have:

$$\begin{aligned} &\mathcal{D}_{JS}(P_s(X, Y) \parallel P_t(X, Y)) \\ &= \frac{1}{2}(\mathcal{D}_{KL}(P_s(X, Y) \parallel P_{s,t}(X, Y))) + \frac{1}{2}(\mathcal{D}_{KL}(P_t(X, Y) \parallel P_{s,t}(X, Y))) \\ &\stackrel{(1)}{=} \frac{1}{2}(\mathcal{D}_{KL}(P_s(Y) \parallel P_{s,t}(Y))) + \frac{1}{2}(\mathbb{E}_{P_s}[\mathcal{D}_{KL}(P_s(X|Y) \parallel P_{s,t}(X|Y))]) \\ &\quad + \frac{1}{2}(\mathcal{D}_{KL}(P_t(Y) \parallel P_{s,t}(Y))) + \frac{1}{2}(\mathbb{E}_{P_t}[\mathcal{D}_{KL}(P_t(X|Y) \parallel P_{s,t}(X|Y))]) \\ &= \mathcal{D}_{JS}(P_s(Y) \parallel P_t(Y)) + \frac{1}{2}(\mathbb{E}_{P_s}[\mathcal{D}_{KL}(P_s(X|Y) \parallel P_{s,t}(X|Y))]) + \frac{1}{2}(\mathbb{E}_{P_t}[\mathcal{D}_{KL}(P_t(X|Y) \parallel P_{s,t}(X|Y))]) \\ &\leq \mathcal{D}_{JS}(P_s(Y) \parallel P_t(Y)) + \frac{1}{2}(\mathbb{E}_{P_s}[\mathcal{D}_{KL}(P_s(X|Y) \parallel P_{s,t}(X|Y))]) + \frac{1}{2}(\mathbb{E}_{P_s}[\mathcal{D}_{KL}(P_t(X|Y) \parallel P_{s,t}(X|Y))]) \\ &\quad + \frac{1}{2}(\mathbb{E}_{P_t}[\mathcal{D}_{KL}(P_t(X|Y) \parallel P_{s,t}(X|Y))]) + \frac{1}{2}(\mathbb{E}_{P_t}[\mathcal{D}_{KL}(P_s(X|Y) \parallel P_{s,t}(X|Y))]) \\ &= \mathcal{D}_{JS}(P_s(Y) \parallel P_t(Y)) + \mathbb{E}_{P_s}[\mathcal{D}_{JS}(P_s(X|Y) \parallel P_t(X|Y))] + \mathbb{E}_{P_t}[\mathcal{D}_{JS}(P_s(X|Y) \parallel P_t(X|Y))] \end{aligned}$$

We have  $\stackrel{(1)}{=}$  by applying Eq. (9) for  $\mathcal{D}_{KL}(P_s(X, Y) \parallel P_{s,t}(X, Y))$  and  $\mathcal{D}_{KL}(P_t(X, Y) \parallel P_{s,t}(X, Y))$ .

**Lemma 3** Given two domains  $s$  and  $t$  associated with two distributions  $P_s(X, Y)$  and  $P_t(X, Y)$ , respectively, let  $f : \mathcal{X} \rightarrow \mathcal{Z}$  be the representation mapping from input space  $\mathcal{X}$  to representation space  $\mathcal{Z}$ . If the shift between domains  $s$  and  $t$  is covariate shift (i.e.,  $P_s(Y|X) = P_t(Y|X)$ ), and Assumption 2 holds for the representation  $Z$ , then the shift between these two domains in representation space is also covariate shift (i.e.,  $P_s(Y|Z) = P_t(Y|Z)$ ).

**Proof of Lemma 3** We have:

$$\begin{aligned}
\log p_s(y|x) &= \log \left( \int p_s(y, z|x) dz \right) \\
&= \log \left( \int p_s(y|z) p(z|x) dz \right) \\
&= \log \left( \mathbb{E}_{P(Z|x)} [p_s(y|Z)] \right) \\
&\stackrel{(1)}{\geq} \mathbb{E}_{P(Z|x)} [\log p_s(y|Z)]
\end{aligned} \tag{10}$$

We have  $\stackrel{(1)}{\geq}$  by using Jensen's inequality. Taking expectation w.r.t.  $P_t(X, Y)$  over both sides, we have:

$$\begin{aligned}
&\mathbb{E}_{P_t(X, Y)} [\log p_s(Y|X) - \mathbb{E}_{P(Z|X)} [\log p_s(Y|Z)]] \\
&= \int \int (\log p_s(y|x) - \mathbb{E}_{P(Z|x)} [\log p_s(Y|Z)]) p_t(x, y) dx dy \\
&= \int \int (\log p_s(y|x) - \mathbb{E}_{P(Z|x)} [\log p_s(Y|Z)]) p_s(x, y) \frac{p_t(x, y)}{p_s(x, y)} dx dy \\
&= \mathbb{E}_{P_s(X, Y)} \left[ (\log p_s(Y|X) - \mathbb{E}_{P(Z|X)} [\log p_s(Y|Z)]) \frac{p_t(X, Y)}{p_s(X, Y)} \right] \\
&\stackrel{(1)}{\leq} \left( \max_{x, y} \frac{p_t(x, y)}{p_s(x, y)} \right) \mathbb{E}_{P_s(X, Y)} [\log p_s(Y|X) - \mathbb{E}_{P(Z|X)} [\log p_s(Y|Z)]] \\
&= \left( \max_{x, y} \frac{p_t(x, y)}{p_s(x, y)} \right) (\mathbb{E}_{P_s(X, Y)} [\log p_s(Y|X)] - \mathbb{E}_{P_s(Z, Y)} [\log p_s(Y|Z)]) \\
&= \left( \max_{x, y} \frac{p_t(x, y)}{p_s(x, y)} \right) (H_s(Y, X) - H_s(Y, Z)) \\
&= \left( \max_{x, y} \frac{p_t(x, y)}{p_s(x, y)} \right) ((H_s(Y) - H_s(Y, Z)) - (H_s(Y) - H_s(Y, X))) \\
&= \left( \max_{x, y} \frac{p_t(x, y)}{p_s(x, y)} \right) (I_s(Y, Z) - I_s(Y, X)) \\
&\stackrel{(2)}{=} 0
\end{aligned} \tag{11}$$

We have  $\stackrel{(1)}{\leq}$  because  $\log p_s(y|x) - \mathbb{E}_{P(Z|x)} [\log p_s(y|z)] \geq 0$  according to Eq. (10);  $\stackrel{(2)}{=}$  because  $I_s(Y, Z) = I_s(Y, X)$  according to Assumption 2. Based on Eq. (11), we have:

$$\begin{aligned}
\mathbb{E}_{P_t(X, Y)} [\log p_s(Y|X)] &= \mathbb{E}_{P_t(X, Y)} [\mathbb{E}_{P(Z|X)} [\log p_s(Y|Z)]] \\
&= \mathbb{E}_{P_t(Y, Z)} [\log p_s(Y|Z)]
\end{aligned} \tag{12}$$

We also have:

$$\begin{aligned}
\mathbb{E}_{P_t(Y, Z)} [\log P_t(Y|Z)] &= -H_t(Y|Z) \\
&= I_t(Y, Z) - H_t(Y) \\
&\stackrel{(1)}{\leq} I_t(Y, X) - H_t(Y) \\
&= -H_t(Y|X) \\
&= \mathbb{E}_{P_t(X, Y)} [\log P_t(Y|X)]
\end{aligned} \tag{13}$$

We have  $\stackrel{(1)}{\leq}$  by using data processing inequality. Finally, we have:

$$\begin{aligned}
& \mathbb{E}_{P_t(Z)} [\mathcal{D}_{KL} (P_t(Y|Z) \parallel P_s(Y|Z))] \\
& \stackrel{(1)}{=} \mathbb{E}_{P_t(Z)} [\mathcal{D}_{KL} (P_t(Y|Z) \parallel P_s(Y|Z))] - \mathbb{E}_{P_t(X)} [\mathcal{D}_{KL} (P_t(Y|X) \parallel P_s(Y|X))] \\
& = \mathbb{E}_{P_t(Y,Z)} [\log p_t(Y|Z) - \log p_s(Y|Z)] - \mathbb{E}_{P_t(X,Y)} [\log p_t(Y|X) - \log p_s(Y|X)] \\
& = (\mathbb{E}_{P_t(Y,Z)} [\log p_t(Y|Z)] - \mathbb{E}_{P_t(X,Y)} [\log p_t(Y|X)]) + (\mathbb{E}_{P_t(X,Y)} [\log p_s(Y|X)] - \mathbb{E}_{P_t(Y,Z)} [\log p_s(Y|Z)]) \\
& \stackrel{(2)}{=} 0
\end{aligned} \tag{14}$$

We have  $\stackrel{(1)}{=}$  because the shift between two domains w.r.t. input space  $\mathcal{X}$  is covariate shift;  $\stackrel{(2)}{=}$  by using Eq. (12) and Eq. (13) and the fact that KL-divergence is non-negative. Note that Eq. (14) implies that the shift between these two domains w.r.t. representation space  $\mathcal{Z}$  is also covariate shift (i.e.,  $P_s(Y|Z) = P_t(Y|Z)$ ).

**Lemma 4** *Given domain  $d$  associated with a distribution  $P_d(X, Y)$ , then for any  $\delta > 0$ , with probability at least  $1 - \delta$  over sample  $S$  of size  $n$  drawn i.i.d from domain  $d$ , for all  $h \in \mathcal{H} : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ , the expected error of  $h$  in domain  $d$  can be upper bounded:*

$$\mathbb{E} (P_d, h) \leq \widehat{\mathbb{E}} (P_d, h) + 2\mathcal{R}_S(\mathcal{L} \circ \mathcal{H}) + 3C\sqrt{\frac{\log(2/\delta)}{2n}}$$

where  $\mathcal{L} \circ \mathcal{H} = \{(x, y) \rightarrow L(h(x), y) : h \in \mathcal{H}\}$  and  $\mathcal{R}_S(\mathcal{L} \circ \mathcal{H})$  is an empirical Rademacher complexity of the function class  $\mathcal{L} \circ \mathcal{H}$  computed from the sample  $S$ .

**Proof of Lemma 4** We start from the Rademacher bound (Koltchinskii and Panchenko 2000) which is stated as follows.

**Rademacher Bounds.** Let  $\mathcal{F}$  be a family of functions mapping from  $Z$  to  $[0, 1]$ . Then, for any  $0 < \delta < 1$ , with probability at least  $1 - \delta$  over sample  $S = \{z_1, \dots, z_n\}$ , the following holds for all  $f \in \mathcal{F}$ :

$$\mathbb{E} [f^Z] \leq \frac{1}{n} \sum_{i=1}^n f(z_i) + 2\mathcal{R}_S(\mathcal{F}) + 3\sqrt{\frac{\log(2/\delta)}{2n}}$$

where  $\mathcal{R}_S(\mathcal{F})$  is an empirical Rademacher complexity of function class  $\mathcal{F}$  computed from the sample  $S$ .

We then apply this result to our setting with  $Z = (X, Y)$ , the loss function  $L$  bounded by  $C$ , and the function class  $\mathcal{L} \circ \mathcal{H} = \{(x, y) \rightarrow L(h(x), y) : h \in \mathcal{H}\}$ . In particular, we scale the loss function  $L$  to  $[0, 1]$  by dividing by  $C$  and denote the new class of scaled loss functions as  $\mathcal{L} \circ \mathcal{H}/C$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , we have:

$$\begin{aligned}
\frac{\mathbb{E}(P_d, h)}{C} & \leq \frac{\widehat{\mathbb{E}}(P_d, h)}{C} + 2\mathcal{R}_S(\mathcal{L} \circ \mathcal{H}/C) + 3\sqrt{\frac{\log(2/\delta)}{2n}} \\
& \stackrel{(1)}{=} \frac{\widehat{\mathbb{E}}(P_d, h)}{C} + \frac{2}{C}\mathcal{R}_S(\mathcal{L} \circ \mathcal{H}) + 3\sqrt{\frac{\log(2/\delta)}{2n}}
\end{aligned} \tag{15}$$

We have  $\stackrel{(1)}{=}$  by using the property of empirical Rademacher complexity that  $\mathcal{R}_S(\alpha\mathcal{F}) = \alpha\mathcal{R}_S(\mathcal{F})$ . We derive Lemma 4 by multiplying Eq. (15) by  $C$ .

**Lemma 5** *Given a loss function  $L$  satisfied Assumption 1 and a sample  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$  of size  $n$ , then an empirical Rademacher complexity  $\mathcal{R}_S(\mathcal{L} \circ \mathcal{H})$  computed from the sample  $S$  is upper bounded as follows.*

$$\mathcal{R}_S(\mathcal{L} \circ \mathcal{H}) \leq C\sqrt{\frac{2d \log n + 4d \log |\mathcal{Y}|}{n}}$$

where  $d$  is Natarajan dimension of hypothesis class  $\mathcal{H}$ .

**Proof of Lemma 5** We have:

$$\begin{aligned}
\mathcal{R}_S(\mathcal{L} \circ \mathcal{H}) & \stackrel{(1)}{\leq} C\sqrt{\frac{2 \log |\mathcal{L} \circ \mathcal{H}|}{n}} \\
& \stackrel{(2)}{\leq} C\sqrt{\frac{2d \log n + 4d \log |\mathcal{Y}|}{n}}
\end{aligned}$$

We have  $\stackrel{(1)}{\leq}$  by applying Massart's finite lemma (Lemma 5 in Liang (2016)) for function class  $\mathcal{L} \circ \mathcal{H}$  and note that  $\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i)^2 \leq C^2$  because of Assumption 2;  $\stackrel{(2)}{\leq}$  by using the fact that  $|\mathcal{L} \circ \mathcal{H}| \leq |\mathcal{H}|$  and then applying Natarajan lemma (Lemma 29.4 in Shalev-Shwartz and Ben-David (2014)) for hypothesis class  $\mathcal{H}$  with Natarajan dimension  $d$ .



## A.2 Proof of main theorems

**Proof of Theorem 1** We have:

$$\begin{aligned}
\mathbb{E}(P_t, h \circ f_t) &= \mathbb{E}_{P_t(Y,Z)} [L(h(Z), Y)] \\
&= \int \int L(h(z), y) p_t(y, z) dy dz \\
&= \int \int L(h(z), y) (p_t(y, z | y \in \mathcal{Y}^s) p_t(y \in \mathcal{Y}^s) + p_t(y, z | y \notin \mathcal{Y}^s) p_t(y \notin \mathcal{Y}^s)) dy dz \\
&= \lambda \mathbb{E}_{P_{t,k}(Y,Z)} [L(h(Z), Y)] + (1 - \lambda) \mathbb{E}_{P_{t,u}(Y,Z)} [L(h(Z), Y)] \\
&= \lambda \mathbb{E}(P_{t,k}, h \circ f_t) + (1 - \lambda) \mathbb{E}(P_{t,u}, h \circ f_t)
\end{aligned} \tag{16}$$

Applying Lemma 1 for the two distributions  $P_{t,k}(Z, Y)$  and  $P_s(Z, Y)$ , we have:

$$\mathbb{E}(P_{t,k}, h \circ f_t) \leq \mathbb{E}(P_s, h \circ f_s) + \sqrt{2}C (\mathcal{D}_{JS}(P_s(Y, Z), P_{t,k}(Y, Z)))^{1/2} \tag{17}$$

Let  $P_{t,k}^u$  is the distribution induced from  $P_{t,k}$  by the mapping  $f_t^u$ . Then, we have:

$$\begin{aligned}
\mathbb{E}(P_t^u, h \circ f_t) &= \mathbb{E}_{P_t^u(Y,Z)} [L(h(Z), Y)] \\
&= \int \int L(h(z), y) p_t^u(y, z) dy dz \\
&\stackrel{(1)}{=} \int \int L(h(z), y) (p_{t,k}^u(y, z) p_t(y \in \mathcal{Y}^s) + p_{t,u}(y, z) p_t(y \notin \mathcal{Y}^s)) dy dz \\
&= \lambda \mathbb{E}_{P_{t,k}^u(Y,Z)} [L(h(Z), Y)] + (1 - \lambda) \mathbb{E}_{P_{t,u}(Y,Z)} [L(h(Z), Y)] \\
&= \lambda \mathbb{E}(P_{t,k}^u, h \circ f_t) + (1 - \lambda) \mathbb{E}(P_{t,u}, h \circ f_t) \\
&\stackrel{(2)}{\geq} \lambda \left( \mathbb{E}(P_s^u, h \circ f_s) - \sqrt{2}C (\mathcal{D}_{JS}(P_{t,k}^u(Y, Z) \| P_s^u(Y, Z)))^{1/2} \right) + (1 - \lambda) \mathbb{E}(P_{t,u}, h \circ f_t) \\
&\stackrel{(3)}{\geq} \lambda \left( \mathbb{E}(P_s^u, h \circ f_s) - \sqrt{2}C \left( \mathcal{D}_{JS}(P_{t,k}^u(Z) \| P_s^u(Z)) + \mathbb{E}_{P_{t,k}^u(Z)} [P_{t,k}^u(Y|Z) \| P_s^u(Y|Z)] \right. \right. \\
&\quad \left. \left. + \mathbb{E}_{P_s^u(Z)} [P_{t,k}^u(Y|Z) \| P_s^u(Y|Z)] \right)^{1/2} \right) + (1 - \lambda) \mathbb{E}(P_{t,u}, h \circ f_t) \\
&\stackrel{(4)}{=} \lambda \mathbb{E}(P_s^u, h \circ f_s) - \sqrt{2}\lambda C (\mathcal{D}_{JS}(P_{t,k}^u(Z) \| P_s^u(Z)))^{1/2} + (1 - \lambda) \mathbb{E}(P_{t,u}, h \circ f_t)
\end{aligned} \tag{18}$$

We have  $\stackrel{(1)}{=}$  by using the fact that  $P_t(Y, Z) = \lambda P_{t,k}(Y, Z) + (1 - \lambda) P_{t,u}(Y, Z)$  and  $f_t^u$  is the mapping such that  $f_t^u(X^t, Y) = (X^t, unk)$ ;  $\stackrel{(2)}{\geq}$  by using Lemma 1;  $\stackrel{(3)}{\geq}$  by using Lemma 2;  $\stackrel{(4)}{=}$  because  $P_{t,k}^u(Y|Z) \| P_s^u(Y|Z)$  (support of  $Y = \{unk\}$ ). Finally, by combining Eq. (16), Eq. (17), and Eq. (18), we have:

$$\begin{aligned}
\mathbb{E}(P_t, h \circ f_t) &\leq \underbrace{\lambda \mathbb{E}(P_s, h \circ f_s)}_{\text{source error}} + \underbrace{\mathbb{E}(P_t^u, h \circ f_t) - \lambda \mathbb{E}(P_s^u, h \circ f_s)}_{\text{open-set difference}} \\
&\quad + \underbrace{\sqrt{2}\lambda C \left( (\mathcal{D}_{JS}(P_s(Z) \| P_{t,k}(Z)))^{\frac{1}{2}} + (\mathcal{D}_{JS}(P_s(Z, Y) \| P_{t,k}(Z, Y)))^{\frac{1}{2}} \right)}_{\text{domain distance}}
\end{aligned}$$

**Proof of Proposition 1** We have:

$$\begin{aligned}
\mathbb{E}(P_t, h \circ f_t) &= \mathbb{E}_{P_t(Y,Z)} [L(h(Z), Y)] \\
&= \int \int L(h(z), y) p_t(y, z) dy dz \\
&= \int \int L(h(z), y) (p_t(y, z | y \in \mathcal{Y}^s) p_t(y \in \mathcal{Y}^s) + p_t(y, z | y \notin \mathcal{Y}^s) p_t(y \notin \mathcal{Y}^s)) dy dz \\
&= \lambda \mathbb{E}_{P_{t,k}(Y,Z)} [L(h(Z), Y)] + (1 - \lambda) \mathbb{E}_{P_{t,u}(Y,Z)} [L(h(Z), Y)] \\
&= \lambda \mathbb{E}(P_{t,k}, h \circ f_t) + (1 - \lambda) \mathbb{E}(P_{t,u}, h \circ f_t)
\end{aligned} \tag{19}$$

Next, applying Lemma 1 for the two distributions  $P_{t,u}(Y, Z)$  and  $P_s^u(Y, Z)$ , we have:

$$\begin{aligned}
\mathbb{E}(P_s^u, h \circ f_s) &\leq \mathbb{E}(P_{t,u}, h \circ f_t) + \sqrt{2}C (\mathcal{D}_{JS}(P_s^u(Y, Z) \parallel P_{t,u}(Y, Z)))^{1/2} \\
&\stackrel{(1)}{\leq} \mathbb{E}(P_{t,u}, h \circ f_t) + \sqrt{2}C (\mathcal{D}_{JS}(P_s^u(Z) \parallel P_{t,u}(Z))) \\
&\quad + \mathbb{E}_{P_s^u(Z)} [\mathcal{D}_{JS}(P_s^u(Y|Z) \parallel P_{t,u}(Y|Z))] + \mathbb{E}_{P_{t,u}(Z)} [\mathcal{D}_{JS}(P_s^u(Y|Z) \parallel P_{t,u}(Y|Z))]^{1/2} \\
&\stackrel{(2)}{=} \mathbb{E}(P_{t,u}, h \circ f_t) + \sqrt{2}C (\mathcal{D}_{JS}(P_s(Z) \parallel P_{t,u}(Z)))^{1/2}
\end{aligned} \tag{20}$$

We have  $\stackrel{(1)}{\leq}$  by applying Lemma 2;  $\stackrel{(2)}{=}$  because  $P_s^u(Y|Z) = P_{t,u}(Y|Z)$  (support of  $Y = \{unk\}$ ) and  $P_s^u(Z) = P_s(Z)$ . Combining Eq. (19) and Eq. (20), we have:

$$\mathbb{E}(P_t, h \circ f_t) \geq \lambda \mathbb{E}(P_{t,u}, h \circ f_t) + (1 - \lambda) \mathbb{E}(P_s^u, h \circ f_s) - \sqrt{2}(1 - \lambda)C (\mathcal{D}_{JS}(P_s(Z) \parallel P_{t,u}(Z)))^{1/2}$$

**Proof of Proposition 2** Before giving the proof, we first introduce following definition about adversarial learning for invariant representation.

**Definition 1** Given dataset  $D_s = \{x_i^s\}_{i=1}^{n_s}$  and  $D_t = \{x_i^t\}_{i=1}^{n_t}$  associated with the distributions  $P_s(X_s)$  and  $P_t(X_t)$ , respectively, the goal of adversarial learning approach for invariant representation is to achieve  $\hat{L}_{adv} = \inf_{\alpha, \beta} \sup_{\gamma} \left( \frac{1}{n_s} \sum_{i=1}^{n_s} \log(D_{\gamma}(F_{\alpha}(x_i^s))) + \frac{1}{n_t} \sum_{i=1}^{n_t} \log(1 - D_{\gamma}(F_{\beta}(x_i^t))) \right)$  where  $F_{\alpha}, F_{\beta}$  are the mappings from the feature spaces  $\mathcal{X}^s, \mathcal{X}^t$  to the representation space  $\mathcal{Z}$  parameterized by  $\alpha \in \mathcal{A}$  and  $\beta \in \mathcal{B}$ , and  $D_{\gamma}$  are the discriminator parameterized by  $\gamma \in \Gamma$ .

Then Proposition 2 are formally state as follows.

**Proposition 2 (Formal).** Let  $\alpha^*, \beta^*, \gamma^*$  be the parameters learned by optimizing  $L_{adv}$  and  $\hat{\alpha}, \hat{\beta}, \hat{\gamma}$  be the parameters learned by optimizing  $\hat{L}_{adv}$ . We have:

$$\begin{aligned}
\mathbb{E} \left[ \mathcal{D}_{JS} \left( P_{\hat{\alpha}}(Z) \parallel P_{\hat{\beta}}(Z) \right) \right] &\leq \mathcal{D}_{JS}(P_{\alpha^*}(Z) \parallel P_{\beta^*}(Z)) \\
&\quad + \mathcal{O} \left( \left( \frac{1}{\sqrt{n_s}} + \frac{1}{\sqrt{n_t}} \right) \times C(\mathcal{A}, \mathcal{B}, \Gamma) \right)
\end{aligned}$$

where  $C(\mathcal{A}, \mathcal{B}, \Gamma)$  is a constant specified by the parameter spaces  $\mathcal{A}, \mathcal{B}, \Gamma$ , and  $L_{adv} = \inf_{\alpha, \beta} \sup_{\gamma} \int_{\mathcal{X}^s} \log(D_{\gamma}(F_{\alpha}(x^s))) p_s(x^s) dx^s + \int_{\mathcal{X}^t} \log(1 - D_{\gamma}(F_{\beta}(x^t))) p_t(x^t) dx^t$  is the objective function of adversarial learning for invariant representation with infinite data.

The proof for Proposition 2 is based on the proof provided for GAN model by Biau et al. (2020). Let  $L(\alpha, \beta, \gamma) =$

$\int_{\mathcal{Z}} (\log(D_\gamma(z)) p_\alpha(z) + \log(1 - D_\gamma(z)) p_\beta(z)) dz$ , we have:

$$\begin{aligned}
2\mathcal{D}_{JS}(P_{\hat{\alpha}}(Z) \parallel P_{\hat{\beta}}(Z)) &= L(\hat{\alpha}, \hat{\beta}, \hat{\gamma}) + \log(4) \\
&\leq \sup_{\gamma} L(\hat{\alpha}, \hat{\beta}, \gamma) + \log(4) \\
&\leq \sup_{\gamma} \left( \widehat{L}(\hat{\alpha}, \hat{\beta}, \gamma) + \left| \widehat{L}(\hat{\alpha}, \hat{\beta}, \gamma) - L(\hat{\alpha}, \hat{\beta}, \gamma) \right| \right) + \log(4) \\
&\leq \sup_{\gamma} \widehat{L}(\hat{\alpha}, \hat{\beta}, \gamma) + \sup_{\gamma} \left| \widehat{L}(\hat{\alpha}, \hat{\beta}, \gamma) - L(\hat{\alpha}, \hat{\beta}, \gamma) \right| + \log(4) \\
&\leq \inf_{\alpha, \beta} \sup_{\gamma} \widehat{L}(\alpha, \beta, \gamma) + \sup_{\alpha, \beta, \gamma} \left| \widehat{L}(\alpha, \beta, \gamma) - L(\alpha, \beta, \gamma) \right| + \log(4) \\
&\leq \inf_{\alpha, \beta} \sup_{\gamma} L(\alpha, \beta, \gamma) + \left| \inf_{\alpha, \beta} \sup_{\gamma} \widehat{L}(\alpha, \beta, \gamma) - \inf_{\alpha, \beta} \sup_{\gamma} L(\alpha, \beta, \gamma) \right| \\
&\quad + \sup_{\alpha, \beta, \gamma} \left| \widehat{L}(\alpha, \beta, \gamma) - L(\alpha, \beta, \gamma) \right| + \log(4) \\
&\stackrel{(1)}{\leq} \inf_{\alpha, \beta} \sup_{\gamma} L(\alpha, \beta, \gamma) + \sup_{\alpha, \beta} \left| \sup_{\gamma} \widehat{L}(\alpha, \beta, \gamma) - \sup_{\gamma} L(\alpha, \beta, \gamma) \right| \\
&\quad + \sup_{\alpha, \beta, \gamma} \left| \widehat{L}(\alpha, \beta, \gamma) - L(\alpha, \beta, \gamma) \right| + \log(4) \\
&\stackrel{(2)}{\leq} \inf_{\alpha, \beta} \sup_{\gamma} L(\alpha, \beta, \gamma) + 2 \sup_{\alpha, \beta, \gamma} \left| \widehat{L}(\alpha, \beta, \gamma) - L(\alpha, \beta, \gamma) \right| + \log(4) \\
&= 2\mathcal{D}_{JS}(P_{\alpha^*}(Z) \parallel P_{\beta^*}(Z)) + 2 \sup_{\alpha, \beta, \gamma} \left| \widehat{L}(\alpha, \beta, \gamma) - L(\alpha, \beta, \gamma) \right|
\end{aligned}$$

We have  $\stackrel{(1)}{\leq}$  by using inequality  $|\inf A - \inf B| \leq \sup |A - B|$ ,  $\stackrel{(2)}{\leq}$  by using inequality  $|\sup A - \sup B| \leq \sup |A - B|$ . Take the expectation and rearrange the both sides, we have:

$$\begin{aligned}
&\mathbb{E} \left[ \mathcal{D}_{JS}(P_{\hat{\alpha}}(Z) \parallel P_{\hat{\beta}}(Z)) \right] - \mathcal{D}_{JS}(P_{\alpha^*}(Z) \parallel P_{\beta^*}(Z)) \\
&\leq \mathbb{E} \left[ \sup_{\alpha, \beta, \gamma} \left| \widehat{L}(\alpha, \beta, \gamma) - L(\alpha, \beta, \gamma) \right| \right] \\
&= \mathbb{E} \left[ \sup_{\alpha, \beta, \gamma} \left| \frac{1}{n_s} \sum_{i=1}^{n_s} \log(D_\gamma(z_i^s)) + \frac{1}{n_t} \sum_{i=1}^{n_t} \log(1 - D_\gamma(z_i^t)) \right. \right. \\
&\quad \left. \left. - \int_{\mathcal{Z}} (\log(D_\gamma(z)) p_\alpha(z) + \log(1 - D_\gamma(z)) p_\beta(z)) dz \right| \right] \\
&\leq \mathbb{E} \left[ \sup_{\alpha, \beta, \gamma} \left| \underbrace{\frac{1}{n_s} \sum_{i=1}^{n_s} \log(D_\gamma(z_i^s)) - \int_{\mathcal{Z}} (\log(D_\gamma(z)) p_\alpha(z)) dz}_{A_s(\alpha, \beta, \gamma)} \right| \right] \\
&\quad + \mathbb{E} \left[ \sup_{\alpha, \beta, \gamma} \left| \underbrace{\frac{1}{n_t} \sum_{i=1}^{n_t} \log(1 - D_\gamma(z_i^t)) - \int_{\mathcal{Z}} (\log(1 - D_\gamma(z)) p_\beta(z)) dz}_{A_t(\alpha, \beta, \gamma)} \right| \right]
\end{aligned}$$

Note that  $(A_s(\alpha, \beta, \gamma))_{\alpha \in \mathcal{A}, \beta \in \mathcal{B}, \gamma \in \Gamma}$  and  $(A_t(\alpha, \beta, \gamma))_{\alpha \in \mathcal{A}, \beta \in \mathcal{B}, \gamma \in \Gamma}$  are the subgaussian processes in the metric spaces  $(\mathcal{A} \times \mathcal{B} \times \Gamma, C_1 \|\cdot\| / \sqrt{n_s})$  and  $(\mathcal{A} \times \mathcal{B} \times \Gamma, C_1 \|\cdot\| / \sqrt{n_t})$  where  $C_1$  is a constant and  $\|\cdot\|$  is the Euclidean norm on  $\mathcal{A} \times \mathcal{B} \times \Gamma$ .

Then using Dudley's entropy integral, we have:

$$\begin{aligned}
& \mathbb{E} \left[ \mathcal{D}_{JS} \left( P_{\hat{\alpha}}(Z) \parallel P_{\hat{\beta}}(Z) \right) \right] - \mathcal{D}_{JS} \left( P_{\alpha^*}(Z) \parallel P_{\beta^*}(Z) \right) \\
& \leq \mathbb{E} \left[ \sup_{\alpha, \beta, \gamma} A_s(\alpha, \beta, \gamma) \right] + \mathbb{E} \left[ \sup_{\alpha, \beta, \gamma} A_t(\alpha, \beta, \gamma) \right] \\
& \leq 12 \int_0^\infty \left( \sqrt{\log N(\mathcal{A} \times \mathcal{B} \times \Gamma, C \|\cdot\| / \sqrt{n_s}, \epsilon)} + \sqrt{\log N(\mathcal{A} \times \mathcal{B} \times \Gamma, C \|\cdot\| / \sqrt{n_t}, \epsilon)} \right) d\epsilon \\
& = 12C_1 \left( \frac{1}{\sqrt{n_s}} + \frac{1}{\sqrt{n_t}} \right) \int_0^\infty \sqrt{\log N(\mathcal{A} \times \mathcal{B} \times \Gamma, \|\cdot\|, \epsilon)} d\epsilon \\
& \stackrel{(3)}{=} 12C_1 \left( \frac{1}{\sqrt{n_s}} + \frac{1}{\sqrt{n_t}} \right) \int_0^{\text{diam}(\mathcal{A} \times \mathcal{B} \times \Gamma)} \sqrt{\log N(\mathcal{A} \times \mathcal{B} \times \Gamma, \|\cdot\|, \epsilon)} d\epsilon \\
& \stackrel{(4)}{\leq} 12C_1 \left( \frac{1}{\sqrt{n_s}} + \frac{1}{\sqrt{n_t}} \right) \int_0^{\text{diam}(\mathcal{A} \times \mathcal{B} \times \Gamma)} \sqrt{\log \left( \left( \frac{2C_2 \sqrt{\text{dim}(\mathcal{A} \times \mathcal{B} \times \Gamma)}}{\epsilon} \right)^{\text{dim}(\mathcal{A} \times \mathcal{B} \times \Gamma)} \right)} d\epsilon \\
& = \mathcal{O} \left( \left( \frac{1}{\sqrt{n_s}} + \frac{1}{\sqrt{n_t}} \right) \times C(\mathcal{A}, \mathcal{B}, \Gamma) \right)
\end{aligned}$$

where  $\text{diam}(\cdot)$  and  $\text{dim}(\cdot)$  are the diameter and the dimension of the metric space, and  $C(\mathcal{A}, \mathcal{B}, \Gamma)$  is the function of  $\text{diam}(\mathcal{A} \times \mathcal{B} \times \Gamma)$  and  $\text{dim}(\mathcal{A} \times \mathcal{B} \times \Gamma)$ . We have  $\stackrel{(3)}{=}$  because  $N(\mathcal{A} \times \mathcal{B} \times \Gamma, \|\cdot\|, \epsilon) = 1$  for  $\epsilon > \text{diam}(\mathcal{A} \times \mathcal{B} \times \Gamma)$ ,  $\stackrel{(4)}{\leq}$  by using inequality  $N(\mathcal{T}, \|\cdot\|, \epsilon) \leq \left( \frac{2C_2 \sqrt{d}}{\epsilon} \right)^d$  where  $\mathcal{T}$  lied in Euclidean space  $\mathbb{R}^d$  is the set of vectors whose length is at most  $C_2$ .

**Proof of Theorem 2** We have:

$$\begin{aligned}
\mathcal{D}_{JS} (P_s(Z, Y) \parallel P_{t,k}(Z, Y)) & \stackrel{(1)}{\leq} \mathcal{D}_{JS} (P_s(Z, Y) \parallel P_{t,k}(Z, g(Z))) + \mathcal{D}_{JS} (P_{t,k}(Z, g(Z)) \parallel P_{t,k}(Z, Y)) \\
& \stackrel{(2)}{\leq} \mathcal{D}_{JS} (P_s(Z, Y) \parallel P_{t,k}(Z, g(Z))) + \mathcal{D}_{JS} (P_{t,k}(Z) \parallel P_{t,k}(Z)) \\
& \quad + \mathbb{E}_{P_{t,k}(Z)} [\mathcal{D}_{JS} (P_{t,k}(g(Z)|Z) \parallel P_{t,k}(Y|Z))] \\
& = \mathcal{D}_{JS} (P_s(Z, Y) \parallel P_{t,k}(Z, g(Z))) + \mathcal{N} (P_{t,k}, g)
\end{aligned} \tag{21}$$

We have  $\stackrel{(1)}{\leq}$  by using triangle inequality for JS-divergence;  $\stackrel{(2)}{\leq}$  by applying Lemma 2. According to Theorem 1, we have:

$$\begin{aligned}
\mathbb{E} (P_t, h \circ f_t) & \leq \lambda \mathbb{E} (P_s, h \circ f_s) + \mathbb{E} (P_t^u, h \circ f_t) - \lambda \mathbb{E} (P_s^u, h \circ f_s) \\
& \quad + \sqrt{2} \lambda C \left( (\mathcal{D}_{JS} (P_s(Z) \parallel P_{t,k}(Z)))^{\frac{1}{2}} + (\mathcal{D}_{JS} (P_s(Z, Y) \parallel P_{t,k}(Z, Y)))^{\frac{1}{2}} \right) \\
& \stackrel{(1)}{\leq} \lambda \mathbb{E} (P_s, h \circ f_s) + \mathbb{E} (P_t^u, h \circ f_t) - \lambda \mathbb{E} (P_s^u, h \circ f_s) \\
& \quad + \sqrt{2} \lambda C \left( (\mathcal{D}_{JS} (P_s(Z) \parallel P_{t,k}(Z)))^{\frac{1}{2}} + (\mathcal{D}_{JS} (P_s(Z, Y) \parallel P_{t,k}(Z, g(Z))))^{\frac{1}{2}} + (\mathcal{N} (P_{t,k}, g))^{\frac{1}{2}} \right)
\end{aligned} \tag{22}$$

We have  $\stackrel{(1)}{\leq}$  by applying Eq. (21) and using inequality  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ . Finally, by applying Lemma 4 for  $\mathbb{E} (P_s, h \circ f_s)$ ,  $\mathbb{E} (P_t^u, h \circ f_t)$ , and  $\mathbb{E} (P_s^u, h \circ f_s)$  in Eq. (22), then with probability at least  $1 - \delta$  over the choice of source and target datasets  $D_s$  and  $D_t$ , we have:

$$\begin{aligned}
\mathbb{E}(P_t, h \circ f_t) &\leq \lambda \widehat{\mathbb{E}}(P_s, h \circ f_s) + \widehat{\mathbb{E}}(P_t^u, h \circ f_t) - \lambda \widehat{\mathbb{E}}(P_s^u, h \circ f_s) \\
&\quad + \sqrt{2} \lambda C \left( (\mathcal{D}_{JS}(P_s(Z) \parallel P_{t,k}(Z)))^{\frac{1}{2}} + (\mathcal{D}_{JS}(P_s(Z, Y) \parallel P_{t,k}(Z, g(Z))))^{\frac{1}{2}} + (\mathbb{N}(P_{t,k}, g))^{\frac{1}{2}} \right) \\
&\quad + 2\lambda \mathcal{R}_{D_s}(\mathcal{L} \circ \mathcal{H} \circ \mathcal{F}_s) + 2\mathcal{R}_{D_t^u}(\mathcal{L} \circ \mathcal{H} \circ \mathcal{F}_t) + 2\lambda \mathcal{R}_{D_s^u}(\mathcal{L} \circ \mathcal{H} \circ \mathcal{F}_s) + \mathcal{O} \left( C \sqrt{\log(1/\delta)} \left( \frac{\lambda}{\sqrt{n_s}} + \frac{1}{\sqrt{n_t}} \right) \right) \\
&\stackrel{(1)}{\leq} \lambda \widehat{\mathbb{E}}(P_s, h \circ f_s) + \widehat{\mathbb{E}}(P_t^u, h \circ f_t) - \lambda \widehat{\mathbb{E}}(P_s^u, h \circ f_s) \\
&\quad + \sqrt{2} \lambda C \left( (\mathcal{D}_{JS}(P_s(Z) \parallel P_{t,k}(Z)))^{\frac{1}{2}} + (\mathcal{D}_{JS}(P_s(Z, Y) \parallel P_{t,k}(Z, g(Z))))^{\frac{1}{2}} + (\mathbb{N}(P_{t,k}, g))^{\frac{1}{2}} \right) \\
&\quad + \mathcal{O} \left( \lambda C \sqrt{\frac{d_s \log n_s + d_s \log |\mathcal{Y}^t| + \log \frac{1}{\delta}}{n_s}} + C \sqrt{\frac{d_t \log n_t + d_t \log |\mathcal{Y}^t| + \log \frac{1}{\delta}}{n_t}} \right) \tag{23}
\end{aligned}$$

where  $D_s^u$  and  $D_t^u$  are datasets induced from  $D_s$  and  $D_t$  using mappings  $f_s^u$  and  $f_t^u$ , respectively. We have  $\stackrel{(1)}{\leq}$  by applying Lemma 5 for  $\mathcal{R}_{D_s}(\mathcal{L} \circ \mathcal{H} \circ \mathcal{F}_s)$ ,  $\mathcal{R}_{D_t^u}(\mathcal{L} \circ \mathcal{H} \circ \mathcal{F}_t)$ , and  $\mathcal{R}_{D_s^u}(\mathcal{L} \circ \mathcal{H} \circ \mathcal{F}_s)$ .

**Proof of Proposition 3** We have:

$$\begin{aligned}
\mathbb{E}(P_t, h \circ f) &\stackrel{(1)}{\leq} \mathbb{E}(P_s, h \circ f) + \sqrt{2} C (\mathcal{D}_{JS}(P_s(Y, Z) \parallel P_t(Y, Z)))^{1/2} \\
&\stackrel{(2)}{\leq} \mathbb{E}(P_s, h \circ f) + \sqrt{2} C (\mathcal{D}_{JS}(P_s(Z) \parallel P_t(Z)) + 2\mathbb{E}_{P_{s,t}}[\mathcal{D}_{JS}(P_s(Y|Z) \parallel P_t(Y|Z))])^{1/2} \\
&\stackrel{(3)}{=} \mathbb{E}(P_s, h \circ f) + \sqrt{2} C (\mathcal{D}_{JS}(P_s(Z) \parallel P_t(Z)) + 2\mathbb{E}_{P_{s,t}}[\mathcal{D}_{JS}(P_s(Y|Z) \parallel P_s(Y|Z))])^{1/2} \\
&= \mathbb{E}(P_s, h \circ f) + \sqrt{2} C (\mathcal{D}_{JS}(P_s(Z) \parallel P_t(Z)))^{1/2}
\end{aligned}$$

We have  $\stackrel{(1)}{\leq}$  by applying Lemma 1 for two distributions  $P_s(Z, Y)$  and  $P_t(Z, Y)$ , and note that  $\mathbb{E}(P_d, h \circ f) = \mathbb{E}_{P_d(X, Y)}[L(h(f(X)), Y)] = \mathbb{E}_{P_d(Z, Y)}[L(h(Z), Y)] = \mathbb{E}(P_d, h)$  for  $d \in \{s, t\}$ ;  $\stackrel{(2)}{\leq}$  by applying Lemma 2 for  $\mathcal{D}_{JS}(P_s(Y, Z) \parallel P_t(Y, Z))$ ;  $\stackrel{(3)}{=}$  by applying Lemma 3.

## B Model Details

In this section, we provide a comprehensive overview of the architectures utilized in our experiment, along with the pseudocode for training our proposed method.

### B.1 Backbone Architecture Details

To ensure fair comparisons between methods, we use the same backbone architectures across all approaches. Most methods involve representation mapping and classifier networks, except for SSAN and SCT, which include an additional discriminator network, and KPG, which employs optimal transport to transform data within the input space. For representation mapping, we implement multi-layer feed-forward neural networks on the CIFAR10 & ILSVRC2012, Wikipedia, Multilingual Reuters Collection, NUSWIDE & ImageNet, Office & Caltech256, and ImageCLEF-DA datasets, while using ResNet/WideResNet to encode paper and digital ECG data in the PTB-XL dataset. A linear classifier is applied across all datasets. The specifics of these networks are detailed in Tables 3-6 below.

Table 3: Details of backbone networks used in CIFAR10 & ILSVRC2012, Wikipedia, Multilingual Reuters Collection, NUSWIDE & ImageNet, Office & Caltech256, and ImageCLEF-DA datasets.  $\mathbf{d\_source}$ ,  $\mathbf{d\_target}$ , and  $\mathbf{n\_output}$  represent the dimensions of the source feature space, target feature space, and output space, respectively. For RL-OSHeDA and OPDA,  $\mathbf{n\_output}$  is set to  $|\mathcal{Y}^t|$  while for other methods,  $\mathbf{n\_output}$  is set to  $|\mathcal{Y}^s|$ .

Networks	Layers
Representation Mapping $f_s$	Linear(input dim= $\mathbf{d\_source}$ , output dim= $(\mathbf{d\_source} + 256)/2$ )
	LeakyReLU(negative_slope=0.2)
	Linear(input dim= $(\mathbf{d\_source} + 256)/2$ , output dim=256)
Representation Mapping $f_t$	LeakyReLU(negative_slope=0.2)
	Normalize(p=2)
	Linear(input dim= $\mathbf{d\_target}$ , output dim= $(\mathbf{d\_target} + 256)/2$ )
Representation Mapping $f_t$	LeakyReLU(negative_slope=0.2)
	Linear(input dim= $(\mathbf{d\_target} + 256)/2$ , output dim=256)
	LeakyReLU(negative_slope=0.2)
Classifier $h$	Normalize(p=2)
	Linear(input dim=256, output dim= $\mathbf{n\_output}$ )
	LeakyReLU(negative_slope=0.2)



Table 4: Details of backbone networks used in PTB-XL dataset. The representation mapping networks are constructed from multiple ResNetBlock1d (for digital ECG) and ResNetBlock2d (for paper ECG) sub-modules. **n\_output** represent the dimensions of output space. For RL-OSHeDA and OPDA, **n\_output** is set to  $|\mathcal{Y}^t|$  while for other methods, **n\_output** is set to  $|\mathcal{Y}^s|$ .

Networks	Layers
ResNetBlock1d(i_ch, o_ch)	Conv1d(input channel= <b>i_ch</b> , output channel= <b>o_ch</b> , kernel=3, stride= <b>o_ch</b> / <b>i_ch</b> , padding=1)
	BatchNorm1d
	ReLU
ResNetBlock1d(i_ch, o_ch)	Conv1d(input channel= <b>o_ch</b> , output channel= <b>o_ch</b> , kernel=3, padding=1)
	BatchNorm1d
	ReLU
ResNetBlock2d(i_ch, o_ch)	Conv2d(input channel= <b>i_ch</b> , output channel= <b>o_ch</b> , kernel=3, stride= <b>o_ch</b> / <b>i_ch</b> , padding=1)
	BatchNorm2d
	ReLU
	Conv2d(input channel= <b>o_ch</b> , output channel= <b>o_ch</b> , kernel=3, padding=1)
ResNetBlock2d(i_ch, o_ch)	BatchNorm2d
	ReLU
	Conv1d(input channel=1, output channel=64, kernel=7, stride=2, padding=3)
	BatchNorm1d
ResNet1d	ReLU
	MaxPool1d
	ResNetBlock1d(i_ch=64, o_ch=64) $\times$ 2
	ResNetBlock1d(i_ch=64, o_ch=128)
	ResNetBlock1d(i_ch=128, o_ch=128)
	ResNetBlock1d(i_ch=128, o_ch=256)
	ResNetBlock1d(i_ch=256, o_ch=256)
	ResNetBlock1d(i_ch=256, o_ch=512)
	ResNetBlock1d(i_ch=512, o_ch=512)
	BatchNorm1d
	ReLU
	AdaptiveAvgPool1d(output size=7)
Linear(input dim=3584, output dim=256)	
Representation Mapping $f_s$	Concatenation of 12 ResNet1d modules
	Linear(input dim=256 $\times$ 12, output dim=256)
Representation Mapping $f_t$	Conv2d(input channel=3, output channel=64, kernel=7, stride=2, padding=3)
	BatchNorm2d
	ReLU
	MaxPool2d
	ResNetBlock2d(i_ch=64, o_ch=64) $\times$ 2
	ResNetBlock2d(i_ch=64, o_ch=128)
	ResNetBlock2d(i_ch=128, o_ch=128)
	ResNetBlock2d(i_ch=128, o_ch=256)
	ResNetBlock2d(i_ch=256, o_ch=256)
	ResNetBlock2d(i_ch=256, o_ch=512)
	ResNetBlock2d(i_ch=512, o_ch=512)
	AdaptiveAvgPool2d(output size=1)
Linear(input dim=512, output dim=256)	
Classifier $h$	Linear(input dim=256, output dim= <b>n_output</b> )
	LeakyReLU(negative_slope=0.2)

Table 5: Details of discriminator network used in SSAN and SCT methods.

Networks	Layers
Discriminator $D$	Linear(input dim=256, output dim=1)
	Sigmoid

Table 6: Details of the WideResNet backbone. This network is employed in DS3L for the PTB-XL dataset, in place of the standard ResNet backbone.

Networks	Layers
ResNetUnit(i_ch, o_ch)	BatchNorm2d
	LeakyReLU(negative_slope=0.1)
	Conv2d(input channel= <b>i_ch</b> , output channel= <b>o_ch</b> , kernel=3, stride= <b>o_ch / i_ch</b> , padding=1)
	BatchNorm2d
	LeakyReLU(negative_slope=0.1)
ResNetBlock(i_ch, o_ch)	Conv2d(input channel= <b>o_ch</b> , output channel= <b>o_ch</b> , kernel=3, stride=1, padding=1)
	ResNetUnit(i_ch= <b>i_ch</b> , o_ch= <b>o_ch</b> )
	ResNetUnit(i_ch= <b>o_ch</b> , o_ch= <b>o_ch</b> ) $\times$ 3
Representation Mapping $f_t$	Conv2d(input channel=3, output channel=16, kernel=3, padding=1)
	ResNetBlock(i_ch=16, o_ch=32)
	ResNetBlock(i_ch=32, o_ch=64)
	ResNetBlock(i_ch=64, o_ch=128)
	BatchNorm2d
	LeakyReLU(negative_slope=0.1)
	AveragePool2d(output size=1)
	Linear(input dim=128, output dim=256)

## B.2 Training Details

We utilize a two-stage learning approach for our proposed method, RL-OSHeDA. Specifically, in stage 1, we update the model parameters by optimizing  $L_{cls}$  as defined in Eq. (2). In stage 2, we update the model parameters by optimizing  $L$  as specified in Eq. (1), with the assistance of pseudo-labels. The details of this training process are outlined in Algorithm 1 below.

Algorithm 1: Two-stage learning process for RL-OSHeDA

---

```

1: Inputs: Source datasets  $D_s = \{x_i^s, y_i^s\}_{i=1}^{n_s}$ , labeled target datasets  $D_{t_l} = \{x_i^t, y_i^t\}_{i=1}^{n_{t_l}}$ , unlabeled target datasets  $D_{t_u} = \{x_i^t\}_{i=1}^{n_{t_u}}$ , known class prior  $\lambda$ , threshold  $T$ , number of epochs  $T_{max}$ 
2: Outputs: Trained network parameters  $\theta_{f_s}, \theta_{f_t}, \theta_h$ 
3: Initialize network parameters  $\theta_{f_s}, \theta_{f_t}, \theta_h$ 
4: for  $epoch = 1$  to  $T_{max}$  do
5:   Sample source minibatch  $B_s$ 
6:   Sample labeled target minibatch  $B_{t_l}$ 
7:   Sample unlabeled target minibatch  $B_{t_u}$ 
8:   if  $epoch < T$  then
9:     /* Stage 1 */
10:    Calculate  $L_{cls}$  in Eq. (2) using  $B_s, B_{t_l}, B_{t_u}, \lambda$ 
11:    Update parameters:
12:     $\theta_{f_s} = \theta_{f_s} - \gamma \nabla_{\theta_{f_s}} L_{cls}$ 
13:     $\theta_{f_t} = \theta_{f_t} - \gamma \nabla_{\theta_{f_t}} L_{cls}$ 
14:     $\theta_h = \theta_h - \gamma \nabla_{\theta_h} L_{cls}$ 
15:   else
16:     /* Stage 2 */
17:     Generate pseudo-labels for  $B_{t_u}$ 
18:     Calculate  $L_{cls}$  in Eq. (2) using  $B_s, B_{t_l}, B_{t_u}, \lambda$ 
19:     Calculate  $L_{inv}$  in Eq. (3) using  $B_s, B_{t_l}, B_{t_u}$ , pseudo-labels
20:     Calculate  $L_{seg}$  in Eq. (4) using  $B_s, B_{t_l}, B_{t_u}$ , pseudo-labels
21:     Calculate  $L_{osd}$  in Eq. (5) using  $B_s, B_{t_l}, B_{t_u}, \lambda$ 
22:     Update parameters:
23:      $\theta_{f_s} = \theta_{f_s} - \gamma \nabla_{\theta_{f_s}} (L_{cls} + L_{inv} - L_{seg} + L_{osd})$ 
24:      $\theta_{f_t} = \theta_{f_t} - \gamma \nabla_{\theta_{f_t}} (L_{cls} + L_{inv} - L_{seg} + L_{osd})$ 
25:      $\theta_h = \theta_h - \gamma \nabla_{\theta_h} (L_{cls} + L_{inv} - L_{seg} + L_{osd})$ 
26:   end if
27: end for
28: return Trained parameters  $\theta_{f_s}, \theta_{f_t}, \theta_h$ 

```

---

## C Experimental Details

### C.1 Datasets

We conduct experiments using seven datasets covering the clinical, computer vision, and natural language processing domains. Detailed descriptions of these datasets are provided below, with corresponding data statistics presented in Tables 7-13.

**CIFAR-10 (Krizhevsky 2009) & ILSVRC2012 (Russakovsky et al. 2015).** These two datasets are used for the image-to-image adaptation task. Big Transfer-M with ResNet-50 and ResNet-101 (Kolesnikov et al. 2019) are utilized to extract features from the images. In the target domain (ILSVRC2012), 4 out of 8 shared classes are designated as unknown classes. For the source (CIFAR-10) and unlabeled target data, 50 instances are randomly selected for each class, and we randomly choose 1, 3, or 5 instances per class as labeled target data. This process results in 6 DA tasks.

**Wikipedia (Rasiwasia et al. 2010).** This dataset, consisting of text-image pairs, is used for image-to-text and text-to-image adaptation tasks. Big Transfer-M with ResNet-101 and Big Bird (Zaheer et al. 2020) are used to extract features for image and text, respectively. 5 out of 10 classes are designated as unknown classes. All data in the source domain are used as labeled source data. For the target domain, we randomly select 5 instances per class as labeled target data and randomly select 50 instances per class from the remaining data as unlabeled target data. This process results in 2 DA tasks.

**Multilingual Reuters Collection (Amini, Usunier, and Goutte 2009).** This dataset, consisting of articles in 5 languages, is used for text-to-text adaptation tasks. Bag-of-Words with TF-IDF, followed by Principal Component Analysis, is used to generate features for each article. English, French, Italian, and German are used as the source domains, while Spanish is used as the target domain. 3 out of 6 classes are designated as unknown classes. For the source and unlabeled target datasets, 100 and 500 instances are randomly selected for each class, respectively, and we randomly choose 20 instances per class as labeled target data. This process results in 4 DA tasks.

**NUS-WIDE (Chua et al. 2009) and ImageNet (Deng et al. 2009).** These datasets are used for the text-to-image adaptation task. We utilize the tag information from NUS-WIDE as the source domain (text) and the image data from ImageNet as the target domain (image). 4 out of 8 shared classes between the two datasets are designated as unknown classes. Features for NUS-WIDE tags are extracted using a pre-trained 5-layer neural network, while DeCAF6 (Donahue et al. 2014) features are used for images in ImageNet. For NUS-WIDE, 100 instances per class are selected, whereas for ImageNet, 3 instances per class are sampled as labeled target data, with all remaining images used as unlabeled target data.

**Office (Saenko et al. 2010) and Caltech-256 (Griffin et al. 2007).** These datasets, which include 4 domains Amazon, Webcam, DSLR from Office, and Caltech from Caltech-256, are used for the image-to-image adaptation task. Amazon, Webcam, and Caltech are used as source domains while Amazon, Webcam, DSLR, and Caltech are used as target domains. SURF (Saenko et al. 2010) and DeCAF6 are utilized as 2 different feature sets for images in these datasets. 5 out of 10 classes are designated as unknown classes, and 3 instances per class are sampled as labeled target data. This process results in 18 DA tasks.

**ImageCLEF-DA (Griffin et al. 2007).** This data, which include 4 domains Caltech, ImageNet, Bing, and PascalVOC, are used for the image-to-image adaptation task. ResNet50 and VGG-19 (Simonyan and Zisserman 2014) are utilized as 2 different feature sets for images in this dataset. 6 out of 12 classes are designated as unknown classes, and 3 instances per class are sampled as labeled target data. This process results in 24 DA tasks.

**PTB-XL (Wagner et al. 2020).** This dataset is used for the digital-to-paper electrocardiogram (ECG) adaptation task. Frequent classes including NORM, Old MI, STTC, and CD are designated as known classes, while the remaining classes are considered unknown. For each known class, we sample 2,000 instances per class to construct source (1,000 instances) and target (1,000 instances) datasets. All instances from the unknown class are used for the target dataset. This process results in 1 DA task.

Table 7: Data statistics for each domain adaptation task in CIFAR10 & ILSVRC2012 dataset.

	CIFAR10 & ILSVRC2012								
	Source Dataset			Labeled Target Dataset			Unlabeled Target Dataset		
	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim
ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (1)	4	2000	2048	4	4	2048	5	4000	2048
ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (3)	4	2000	2048	4	12	2048	5	4000	2048
ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (5)	4	2000	2048	4	20	2048	5	4000	2048
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (1)	4	2000	2048	4	4	2048	5	4000	2048
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (3)	4	2000	2048	4	12	2048	5	4000	2048
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (5)	4	2000	2048	4	20	2048	5	4000	2048

Table 8: Data statistics for each domain adaptation task in Wikipedia dataset.

Wikipedia									
	Source Dataset			Labeled Target Dataset			Unlabeled Target Dataset		
	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim
Image $\Rightarrow$ Text	5	1472	768	5	25	2048	6	500	2048
Text $\Rightarrow$ Image	5	1472	2048	5	25	768	6	500	768

Table 9: Data statistics for each domain adaptation task in Multilingual Reuters Collection dataset.

Multilingual Reuters Collection									
	Source Dataset			Labeled Target Dataset			Unlabeled Target Dataset		
	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim
English $\Rightarrow$ Spanish	3	300	1131	3	60	807	4	2910	807
French $\Rightarrow$ Spanish	3	300	1230	3	60	807	4	2910	807
German $\Rightarrow$ Spanish	3	300	1417	3	60	807	4	2910	807
Italian $\Rightarrow$ Spanish	3	300	1041	3	60	807	4	2910	807

Table 10: Data statistics for each domain adaptation task in NUSWIDE &amp; ImageNet dataset.

NUSWIDE & ImageNet									
	Source Dataset			Labeled Target Dataset			Unlabeled Target Dataset		
	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim
Text $\Rightarrow$ Image	4	400	64	4	12	4096	5	800	4096

Table 11: Data statistics for each domain adaptation task in Office &amp; Caltech256 dataset.

Office & Caltech256									
	Source Dataset			Labeled Target Dataset			Unlabeled Target Dataset		
	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim
Amazon (DeCAF6) $\Rightarrow$ Caltech (SURF)	5	100	4096	5	15	800	6	1093	800
Amazon (DeCAF6) $\Rightarrow$ DSLR (SURF)	5	100	4096	5	15	800	6	127	800
Amazon (DeCAF6) $\Rightarrow$ Webcam (SURF)	5	100	4096	5	15	800	6	265	800
Amazon (SURF) $\Rightarrow$ Caltech (DeCAF6)	5	100	800	5	15	4096	6	1093	4096
Amazon (SURF) $\Rightarrow$ DSLR (DeCAF6)	5	100	800	5	15	4096	6	127	4096
Amazon (SURF) $\Rightarrow$ Webcam (DeCAF6)	5	100	800	5	15	4096	6	265	4096
Caltech (DeCAF6) $\Rightarrow$ Amazon (SURF)	5	100	4096	5	15	800	6	928	800
Caltech (DeCAF6) $\Rightarrow$ DSLR (SURF)	5	100	4096	5	15	800	6	127	800
Caltech (DeCAF6) $\Rightarrow$ Webcam (SURF)	5	100	4096	5	15	800	6	265	800
Caltech (SURF) $\Rightarrow$ Amazon (DeCAF6)	5	100	800	5	15	4096	6	928	4096
Caltech (SURF) $\Rightarrow$ DSLR (DeCAF6)	5	100	800	5	15	4096	6	127	4096
Caltech (SURF) $\Rightarrow$ Webcam (DeCAF6)	5	100	800	5	15	4096	6	265	4096
Webcam (DeCAF6) $\Rightarrow$ Amazon (SURF)	5	100	4096	5	15	800	6	928	800
Webcam (DeCAF6) $\Rightarrow$ Caltech (SURF)	5	100	4096	5	15	800	6	1093	800
Webcam (DeCAF6) $\Rightarrow$ DSLR (SURF)	5	100	4096	5	15	800	6	127	800
Webcam (SURF) $\Rightarrow$ Amazon (DeCAF6)	5	100	800	5	15	4096	6	928	4096
Webcam (SURF) $\Rightarrow$ Caltech (DeCAF6)	5	100	800	5	15	4096	6	1093	4096
Webcam (SURF) $\Rightarrow$ DSLR (DeCAF6)	5	100	800	5	15	4096	6	127	4096

Table 12: Data statistics for each domain adaptation task in ImageCLEF-DA dataset.

	ImageCLEF-DA								
	Source Dataset			Labeled Target Dataset			Unlabeled Target Dataset		
	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim
Bing (Reset-50) ⇒ Caltech (VGG-19)	6	120	2048	6	18	4096	7	564	4096
Bing (Reset-50) ⇒ ImageNet (VGG-19)	6	120	2048	6	18	4096	7	564	4096
Bing (Reset-50) ⇒ PascalVOC (VGG-19)	6	120	2048	6	18	4096	7	564	4096
Bing (VGG-19) ⇒ Caltech (Reset-50)	6	120	4096	6	18	2048	7	564	2048
Bing (VGG-19) ⇒ ImageNet (Reset-50)	6	120	4096	6	18	2048	7	564	2048
Bing (VGG-19) ⇒ PascalVOC (Reset-50)	6	120	4096	6	18	2048	7	564	2048
Caltech (Reset-50) ⇒ Bing (VGG-19)	6	120	2048	6	18	4096	7	564	4096
Caltech (Reset-50) ⇒ ImageNet (VGG-19)	6	120	2048	6	18	4096	7	564	4096
Caltech (Reset-50) ⇒ PascalVOC (VGG-19)	6	120	2048	6	18	4096	7	564	4096
Caltech (VGG-19) ⇒ Bing (Reset-50)	6	120	4096	6	18	2048	7	564	2048
Caltech (VGG-19) ⇒ ImageNet (Reset-50)	6	120	4096	6	18	2048	7	564	2048
Caltech (VGG-19) ⇒ PascalVOC (Reset-50)	6	120	4096	6	18	2048	7	564	2048
ImageNet (Reset-50) ⇒ Bing (VGG-19)	6	120	2048	6	18	4096	7	564	4096
ImageNet (Reset-50) ⇒ Caltech (VGG-19)	6	120	2048	6	18	4096	7	564	4096
ImageNet (Reset-50) ⇒ PascalVOC (VGG-19)	6	120	2048	6	18	4096	7	564	4096
ImageNet (VGG-19) ⇒ Bing (Reset-50)	6	120	4096	6	18	2048	7	564	2048
ImageNet (VGG-19) ⇒ Caltech (Reset-50)	6	120	4096	6	18	2048	7	564	2048
ImageNet (VGG-19) ⇒ PascalVOC (Reset-50)	6	120	4096	6	18	2048	7	564	2048
PascalVOC (Reset-50) ⇒ Bing (VGG-19)	6	120	2048	6	18	4096	7	564	4096
PascalVOC (Reset-50) ⇒ Caltech (VGG-19)	6	120	2048	6	18	4096	7	564	4096
PascalVOC (Reset-50) ⇒ ImageNet (VGG-19)	6	120	2048	6	18	4096	7	564	4096
PascalVOC (VGG-19) ⇒ Bing (Reset-50)	6	120	4096	6	18	2048	7	564	2048
PascalVOC (VGG-19) ⇒ Caltech (Reset-50)	6	120	4096	6	18	2048	7	564	2048
PascalVOC (VGG-19) ⇒ ImageNet (Reset-50)	6	120	4096	6	18	2048	7	564	2048

Table 13: Data statistics for each domain adaptation task in PTB-XL dataset.

	PTB-XL								
	Source Dataset			Labeled Target Dataset			Unlabeled Target Dataset		
	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim	# of classes	# of instances	feature dim
Digital ECG ⇒ Paper ECG	4	4000	12×5000	4	80	3×224×224	5	4602	3×224×224

## C.2 Baselines

We experiment with diverse baselines from heterogeneous domain adaptation, open-set domain adaptation, open-set semi-supervised learning, supervised learning, and semi-supervised learning. The details of these baselines are as follows.

**Heterogeneous Domain Adaptation.** Heterogeneous domain adaptation methods are trained on both source and target data. During inference, these methods classify instances as unknown using the same method as our pseudo-label model  $g$  (see Section 5.2).

- SSAN (Li et al. 2020): This method maps heterogeneous source features into a shared representation space and then makes predictions from this space. To adapt from the source to the target domains, it aligns the marginal and label-conditional representation distributions between the two domains using Maximum Mean Discrepancy (MMD). Additionally, it employs pseudo-labels generated from geometric similarity and hard predictions made by the classifier.
- STN (Yao et al. 2019): This method is similar to SSAN, but it calculates the MMD distance using soft labels (i.e., softmax probabilities) rather than hard labels.
- SCT (Zhao et al. 2022): This method is similar to SSAN but differs in that it aligns marginal and label-conditional representation distributions between source and target domains using cosine similarity. Additionally, it generates pseudo-labels based on geometric distances from the source data.
- KPG (Gu et al. 2022): This method utilizes partial optimal transport and the Gromov-Wasserstein distance to map features from the source domain to the target domain. An SVM, trained on the transported source data and labeled target data, is then used for making predictions.

### Open-Set Domain Adaptation.

- OPDA (Saito et al. 2018): This method is trained exclusively on target data. It employs adversarial training to train a classifier that distinguishes between labeled and unlabeled target samples, while a generator is trained to push the unlabeled samples away from the decision boundary. This setup provides the generator with two options: aligning unlabeled samples with labeled ones or classifying them as unknown. Consequently, this approach enables the extraction of features that effectively differentiate between known and unknown target samples. Unlike other baselines, OPDA can directly classify instances as unknown during inference based on the classifier’s output.

## Open-Set Semi-Supervised Learning.

- DS3L (Guo et al. 2020): This method is trained exclusively on target data and selectively uses unlabeled data while monitoring its impact to mitigate performance risks. Specifically, DS3L weakens the influence of unlabeled data with unknown classes to enhance distribution matching and maintain strong generalization. Simultaneously, it reinforces the use of labeled data to prevent performance degradation. These considerations are integrated into a unified bi-level optimization framework. During inference, DS3L classifies instances as unknown using the same method as our pseudo-label model  $g$ .

## Supervised Learning.

- Supervised Learning (SL): We train the model directly on the labeled target dataset by minimizing the classification loss. During inference, SL classifies instances as unknown using the same method as our pseudo-label model  $g$ .

## Semi-Supervised Learning.

- Pseudo Labeling (PL): Similar to supervised learning, but we use the model to generate pseudo-labels for the unlabeled target data and then incorporate these pseudo-labeled examples into the training process. During inference, PL classifies instances as unknown using the same method as our pseudo-label model  $g$ .

## C.3 Implementation Details

Data, model implementation, and training script are included in the code & data supplementary material. We train each model on each domain adaptation task with 10 different random seeds and report the average prediction performances. All experiments are conducted on a machine with 24-Core CPU, 4 RTX A4000 GPUs, and 128G RAM.

## C.4 Additional Results

In this section, we provide a comprehensive overview of the results for domain adaptation tasks across seven datasets. Detailed results are presented in Tables 14-20. Additionally, Table 21 includes the results of statistical tests used to assess the significance of our method in comparison to the baseline approaches across all domain adaptation tasks.

Table 14: Prediction performances of RL-OSHeDA as well as baselines for all domain adaptation tasks in CIFAR10 & ILSVRC2012 dataset.

	CIFAR10 & ILSVRC2012								
	DS3L			KPG			OPDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (1)	52.69 $\pm$ 0.74	48.53 $\pm$ 0.93	57.77 $\pm$ 1.10	56.92 $\pm$ 0.48	53.95 $\pm$ 0.00	60.33 $\pm$ 1.11	45.42 $\pm$ 0.75	38.96 $\pm$ 0.90	54.83 $\pm$ 1.07
ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (3)	65.87 $\pm$ 0.72	64.34 $\pm$ 0.99	67.49 $\pm$ 1.03	57.25 $\pm$ 0.48	54.98 $\pm$ 0.00	59.77 $\pm$ 1.12	57.34 $\pm$ 0.81	53.38 $\pm$ 1.06	62.25 $\pm$ 1.10
ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (5)	66.98 $\pm$ 0.73	66.02 $\pm$ 1.01	67.97 $\pm$ 1.08	57.32 $\pm$ 0.52	55.73 $\pm$ 0.00	59.02 $\pm$ 1.13	58.35 $\pm$ 0.79	54.74 $\pm$ 1.10	62.65 $\pm$ 1.13
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (1)	52.70 $\pm$ 0.78	47.72 $\pm$ 1.11	59.02 $\pm$ 1.07	52.68 $\pm$ 0.51	47.70 $\pm$ 0.00	59.53 $\pm$ 1.14	44.22 $\pm$ 0.76	36.41 $\pm$ 0.94	56.61 $\pm$ 1.15
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (3)	62.04 $\pm$ 0.73	59.69 $\pm$ 0.98	64.74 $\pm$ 1.06	58.77 $\pm$ 0.51	56.53 $\pm$ 0.00	61.26 $\pm$ 1.09	53.16 $\pm$ 0.73	47.59 $\pm$ 0.94	60.50 $\pm$ 1.08
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (5)	68.67 $\pm$ 0.72	67.95 $\pm$ 0.99	69.42 $\pm$ 1.01	60.70 $\pm$ 0.53	59.52 $\pm$ 0.00	61.92 $\pm$ 1.10	61.28 $\pm$ 0.77	58.23 $\pm$ 1.08	64.73 $\pm$ 1.07
	PL			SCT			SSAN		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
	ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (1)	36.84 $\pm$ 0.61	29.92 $\pm$ 0.67	49.36 $\pm$ 1.12	51.80 $\pm$ 0.78	47.85 $\pm$ 1.03	56.66 $\pm$ 1.16	55.23 $\pm$ 0.76	53.41 $\pm$ 0.97
ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (3)	47.34 $\pm$ 0.57	42.81 $\pm$ 0.64	53.82 $\pm$ 1.10	62.68 $\pm$ 0.75	61.22 $\pm$ 0.97	64.28 $\pm$ 1.11	62.34 $\pm$ 0.71	61.34 $\pm$ 0.98	63.39 $\pm$ 1.05
ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (5)	43.77 $\pm$ 0.55	38.15 $\pm$ 0.49	53.93 $\pm$ 1.10	65.04 $\pm$ 0.74	64.19 $\pm$ 1.00	65.92 $\pm$ 1.07	63.22 $\pm$ 0.74	62.71 $\pm$ 1.04	63.74 $\pm$ 1.09
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (1)	35.49 $\pm$ 0.34	27.96 $\pm$ 0.27	51.08 $\pm$ 1.12	49.74 $\pm$ 0.78	45.12 $\pm$ 1.07	55.76 $\pm$ 1.07	54.36 $\pm$ 0.78	51.13 $\pm$ 1.11	58.58 $\pm$ 1.11
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (3)	44.72 $\pm$ 0.51	39.19 $\pm$ 0.52	52.61 $\pm$ 1.05	60.24 $\pm$ 0.71	58.16 $\pm$ 0.94	62.55 $\pm$ 1.08	60.57 $\pm$ 0.70	59.28 $\pm$ 0.97	62.06 $\pm$ 1.03
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (5)	48.35 $\pm$ 0.51	44.71 $\pm$ 0.38	53.07 $\pm$ 1.10	68.19 $\pm$ 0.72	67.61 $\pm$ 0.96	68.79 $\pm$ 1.02	66.58 $\pm$ 0.70	66.19 $\pm$ 0.95	66.97 $\pm$ 1.03
	STN			SL			OSHeDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
	ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (1)	55.22 $\pm$ 0.75	51.92 $\pm$ 1.01	59.13 $\pm$ 1.08	52.90 $\pm$ 0.74	48.24 $\pm$ 0.96	58.70 $\pm$ 1.14	<b>60.45<math>\pm</math>0.79</b>	<b>56.70<math>\pm</math>1.04</b>
ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (3)	63.51 $\pm$ 0.69	60.68 $\pm$ 0.90	66.81 $\pm$ 1.10	63.97 $\pm$ 0.69	62.39 $\pm$ 0.89	65.65 $\pm$ 1.07	<b>75.41<math>\pm</math>0.70</b>	<b>71.26<math>\pm</math>0.94</b>	<b>80.43<math>\pm</math>0.95</b>
ImageNet (ResNet-101) $\Rightarrow$ CIFAR (ResNet-50) (5)	65.62 $\pm$ 0.75	63.47 $\pm$ 1.00	67.97 $\pm$ 1.09	66.18 $\pm$ 0.72	65.18 $\pm$ 1.02	67.21 $\pm$ 1.07	<b>78.24<math>\pm</math>0.68</b>	<b>73.49<math>\pm</math>1.00</b>	<b>83.70<math>\pm</math>0.86</b>
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (1)	55.27 $\pm$ 0.74	51.39 $\pm$ 1.03	60.18 $\pm$ 1.09	51.63 $\pm$ 0.81	46.92 $\pm$ 1.14	57.71 $\pm$ 1.09	<b>62.43<math>\pm</math>0.77</b>	<b>55.94<math>\pm</math>1.10</b>	<b>71.40<math>\pm</math>1.08</b>
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (3)	62.09 $\pm$ 0.72	59.28 $\pm$ 1.03	65.30 $\pm$ 0.99	61.38 $\pm$ 0.74	59.29 $\pm$ 1.01	63.69 $\pm$ 1.09	<b>74.82<math>\pm</math>0.66</b>	<b>70.87<math>\pm</math>0.91</b>	<b>79.50<math>\pm</math>0.96</b>
ImageNet (ResNet-50) $\Rightarrow$ CIFAR (ResNet-101) (5)	67.84 $\pm$ 0.64	66.06 $\pm$ 0.89	69.81 $\pm$ 0.92	68.38 $\pm$ 0.74	67.72 $\pm$ 1.01	69.06 $\pm$ 1.00	<b>82.65<math>\pm</math>0.60</b>	<b>79.06<math>\pm</math>0.89</b>	<b>86.67<math>\pm</math>0.79</b>



Table 15: Prediction performances of RL-OSHeDA as well as baselines for all domain adaptation tasks in Multilingual Reuters Collection dataset.

Multilingual Reuters Collection									
	DS3L			KPG			OPDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
English ⇒ Spanish	59.35±0.96	52.92±1.31	67.57±1.20	11.48±0.08	8.80±0.00	17.12±0.95	56.42±1.02	49.09±1.32	66.41±1.23
French ⇒ Spanish	59.35±0.95	52.92±1.25	67.57±1.26	11.32±0.08	8.65±0.00	17.02±0.96	56.37±0.94	49.14±1.24	66.14±1.22
German ⇒ Spanish	59.35±0.92	52.92±1.24	67.57±1.27	11.11±0.07	8.43±0.00	17.00±0.97	55.32±0.85	48.16±1.05	65.07±1.21
Italian ⇒ Spanish	59.35±0.94	52.92±1.30	67.57±1.22	11.17±0.07	8.48±0.00	17.02±0.97	55.27±1.02	47.51±1.30	66.15±1.24
	PL			SCT			SSAN		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
English ⇒ Spanish	42.85±0.81	34.56±1.02	57.86±1.29	60.99±0.92	54.70±1.25	68.96±1.24	58.28±0.94	51.99±1.21	66.36±1.24
French ⇒ Spanish	42.85±0.81	34.56±1.01	57.86±1.29	61.03±0.94	<b>54.94±1.32</b>	68.67±1.20	57.59±0.91	51.04±1.24	66.07±1.21
German ⇒ Spanish	42.85±0.78	34.56±0.98	57.86±1.28	61.58±0.96	<b>55.49±1.33</b>	69.21±1.22	58.65±0.97	52.77±1.29	66.04±1.29
Italian ⇒ Spanish	42.85±0.82	34.56±0.98	57.86±1.30	61.07±0.95	54.69±1.31	69.17±1.18	58.47±0.92	52.15±1.25	66.58±1.20
	STN			SL			OSHeDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
English ⇒ Spanish	59.56±0.97	53.26±1.33	67.59±1.20	58.53±0.96	52.14±1.32	66.74±1.19	<b>65.85±0.90</b>	<b>55.09±1.25</b>	<b>81.97±1.07</b>
French ⇒ Spanish	59.19±0.98	52.95±1.29	67.13±1.27	58.53±0.95	52.14±1.28	66.74±1.22	<b>65.34±0.97</b>	54.52±1.21	<b>81.64±0.85</b>
German ⇒ Spanish	59.12±0.96	52.70±1.28	67.34±1.28	58.53±0.98	52.14±1.38	66.74±1.19	<b>65.54±0.88</b>	54.63±1.17	<b>82.22±0.94</b>
Italian ⇒ Spanish	58.95±0.95	52.72±1.32	66.89±1.19	58.53±0.96	52.14±1.31	66.74±1.25	<b>64.81±0.91</b>	<b>53.63±1.19</b>	<b>82.07±0.98</b>

Table 16: Prediction performances of RL-OSHeDA as well as baselines for all domain adaptation tasks in NUSWIDE & ImageNet dataset.

NUSWIDE & ImageNet									
	DS3L			KPG			OPDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
Text ⇒ Image	67.61±1.65	66.17±2.22	69.20±2.30	55.18±1.17	52.60±0.00	58.10±2.45	71.06±1.44	66.60±1.98	76.38±2.09
	PL			SCT			SSAN		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
Text ⇒ Image	42.43±0.26	34.05±0.00	61.15±2.10	70.42±1.49	68.00±2.20	73.10±1.99	67.98±1.49	66.25±2.04	69.85±2.21
	STN			SL			OSHeDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
Text ⇒ Image	67.75±1.23	64.80±1.42	71.08±2.16	69.41±1.64	66.62±2.26	72.57±2.23	<b>80.01±1.30</b>	<b>74.65±2.01</b>	<b>86.35±0.81</b>

Table 17: Prediction performances of RL-OSHeDA as well as baselines for all domain adaptation tasks in Wikipedia dataset.

Wikipedia									
	DS3L			KPG			OPDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
Image ⇒ Text	35.90±2.42	26.96±2.70	54.68±3.09	25.37±0.40	16.80±0.00	51.92±3.19	39.60±2.46	29.48±2.60	<b>61.16±3.26</b>
Text ⇒ Image	76.10±1.61	74.48±2.02	77.80±2.51	24.27±0.44	16.00±0.00	53.12±3.16	65.72±1.38	62.40±1.03	69.68±2.98
	PL			SCT			SSAN		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
Image ⇒ Text	27.11±2.26	18.28±2.10	54.00±3.18	38.51±2.45	29.28±2.77	56.56±3.00	39.52±2.03	30.16±2.18	57.52±3.10
Text ⇒ Image	56.64±1.04	52.00±0.85	62.80±3.02	78.31±1.58	76.44±2.00	80.28±2.47	77.23±1.48	75.36±1.73	79.20±2.50
	STN			SL			OSHeDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
Image ⇒ Text	39.96±2.11	30.24±2.21	59.20±3.16	37.80±2.48	28.68±2.70	55.68±3.07	<b>40.27±2.40</b>	<b>31.60±2.79</b>	56.68±3.22
Text ⇒ Image	75.54±1.70	72.56±2.05	78.80±2.79	76.40±1.47	74.52±1.68	78.40±2.46	<b>85.93±1.39</b>	<b>82.92±2.11</b>	<b>89.40±1.51</b>

Table 18: Prediction performances of RL-OSHeDA as well as baselines for all domain adaptation tasks in PTB-XL dataset.

PTB-XL									
	DS3L			KPG			OPDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
Digital ECG ⇒ Paper ECG	30.30±1.19	34.95±0.58	26.74±1.82	N/A	N/A	N/A	31.47±1.22	36.35±0.63	27.74±1.86
	PL			SCT			SSAN		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
Digital ECG ⇒ Paper ECG	26.18±1.37	36.43±0.55	20.43±1.66	26.23±1.65	<b>46.48±0.71</b>	18.27±1.60	25.16±1.47	40.40±0.65	18.27±1.54
	STN			SL			OSHeDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
Digital ECG ⇒ Paper ECG	20.64±0.96	22.23±0.26	19.27±1.65	25.74±1.55	44.50±0.76	18.11±1.52	<b>47.48±1.25</b>	44.30±1.39	<b>51.16±1.86</b>



Table 20: Prediction performances of RL-OSHeDA as well as baselines for all domain adaptation tasks in Office & Caltech256 dataset.

	Office & Caltech256								
	DS3L			KPG			OPDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
Amazon (DeCAF6) ⇒ Caltech (SURF)	46.08±1.41	37.41±1.60	60.15±2.05	36.51±0.69	30.20±0.00	46.22±2.22	43.26±1.58	33.61±1.87	60.84±2.01
Amazon (DeCAF6) ⇒ DSLR (SURF)	70.73±4.37	61.96±6.58	82.57±4.50	55.39±1.70	46.05±0.00	69.59±5.04	61.74±4.41	50.73±4.98	79.32±4.97
Amazon (DeCAF6) ⇒ Webcam (SURF)	71.64±2.73	67.43±4.05	76.62±3.64	61.25±1.62	56.91±0.00	66.34±3.81	61.79±2.88	53.28±3.61	74.00±4.06
Amazon (SURF) ⇒ Caltech (DeCAF6)	68.29±1.31	65.51±1.83	71.37±1.98	36.27±0.91	34.10±0.00	38.80±2.06	58.80±1.40	52.44±2.00	67.48±1.85
Amazon (SURF) ⇒ DSLR (DeCAF6)	86.97±3.62	83.57±5.92	90.68±3.41	24.18±0.94	16.86±0.00	44.32±5.69	81.76±3.69	75.07±5.53	90.00±3.28
Amazon (SURF) ⇒ Webcam (DeCAF6)	83.31±1.92	81.56±2.76	85.17±2.58	29.59±1.34	24.79±0.00	36.97±3.90	76.18±2.55	70.34±3.91	83.59±2.91
Caltech (DeCAF6) ⇒ Amazon (SURF)	59.62±1.55	54.12±2.14	66.47±2.16	42.35±0.87	36.87±0.00	49.89±2.27	53.67±1.48	46.16±1.95	64.31±2.24
Caltech (DeCAF6) ⇒ DSLR (SURF)	70.73±4.53	61.96±6.67	82.57±4.46	37.67±1.36	28.80±0.00	54.86±5.40	61.74±4.48	50.73±5.07	79.32±5.22
Caltech (DeCAF6) ⇒ Webcam (SURF)	71.64±2.92	67.43±4.31	76.62±3.72	44.51±1.41	39.21±0.00	51.52±4.14	61.79±2.88	53.28±3.61	74.00±4.20
Caltech (SURF) ⇒ Amazon (DeCAF6)	83.06±1.13	82.42±1.58	83.72±1.59	11.18±0.34	8.48±0.00	16.55±1.69	78.11±1.16	74.73±1.72	81.85±1.58
Caltech (SURF) ⇒ DSLR (DeCAF6)	86.97±3.55	83.57±5.96	90.68±3.48	3.93±0.02	2.20±0.00	30.95±5.32	81.76±3.67	75.07±5.49	90.00±3.16
Caltech (SURF) ⇒ Webcam (DeCAF6)	83.31±1.96	81.56±2.87	85.17±2.68	8.19±0.50	5.21±0.00	22.34±3.45	76.18±2.50	70.34±3.83	83.59±2.93
Webcam (DeCAF6) ⇒ Amazon (SURF)	59.62±1.54	54.12±2.11	66.47±2.23	55.55±0.93	50.11±0.00	62.37±2.24	53.67±1.54	46.16±2.03	64.31±2.23
Webcam (DeCAF6) ⇒ Caltech (SURF)	46.08±1.44	37.41±1.61	60.15±2.09	37.74±0.67	31.18±0.00	47.86±2.25	43.26±1.53	33.61±1.80	60.84±2.14
Webcam (DeCAF6) ⇒ DSLR (SURF)	70.73±4.54	61.96±6.79	82.57±4.39	57.27±1.65	47.80±0.00	71.49±4.83	61.74±4.44	50.73±5.08	79.32±5.20
Webcam (SURF) ⇒ Amazon (DeCAF6)	83.06±1.13	82.42±1.60	83.72±1.64	24.29±0.95	21.34±0.00	28.28±2.22	78.11±1.17	74.73±1.76	81.85±1.57
Webcam (SURF) ⇒ Caltech (DeCAF6)	68.29±1.37	65.51±1.86	71.37±2.03	33.67±0.94	31.11±0.00	36.74±2.18	58.80±1.47	52.44±2.08	67.48±1.87
Webcam (SURF) ⇒ DSLR (DeCAF6)	86.97±3.64	83.57±5.95	90.68±3.32	20.74±0.94	14.02±0.00	40.95±5.73	81.76±3.57	75.07±5.45	90.00±3.16
	PL			SCT			SSAN		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
	Amazon (DeCAF6) ⇒ Caltech (SURF)	28.23±1.18	19.50±1.18	53.47±2.21	<b>47.81±1.35</b>	38.66±1.55	<b>62.75±2.11</b>	43.85±1.73	35.11±1.91
Amazon (DeCAF6) ⇒ DSLR (SURF)	35.22±0.51	23.69±0.16	72.84±5.04	<b>72.68±4.36</b>	<b>64.10±6.65</b>	84.05±4.40	66.44±4.69	57.19±6.61	79.59±4.94
Amazon (DeCAF6) ⇒ Webcam (SURF)	32.23±2.14	22.87±2.21	58.62±4.18	75.46±2.54	71.00±3.65	80.55±3.28	70.98±2.78	65.72±4.19	77.31±3.47
Amazon (SURF) ⇒ Caltech (DeCAF6)	54.21±1.05	46.73±1.21	66.66±1.79	70.29±1.39	67.05±1.94	73.89±1.93	71.94±1.33	68.67±1.78	75.59±2.03
Amazon (SURF) ⇒ DSLR (DeCAF6)	67.80±1.78	58.16±0.00	82.70±4.70	93.35±2.26	91.40±3.62	<b>95.41±2.34</b>	90.18±2.45	87.74±3.70	92.84±2.94
Amazon (SURF) ⇒ Webcam (DeCAF6)	66.08±2.23	60.34±2.59	73.38±3.77	88.13±1.60	85.38±2.29	91.17±2.19	86.68±2.14	83.47±3.14	90.28±2.88
Caltech (DeCAF6) ⇒ Amazon (SURF)	36.43±1.27	27.67±1.35	55.15±2.45	63.17±1.51	57.01±2.17	70.90±2.14	59.90±1.49	55.04±2.02	65.82±2.17
Caltech (DeCAF6) ⇒ DSLR (SURF)	35.22±0.47	23.69±0.00	72.84±4.99	<b>72.83±4.44</b>	64.86±6.74	<b>83.11±4.35</b>	65.39±4.32	56.48±6.14	77.84±5.10
Caltech (DeCAF6) ⇒ Webcam (SURF)	32.23±2.06	22.87±2.15	58.62±4.00	75.66±2.51	71.82±3.76	80.00±3.60	70.76±2.56	65.35±3.72	77.31±3.34
Caltech (SURF) ⇒ Amazon (DeCAF6)	68.53±0.95	63.58±1.11	75.11±1.61	88.53±1.00	87.22±1.45	89.87±1.34	89.04±1.00	87.75±1.51	90.38±1.26
Caltech (SURF) ⇒ DSLR (DeCAF6)	67.80±1.81	58.16±0.32	82.70±4.74	93.50±2.70	91.93±4.64	<b>95.14±2.32</b>	91.36±2.59	88.96±4.08	93.92±2.95
Caltech (SURF) ⇒ Webcam (DeCAF6)	66.08±2.21	60.34±2.51	73.38±3.85	87.79±1.67	84.82±2.51	91.03±2.04	87.13±2.22	83.43±3.32	91.31±2.65
Webcam (DeCAF6) ⇒ Amazon (SURF)	36.43±1.24	27.67±1.32	55.15±2.21	62.57±1.50	56.43±2.06	70.27±2.06	60.23±1.45	55.25±2.03	66.28±2.20
Webcam (DeCAF6) ⇒ Caltech (SURF)	28.23±1.11	19.50±1.08	53.47±2.13	<b>47.38±1.45</b>	<b>38.50±1.70</b>	61.81±2.15	43.79±1.68	34.81±1.89	59.85±2.12
Webcam (DeCAF6) ⇒ DSLR (SURF)	35.22±0.55	23.69±0.22	72.84±5.26	71.98±3.87	63.60±5.95	83.11±3.91	65.60±4.92	56.35±6.71	78.65±5.29
Webcam (SURF) ⇒ Amazon (DeCAF6)	68.53±0.92	63.58±1.04	75.11±1.60	87.85±1.01	86.45±1.50	89.31±1.37	87.85±1.04	86.85±1.53	88.89±1.39
Webcam (SURF) ⇒ Caltech (DeCAF6)	54.21±1.05	46.73±1.23	66.66±1.82	70.55±1.30	66.91±1.90	74.68±1.85	72.99±1.33	69.26±1.82	77.21±1.87
Webcam (SURF) ⇒ DSLR (DeCAF6)	67.80±1.87	58.16±0.57	82.70±4.81	93.41±2.11	91.78±3.48	95.14±2.31	89.06±2.81	86.35±4.49	92.03±3.17
	STN			SL			OSHeDA		
	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>	<i>HOS</i>	<i>OS*</i>	<i>UNK</i>
	Amazon (DeCAF6) ⇒ Caltech (SURF)	47.74±1.61	<b>39.29±1.97</b>	60.99±2.23	45.87±1.48	37.12±1.82	60.25±2.22	46.69±1.39	37.78±1.65
Amazon (DeCAF6) ⇒ DSLR (SURF)	72.52±4.17	63.01±6.13	<b>85.54±3.83</b>	70.36±4.71	62.20±6.92	81.08±4.81	71.88±4.38	63.98±6.39	82.30±4.44
Amazon (DeCAF6) ⇒ Webcam (SURF)	<b>77.24±2.61</b>	<b>72.38±3.91</b>	82.90±3.20	71.76±2.81	67.75±4.18	76.34±3.68	75.72±2.68	68.42±3.90	<b>84.97±2.98</b>
Amazon (SURF) ⇒ Caltech (DeCAF6)	69.07±1.33	65.04±1.92	73.68±1.82	68.20±1.37	66.15±1.93	70.40±1.96	<b>79.28±1.13</b>	<b>72.99±1.74</b>	<b>86.85±0.99</b>
Amazon (SURF) ⇒ DSLR (DeCAF6)	80.02±2.28	74.04±3.82	87.43±2.26	86.88±3.67	84.13±6.05	89.86±3.43	<b>94.59±2.13</b>	<b>95.35±2.31</b>	94.05±3.31
Amazon (SURF) ⇒ Webcam (DeCAF6)	81.48±2.18	77.27±2.87	86.28±3.32	83.33±2.01	80.97±2.76	85.93±2.74	<b>92.45±1.30</b>	<b>90.28±1.94</b>	<b>94.97±1.72</b>
Caltech (DeCAF6) ⇒ Amazon (SURF)	61.69±1.44	55.02±1.93	70.27±2.10	60.31±1.53	54.92±2.12	66.97±2.15	<b>66.11±1.53</b>	<b>57.86±2.16</b>	<b>77.33±1.80</b>
Caltech (DeCAF6) ⇒ DSLR (SURF)	69.10±4.13	59.61±6.27	82.57±4.10	70.36±4.51	62.20±6.67	81.08±4.82	72.05±4.18	<b>65.33±5.65</b>	81.35±5.62
Caltech (DeCAF6) ⇒ Webcam (SURF)	<b>77.24±2.58</b>	<b>72.52±3.98</b>	82.76±3.14	71.76±2.95	67.75±4.23	76.34±3.78	74.86±2.76	67.65±3.92	<b>84.14±3.39</b>
Caltech (SURF) ⇒ Amazon (DeCAF6)	85.54±1.12	83.07±1.61	88.26±1.46	83.90±1.09	83.09±1.57	84.73±1.51	<b>94.36±0.71</b>	<b>91.43±1.13</b>	<b>97.50±0.83</b>
Caltech (SURF) ⇒ DSLR (DeCAF6)	80.56±3.42	75.17±5.36	86.89±3.77	86.88±3.51	84.13±5.87	89.86±3.36	<b>94.32±2.02</b>	<b>96.30±2.63</b>	92.57±2.91
Caltech (SURF) ⇒ Webcam (DeCAF6)	78.34±2.00	74.47±2.27	82.76±3.37	83.33±1.99	80.97±2.79	85.93±2.74	<b>92.00±1.24</b>	<b>88.87±1.77</b>	<b>95.79±1.71</b>
Webcam (DeCAF6) ⇒ Amazon (SURF)	61.26±1.48	54.27±1.91	70.40±2.17	60.31±1.53	54.92±2.18	66.97±2.16	<b>66.00±1.51</b>	<b>57.20±2.12</b>	<b>78.11±1.84</b>
Webcam (DeCAF6) ⇒ Caltech (SURF)	46.97±1.59	38.19±1.94	61.11±2.17	45.87±1.53	37.12±1.87	60.25±2.24	46.72±1.41	37.53±1.73	<b>62.21±2.15</b>
Webcam (DeCAF6) ⇒ DSLR (SURF)	72.18±4.16	63.57±6.06	<b>83.65±4.24</b>	70.36±4.88	62.20±7.04	81.08±4.82	<b>71.91±4.48</b>	<b>63.68±6.80</b>	82.70±4.63
Webcam (SURF) ⇒ Amazon (DeCAF6)	86.87±1.07	84.97±1.55	88.87±1.43	83.90±1.11	83.09±1.58	84.73±1.50	<b>94.44±0.75</b>	<b>91.42±1.27</b>	<b>97.69±0.65</b>
Webcam (SURF) ⇒ Caltech (DeCAF6)	69.43±1.34	65.14±1.84	74.41±1.99	68.20±1.38	66.15±1.87	70.40±1.95	<b>79.49±1.18</b>	<b>73.01±1.81</b>	<b>87.37±1.01</b>
Webcam (SURF) ⇒ DSLR (DeCAF6)	83.46±2.50	79.20±3.97	88.38±2.92	86.88±3.65	84.13±5.95	89.86±3.49	<b>94.36±2.21</b>	<b>95.68±3.54</b>	<b>93.24±2.70</b>

Table 21: Pairwise p-values from the Nemenyi test conducted on 56 domain adaptation tasks. P-values less than 0.05 indicate a statistically significant difference in prediction performance between the two corresponding methods.

	DS3L	KPG	OPDA	PL	SCT	SSAN	STN	SL	RL-OSHeDA
DS3L	1.000	0.001	0.001	0.001	0.525	0.900	0.900	0.900	<b>0.001</b>
KPG	0.001	1.000	0.283	0.900	0.001	0.001	0.001	0.001	<b>0.001</b>
OPDA	0.001	0.283	1.000	0.050	0.001	0.001	0.001	0.001	<b>0.001</b>
PL	0.001	0.900	0.050	1.000	0.001	0.001	0.001	0.001	<b>0.001</b>
SCT	0.525	0.001	0.001	0.001	1.000	0.589	0.062	0.525	<b>0.009</b>
SSAN	0.900	0.001	0.001	0.001	0.589	1.000	0.900	0.900	<b>0.001</b>
STN	0.900	0.001	0.001	0.001	0.062	0.900	1.000	0.900	<b>0.001</b>
SL	0.900	0.001	0.001	0.001	0.525	0.900	0.900	1.000	<b>0.001</b>
RL-OSHeDA	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>	<b>0.009</b>	<b>0.001</b>	<b>0.001</b>	<b>0.001</b>	<b>1.000</b>