

Learning Invariant Representations across Domains and Tasks

Jindong Wang^{1*}, Wenjie Feng², Chang Liu¹, Chaohui Yu³, Mingxuan Du⁴,
Renjun Xu⁵, Tao Qin¹, Tie-Yan Liu¹

¹ Microsoft Research, Beijing, China ² Institute of Data Science, NUS, Singapore

³ Alibaba DAMO Academy ⁴ Hefei University of Technology ⁵ Zhejiang University

{jindong.wang, changliu}@microsoft.com, wenjie.feng@nus.edu.sg, huakun.ych@alibaba-inc.com

Abstract

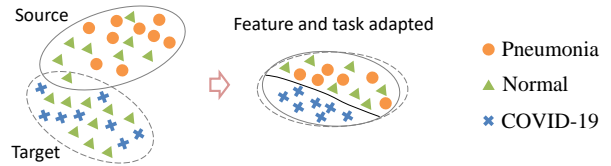
Being expensive and time-consuming to collect massive COVID-19 image samples to train deep classification models, transfer learning is a promising approach by transferring knowledge from the abundant typical pneumonia datasets for COVID-19 image classification. However, negative transfer may deteriorate the performance due to the feature distribution divergence between two datasets and task semantic difference in diagnosing pneumonia and COVID-19 that rely on different characteristics. It is even more challenging when the target dataset has no labels available, i.e., unsupervised task transfer learning.

In this paper, we propose a novel **Task Adaptation Network (TAN)** to solve this unsupervised task transfer problem. In addition to learning transferable features via domain-adversarial training, we propose a novel task semantic adaptor that uses the learning-to-learn strategy to adapt the task semantics. Experiments on three public COVID-19 datasets demonstrate that our proposed method achieves superior performance. Especially on COVID-DA dataset, TAN significantly increases the recall and F1 score by 5.0% and 7.8% compared to recently strong baselines. Moreover, we show that TAN also achieves superior performance on several public domain adaptation benchmarks.

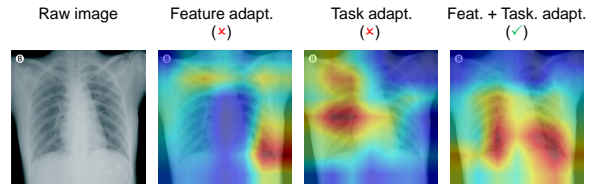
1. Introduction

The COVID-19 pandemic is greatly threatening global public health. In the battle against COVID-19, one critical challenge is to diagnose patients among a large number of people and provide necessary medical treatment so as to prevent further spread of the virus. Nowadays there is a growing trend to use the screening of chest radiography images (CRIs) such as X-ray images [42] for automated computer-aid diagnosis.

The diagnosis of COVID-19 based on chest X-ray images is a standard image classification problem with two classes: the infected and disinfected ones. While deep neu-



(a) Our method performs transfer learning from typical pneumonia to COVID-19 via feature distribution and task semantic adaptation.



(b) Our method gives correct predictions by adapting both feature distributions and task semantics. Attention map shows that our method can capture the critical factors [42] in the image that help detection.

Figure 1. Main idea and results visualization of our proposed TAN method. By adapting both the feature distribution and task semantics, TAN can eventually achieve accurate predictions by finding the most important critical factors.

ral networks (DNNs) have achieved great success for image classification, they often require a large amount of labeled images for training. Unfortunately, for COVID-19, large-scale annotations are costly and time-consuming to collect. Therefore, a straightforward approach is to leverage transfer learning (TL) techniques [23, 44, 25, 45] to transfer knowledge from existing (abundant) typical pneumonia datasets (i.e., the *source* domain) to COVID-19 (i.e., the *target* domain) to facilitate the model learning. In this paper, we mainly focus on the most challenging transfer setting where (1) the target domain has no labels and (2) the labels in the source and target domains are of different semantic meanings. We call this setting *unsupervised task transfer*.

In this unsupervised setting, the standard pretrain-finetune transfer paradigm becomes inapplicable, as there are no labeled images available in the target domain for finetuning. This requires us to conduct unsupervised adap-

tation between two different tasks, i.e., train a model on the labeled source domain and adapt it to the target domain in an unsupervised manner. Unsupervised adaptation presents two critical challenges that can result in *negative transfer* [25] and produces even worse performance than no transfer. The first challenge is the *feature distribution divergence*, which naturally exists since the distribution of these features differs from the source to the target domain. Hence, feature distribution adaptation is necessary. The second challenge is the *task semantic difference* since diagnosing pneumonia and COVID-19 are two related but different tasks that have different preferences on the critical factors [42]. Thus, the task semantics should also be adapted to maximize the transfer performance. While existing domain adaptation (DA) methods are able to adapt feature distributions when the source and target tasks are identical (e.g., both of the source and target domains are classifying monitors under different background), they are not applicable to our problem [7, 50, 18, 4, 37].

In this paper, we propose a **Task Adaptation Network (TAN)** for this unsupervised task transfer problem. The concept of our method is illustrated in Figure 1(a). TAN is able to learn transferable features across domains and tasks. Concretely speaking, TAN firstly adopts domain-adversarial training to reduce the feature distribution divergence between domains. However, only adapting features is not sufficient due to the task semantic difference. TAN devises a novel permutation-invariant task semantic adaptor that uses the learning-to-learn strategy to handle task semantic difference. We design a feature-critic training algorithm that effectively adapts task semantics using the pivot data. Figure 1(b) gives an example of the activation map of TAN to show its effectiveness in finding the critical factors [42] by adapting both feature distribution and task semantics.

To sum up, this paper makes the following contributions:

1. We propose a novel Task Adaptation Network (TAN) for unsupervised task transfer that addresses the feature distribution divergence and task semantic difference. Especially for the challenging task semantic adaptation, we propose a novel task semantic adaptor that leverages the learning-to-learning strategy to adapt cross-domain tasks.
2. Experiments on three public COVID-19 chest X-ray image classification datasets demonstrate that TAN outperforms several state-of-the-art baselines. To be more specific, on the challenging COVID-DA dataset, TAN significantly improves the F1 score and recall by 7.8% and 5.0% respectively, compared to the second best baseline.
3. Moreover, TAN is a general and flexible method that also achieves superior performance on sev-

eral public domain adaptation benchmarks including ImageCLEF-DA, Office-Home, and VisDA-2017.

2. Related Work

Transfer learning (TL) [25] is a useful technology to transfer the knowledge from existing source domains to the target domain, especially when the target domain has sparse or no labels. Such label scarcity problem can be solved using TL by firstly pretraining on a large dataset such as ImageNet [8] and then finetune the pretrained model on downstream tasks. This strategy is widely used in modern computer vision research [23, 36, 45, 44, 9]. In a semi-supervised setting where the target domain has labels, Luo *et al.* [22] proposed a domain and task transfer network to handle the different tasks using task semantic transfer. However, when the target domain has no labels, the pretrain-finetune paradigm is not available.

When two tasks are related, multi-task learning (MTL) [5, 17, 33] can be used to learn transferable features to enhance their learning performance. MTL also works when different tasks have labels available. Meta-learning, or learning-to-learn [2, 10, 31] aims to learn general knowledge from a bunch of tasks and then transfer to unseen tasks. Meta-learning often works under the few-shot setting where different tasks have several labels available and does not explicitly reduce the feature distribution divergence between domains and tasks. Zero-shot learning (ZSL) [24, 34] focuses on classifying all unseen classes which can be seen as a general case to our problem that has one unseen category. In contrast, ZSL does not reduce the distribution divergence across domains.

Domain adaptation (DA) is a specific area of transfer learning [25]. DA aims at building cross-domain models by reducing the distribution divergence of a representation via some divergence measured by such as Maximum Mean Discrepancy (MMD) [14], KL or JS divergence, cosine similarity, and higher-order moments [46, 18, 38, 35, 49, 41]. Another line of work relies on the generative adversarial nets [13] to learn domain-invariant features [50, 29, 37, 11]. While great progress has been made, directly applying DA to our problem is not sufficient since DA generally works when two domains have identical categories. The setting where two domains have different but overlapped tasks is also explored in recent open set DA [26, 30] and partial DA [47] methods. However, their purpose is to recognize the overlapped (common) categories rather than the unshared classes in the target domain.

3. Methodology

We introduce an unsupervised learning model which transfers information from a large labeled source domain, S , to a target domain, T , across different tasks. The goal be-

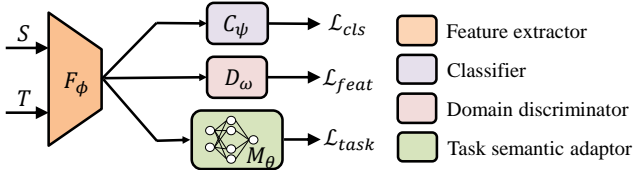


Figure 2. The architecture of the proposed TAN method that consists of four modules: feature extractor F_ϕ , classifier C_ψ , domain discriminator D_ω for feature distribution adaptation, and task semantic adaptor M_θ for task semantic adaptation.

ing to learn a strong transferable target classifier $h : \mathcal{X}^T \rightarrow \mathcal{Y}^T$ that reduces both the feature distribution divergence and task semantic difference.

We assume the source domain contains n^s images, $\mathbf{x}^s \in \mathcal{X}^S$, with associated labels, $\mathbf{y}^s \in \mathcal{Y}^S$. The target domain consists of n^t unlabeled images, $\mathbf{x}^t \in \mathcal{X}^T$. \mathcal{S} and \mathcal{T} share the same feature space, i.e. $\mathcal{X}^S = \mathcal{X}^T$.

Unlike traditional domain adaptation approaches that assume a cross-domain distribution shift under a shared label space ($\mathcal{Y}^S = \mathcal{Y}^T$), we aim to adapt both the feature distribution and task semantics, i.e., we consider the case where the tasks corresponding to source and target spaces are only *similar but not identical*, $\mathcal{Y}^S \neq \mathcal{Y}^T$, like for the pneumonia and COVID-19. Even if these two classification problems can be treated as classifying images into $\{0, 1\}$, their semantics are still different.¹

3.1. Overview

In this paper, we propose a novel Task Adaptation Network (TAN) to adapt the feature distribution and task semantics. We depict our overall model in Figure 2. TAN consists of four modules: feature extractor F_ϕ , classifier C_ψ , domain classifier D_ω , and task semantic adaptor M_θ with ϕ , ψ , ω , and θ as the learnable parameters. Taking as inputs the labeled source examples, TAN learns a latent feature space with F_ϕ and the binary classification network C_ψ by standard supervised learning. Then, by also taking as inputs the unlabeled target domain, TAN adapts the feature distributions and task semantics with the domain classifier D_ω and task semantic adaptor M_θ , respectively.

Our model jointly optimizes over a source classification loss \mathcal{L}_{cls} , a feature distribution adaptation loss \mathcal{L}_{feat} , and a task semantic transfer objective \mathcal{L}_{task} . Thus, the total objective of TAN can be formulated as follows:

$$\mathcal{L} = \mathcal{L}_{cls} + \lambda \mathcal{L}_{feat} + \mu \mathcal{L}_{task}, \quad (1)$$

where the hyperparameters λ and μ determine the influence of the feature distribution adaptation and task semantic adaptation, respectively. We use the cross-entropy loss

¹Although we can formulate both binary classification problems using $\mathcal{Y} = \{0, 1\}$ for each, the semantic of label “1”, i.e. the “true label”, differs.

$\ell^{(CE)}$ to measure the classification error \mathcal{L}_{cls} on the labeled source domain:

$$\mathcal{L}_{cls} = \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{S}} \ell^{(CE)}(C_\psi(F_\phi(\mathbf{x})), y), \quad (2)$$

where \mathbb{E} denotes the expectation. In the following sections, we will elaborate on the feature distribution adaptation and task semantic transfer modules.

3.2. Feature distribution adaptation

Considering the similarity between different tasks in \mathcal{S} and \mathcal{T} , the feature distribution adaptation module aims to reduce the cross-domain feature distribution divergence. Inspired by the well-established domain-adversarial training [11, 12], TAN designs a domain discriminator D_ω to adapt the cross-domain feature distributions. Domain-adversarial training is a two-player game where the domain discriminator D_ω is trained to distinguish the source domain from the target domain, while the feature extractor F_ϕ tries to confuse the domain discriminator by learning domain-invariant features. These two players are trained adversarially, i.e., ω is trained by maximizing the loss on domain discriminator D_ω while ϕ, ω are trained to minimize the loss of the feature extractor which is computed by the domain classifier. This optimization procedure eventually minimizes the difference between the feature distributions on the two domains, measured by the Jensen-Shannon divergence [13]. The adversarial training loss for feature distribution adaptation can be formulated as:

$$\begin{aligned} \mathcal{L}_{feat} = & \mathbb{E}_{\mathbf{x}^s \in \mathcal{S}} \log[D_\omega(F_\phi(\mathbf{x}^s))] \\ & + \mathbb{E}_{\mathbf{x}^t \in \mathcal{T}} \log[1 - D_\omega(F_\phi(\mathbf{x}^t))]. \end{aligned} \quad (3)$$

3.3. Task semantic adaptation

Domain-adversarial training is insufficient for our problem since it learns domain-invariant features across domains regardless of the tasks. In our problem, even if the source and target tasks are similar which makes domain-adversarial training reasonable, these tasks are still not identical. Adapting feature distributions unnecessarily ensures the adaptation across different task semantics. Since features are highly correlated with the tasks, the difference in task semantics would harm the feature adaptation and the performance of knowledge transfer.

On this problem, we propose a task semantic adaptor M_θ that employs a learning-to-learn strategy to adapt task semantics by learning from the domain-adversarial features. Learning-to-learn, or meta-learning [2, 10, 31] aims to effectively leverage the datasets and prior knowledge of a task ensemble in order to rapidly learn new tasks often with a small amount of data. Therefore, since the task semantics are difficult to model, we turn to using the learning-to-learn strategy. The key idea is to let M_θ learn the adaptation ability from domain-adversarial training and then such

ability can be utilized for task adaptation. Therefore, when domain-adversarial training gradually encourages the features to be domain-invariant, the task semantic adaptor M_θ can gradually learn this ability using the learning-to-learn strategy and eventually also enforces the features to be task-invariant. Hence, the task semantics can be adapted.

Technically speaking, the task semantic adaptor M_θ is implemented as an MLP network that can theoretically approximate any continuous functions [6] to enable the superior adaptation. Denoting $\mathbf{F}_\phi^s, \mathbf{F}_\phi^t$ as the source- and target-domain features extracted by F_ϕ , the task semantic adaptation loss \mathcal{L}_{task} can be formulated as:

$$\mathcal{L}_{task} = M_\theta(\mathbf{F}_\phi^s, \mathbf{F}_\phi^t). \quad (4)$$

Unfortunately, it remains challenging to optimize the above equation w.r.t. θ for three reasons. Firstly, what property should M_θ satisfy to ensure that the task semantics can be adapted? Secondly, how to maximally utilize the domain-invariant representations learned by adversarial training to better couple with the feature extractor for more effective training? Thirdly, how to update the adaptor network parameters in training?

Permutation invariance. The task semantic adaptation is supposed to reduce the difference between two tasks, so it should be invariant to permutations of samples that represent each task distribution. Therefore, we design M_θ to be *permutation-invariant* to the rows of its inputs, i.e., it should make no difference for different sample indices like the cases [1, 2, 3] and [3, 2, 1]. To enforce this property, we let M_θ take as inputs the pairwise distance between each element of \mathbf{F}_ϕ^s and \mathbf{F}_ϕ^t , which is permutation-invariant [16]. Then, the task semantic loss can be represented as:

$$M_\theta(\mathbf{F}_\phi^s, \mathbf{F}_\phi^t) = \text{MLP}(\text{Flatten}(\text{Gram}(\mathbf{F}_\phi^s, \mathbf{F}_\phi^t))), \quad (5)$$

where *Gram* denotes the Gram matrix computed by pairwise distance, *Flatten* is a flatten operation, and MLP denotes a multi-layer perceptron.

Pivot data. In a supervised case where target domain has labels, the task adaptor M_θ can be learned easily. However, it becomes challenging in this unsupervised setting. To maximally learn the adaptation ability from domain-adversarial training, M_θ is updated on a *pivot data* \mathcal{P} . It is a selected subset of both the source and target domains to maximally utilize the domain-adversarial training. The pivot data are the data that with high confidence scores during the learning process, so they can be representatives of the domain-adversarial training. As features are getting more domain-invariant, the classification performance on the target domain is gradually better, i.e., the pseudo labels for the target domain is getting more confident. This

pseudo-label training is widely adopted by most transfer learning literature [50, 37]. Therefore, based on the assumption that the adaptation ability can be learned from the domain-adversarial training, the task semantic adaptation could be better learned if M_θ can directly learn from samples with the most confident pseudo labels, i.e., the pivot data. The number of pivot data is important to our problem: less pivot data will bring more confidence and less generalization while more pivot data will do the opposite, which is empirically evaluated in later experiments.

To be more specific, the pivot data can be represented as:

$$\begin{aligned} \mathcal{P} &= \{\mathcal{P}_S^{(c)}, \mathcal{P}_T^{(c)}\}_{c \in \mathcal{Y}}, \\ \mathcal{P}_S^{(c)} &= \{(\mathbf{x}_j, y_j = c)\}_{j=1}^m, \mathcal{P}_T^{(c)} = \{(\mathbf{x}_j, \hat{y}_j = c)\}_{j=1}^m, \end{aligned} \quad (6)$$

where the data pairs are sorted as $\{(\mathbf{x}_j, \hat{y}_j)\}$ in decreasing order of prediction score. \hat{y} is the predicted (pseudo) label on the target domain and c is the class index. We select the top m instances for each class with high prediction scores (softmax probability that does not need the target labels). This selection is iterated in the whole learning process. For the source domain, we directly use its ground-truth labels. In total, we select $m \cdot |\mathcal{Y}|$ pivot data for each domain.

Feature-critic training. Different from classification, there is no supervision information for the target domain that makes it hard to update M_θ . In this paper, we propose a *feature-critic training* strategy [16] to update the task adaptor. For notation brevity, we pack $\Phi = \{\phi, \psi, \omega\}$ for parameters other than θ .

Our feature-critic training is introduced as follows. Let $\Phi(t)$ and $\Phi(t+1)$ denote parameter values in two consecutive learning steps t and $t+1$, respectively. Our key assumption is that as the pseudo labels on pivot data are getting more confident, if $\Phi(t+1)$ is better than $\Phi(t)$ for task adaptation, it should produce lower risks and better classification performance. Therefore, a reasonable feature-critic metric M_θ should evaluate a lower value for $\Phi(t+1)$ than for $\Phi(t)$. We thus update the feature-critic metric M_θ by minimizing difference \mathcal{L}_{val} computed by $\Phi(t)$ and $\Phi(t+1)$:

$$\mathcal{L}_{val} = \mathbb{E}_{\mathbf{x} \sim \mathcal{P}} \sigma(M_\theta(\mathbf{x}; \Phi(t+1)) - M_\theta(\mathbf{x}; \Phi(t))), \quad (7)$$

where $\sigma(\cdot)$ is an activation function.

3.4. Training and inference

The training process consists of two steps: 1) update Φ for the feature extractor, classification layer and domain discriminator and 2) update θ for the task semantic adaptor.

Update Φ . This step is to update Φ for classification and domain-adversarial training. To enforce the update of θ in the next step, we construct an assist model which is a copy

Algorithm 1 Learning algorithm of TAN

Input: Source domain $\mathcal{S} = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^{n_s}$, target domain $\mathcal{T} = \{\mathbf{x}_j^t\}_{j=1}^{n_t}$, learning rate α, β .

Output: $\{\Phi^*\}$.

- 1: Initialize $\Phi(0)$ and $\theta(0)$.
 - 2: **while** not done **do**
 - 3: Build an assist model with its parameter inherited from the main model.
 - 4: **For** mini-batch data B_s, B_t in \mathcal{S}, \mathcal{T} **do**
 - 5: Update Φ by Eq. (8).
 - 6: **End For**
 - 7: Select the data with the highest prediction confidence from \mathcal{T} to construct pivot data \mathcal{P} .
 - 8: Update θ by Eq. (9).
 - 9: **end while**
 - 10: **return** $\{\Phi^*\}$
-

of the main model by inheriting the same architecture and parameters (F_ϕ, C_ψ, D_ω) and use it for calculating the loss by Eq. (1) on all the labeled source domain and unlabeled target domain data. Note that for updating the domain discriminator D_ω , we do not use the mini-max optimization process and instead follow [11] to use the Gradient Reversal Layer for computing efficiency. Therefore, ω can be updated together with ϕ and ψ in a single back-propagation.

After getting the training loss, denote α the learning rate of the main model, then Φ can be updated by:

$$\Phi(t+1) = \Phi(t) - \alpha \nabla_{\Phi} (\mathcal{L}_{cls} + \lambda \mathcal{L}_{feat} + \mu \mathcal{L}_{task})|_{\Phi(t)}. \quad (8)$$

Update θ . This step is to update θ for the task adaptor M_θ using feature-critic training on the pivot data \mathcal{P} . Denote β its learning rate, then, θ can be updated by taking derivative of \mathcal{L}_{val} on θ :

$$\theta(t+1) = \theta(t) - \beta \nabla_{\theta} \mathcal{L}_{val}(\theta; \Phi(t), \Phi(t+1))|_{\theta(t)}, \quad (9)$$

where $\Phi(t)$ and $\Phi(t+1)$ are parameters of the assist and main model, respectively.

The above two steps are used iteratively since the pseudo labels of the pivot data can be more confident and all the losses can be minimized. In our experiments, we observe that the network will converge in dozens of epochs.

The training process of TAN is listed in Algorithm 1 and Figure 3. As for inference, we fix Φ to perform a single forward-pass to get the classification results for the test data.

4. Experiments

4.1. Datasets and setup

We evaluate TAN on public COVID-19 chest X-ray datasets. COVID-DA [51] contains three categories: typical

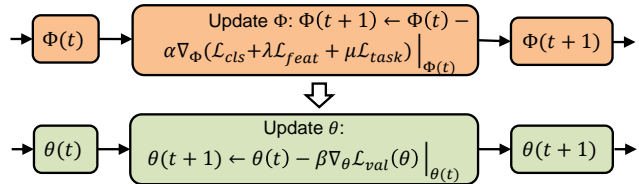


Figure 3. The learning process of TAN is composed of two iterative steps: update Φ and update θ .

Table 1. Statistical information of datasets

Dataset	#normal	#pneumonia	#COVID-19
COVID-DA	9,039	2,306	318
Bacterial	3,270	3,001	1,281
Viral	3,270	1,656	1,281

pneumonia, normal, and COVID-19. Curated-COVID [39] is a larger dataset that consists of four categories: normal, bacteria pneumonia, viral pneumonia, and COVID-19. Since our focus is on investigating the performance on transfer learning from pneumonia to COVID-19, we split Curated-COVID into two datasets where each one only contains one type of pneumonia and COVID-19. We name them *Bacterial* and *Viral* accordingly. Together we have three datasets as shown in Table 1.

For COVID-DA, we follow [51] to construct the source, target, and validation domains. For Bacterial and Viral datasets, we construct their source and target domains by taking all pneumonia (for source) and all COVID-19 (for target) samples accordingly. As for the normal category, we split them evenly into two domains. We further leave 20% of the target domain for validation. Eventually, there are two classes in the source domain: normal and pneumonia; and there are normal and COVID-19 classes in the target and validation datasets. In this case, the source domain does not contain any COVID-19 samples, which makes our problem harder than traditional transfer learning and domain adaptation. The detailed domain split information is presented in appendix.

4.2. Comparison methods

We compare the performance of TAN with three categories of methods: (1) deep and traditional transfer learning baselines, (2) deep diagnostic methods, and (3) unsupervised DA methods.

The deep and traditional TL baselines include: *Pretrain-only*, which trains a network on the source domain, and then directly apply the pretrained model on the target domain. *Target-train*, which is an ideal state and only for comparison since there are no labels for the target domain. We directly use several extra labeled COVID-19 data from the dataset (they are 30% of the target domain data) and train a network on these data. Then, we apply prediction on the target data. *Pretrain-finetune*, which is a standard TL paradigm

Table 2. Results on COVID-DA dataset (typical pneumonia → COVID-19, ResNet-18).

Method	P (%)	R (%)	F1 (%)
Pretrain-only	63.5	66.7	65.0
Target-train (<i>ideal</i>)	91.7	55.0	68.8
Pretrain-finetune (<i>ideal</i>)	56.3	75.0	64.3
DLAD [48]	62.0	73.3	67.2
DANN [11]	61.4	71.7	66.2
MCD [29]	63.2	60.0	61.5
CDAN+TransNorm [43]	85.0	39.2	63.7
MDD [50]	74.0	60.0	67.0
BNM [7]	43.0	75.0	56.0
TAN (Ours)	70.6	80.0	75.0

that finetunes the source pretrained model on the labeled target domain data. Note that *Target-train* and *Pretrain-finetune* methods use labeled data from the target domain, which is not available for TAN and other baselines. They are just for a comparison with the ideal setting of machine learning, i.e. the conventional supervised learning.

We choose DLAD [48] as the deep diagnostic method. The unsupervised DA methods are DANN [11], MCD [29], CDAN+TransNorm [43], MDD [50], and BNM [7]. All methods are using ResNet-18 [15] as the backbone network following [51]. The results of these methods are obtained from [51] to ensure a fair comparison. Note that we do not compare with [51] since it is a semi-supervised method that requires labeled data on the target domain.

For TAN, we use the mini-batch SGD with Nesterov momentum of 0.9 to optimize the main and the meta-network with batch size set as 16. The learning rate α of the main model changes by following [11]: $\alpha_k = \frac{\alpha}{(1+\gamma k)^{-\nu}}$, where k is the training iteration, $\gamma = 0.001$, $\alpha = 0.004$, and decay rate $\nu = 0.75$. The learning rate β for M_θ is set to be 0.0005. M_θ uses a in $-128 - 64 - 1$ MLP structure where in is the dimension of input matching features. We grid search the value of λ and μ in the range $[0.01, 0.05, 0.1, 0.5, 1, 5, 10]$ for the best performance.

During training, the labels of the target domain are not available and they are only be used for evaluation. For this binary classification problem, we use F1, Precision (P), and Recall (R) as the evaluation metrics. We do not use ROC/AUC since we are more interested in the recall and F1 in this specific disease diagnosis problem. The results are the average accuracy of ten trials.

4.3. Results and analysis

The results on COVID-DA dataset are presented in Table 2. Here we use the 95% confidence interval, where the corresponding value of z is 1.96. The computed confidence interval r is around 1.3%. Note that we do not list the accuracy results since all methods achieve similarly high

Table 3. Results on Bacterial dataset (bacterial pneumonia → COVID-19, ResNet-18)

Method	P (%)	R (%)	F1 (%)
Pretrain-only	92.2	91.6	91.9
DAN [18]	87.5	97.6	92.3
DANN [11]	93.1	96.8	94.9
MDD [50]	91.3	97.6	94.4
BNM [7]	90.6	98.9	94.6
TAN (Ours)	92.0	99.1	95.4

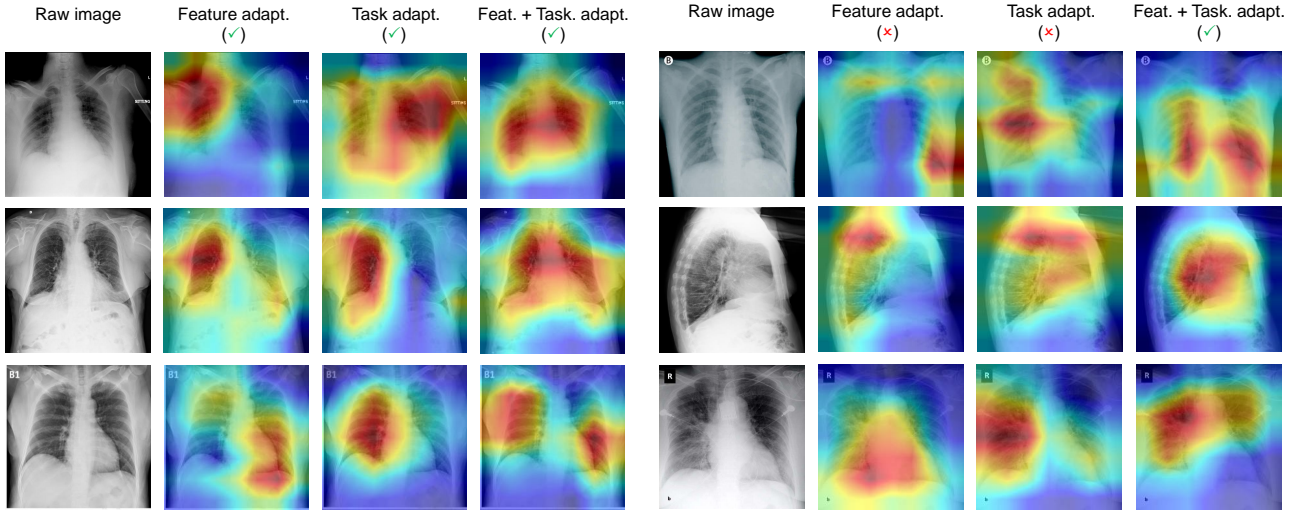
Table 4. Results on Viral dataset (viral pneumonia → COVID-19, ResNet-18)

Method	P (%)	R (%)	F1 (%)
Pretrain-only	85.4	65.9	74.4
DAN [18]	77.4	62.7	72.2
DANN [11]	67.4	83.3	74.5
MDD [50]	92.6	94.6	90.6
BNM [7]	88.7	97.7	93.0
TAN (Ours)	91.5	99.2	95.2

accuracy values in this binary classification problem. On COVID-DA dataset, our proposed TAN achieves a recall of 80.0% and F1 of 75.0%, which significantly outperforms the second best baselines by 5.0% in recall and 7.8% in F1 score. Pretrain-only and Finetune do not get good performance because of the feature distribution divergence and task semantic difference, which drastically limit their performance. Compared to other DA methods (DANN, MCD, CDAN, MDD and BNM), TAN also achieves better recall and F1 score. While the precision of CDAN is better than ours, its recall and F1 is not competent. Compared to the ideal state which is training on labeled target data, it is surprising to find that even working in fully unsupervised setting, our proposed TAN can achieve better recall and F1 score.

The results on Bacteria and Viral datasets are shown in Table 3 and Table 4, respectively. In these two datasets, we do not compare with the two ideal states target-train and pretrain-finetune since their performances are consistently better due to increased COVID-19 samples in these datasets. Our TAN still significantly outperforms other baselines in these large datasets. Same as on the COVID-DA dataset, pretrain-only gives the worst results due to feature distribution divergence and task semantic difference. While domain adaptation methods (DAN, DANN, MDD and BNM) outperform pretrain-only by aligning the feature distributions, they do not adapt the task semantics.

Comparing the results from all three tables, we see that with the numbers of unsupervised target domain samples increase, the transfer learning performance tends to become better, i.e., the results in Table 3 and Table 4 are generally better than the results in Table 2. This is because more representative knowledge can be learned when domains have



(a) All methods give correct predictions

(b) Adapting both feature and tasks give correct predictions

Figure 4. Results visualization using Grad-CAM [32] to show the attention weight of different adaptation modules. In some easy cases in (a), either adapting feature distributions or task semantics alone can achieve correct predictions. However for some hard cases in (b), it requires to adapt both feature distributions and task semantics to get correct predictions.

more samples that facilitates transfer learning. In all situations, TAN consistently achieves the best performance.

4.4. Ablation study

To further evaluate the effectiveness of TAN, we conduct an ablation study in Table 5. For task adaptation, we remove the domain-adversarial training module and directly let the model learn from the classification loss. For feature adaptation, we simply remove the task adaptation module. The results show the following observations. Firstly, the classification loss alone does not get good results, indicating the existence of feature distribution divergence and task semantic difference. Secondly, better performance can be achieved by combining the feature distribution adaptation and task semantic adaptation modules, indicating that both adaptation modules are effective. Thirdly, the best performance is achieved by combining both of the feature adaptation and task semantic adaptation modules, which proves that both of them are important in this problem.

Table 5. Ablation study on TAN

Variants	P (%)	R (%)	F1 (%)
\mathcal{L}_{cls}	63.5	66.7	65.0
$\mathcal{L}_{cls} + \mathcal{L}_{task}$	75.0	65.0	69.6
$\mathcal{L}_{cls} + \mathcal{L}_{feat}$	61.4	71.7	66.2
$\mathcal{L}_{cls} + \mathcal{L}_{feat} + \mathcal{L}_{task}$	70.6	80.0	75.0

4.5. Visualization study

The change of lung is critical factor for diagnosing COVID-19, which could be visualized to study the effectiveness of our method. Therefore, we visualize the atten-

tion maps for several COVID-19 images using the Gradient-weighted Class Activation Mapping (Grad-CAM) [32] in Figure 4. The shadow area in the figures is the lung area and the heat map denotes the activation weights for the model. Specifically, Figure 4(a) shows the cases where all adaptation modules give correct predictions, while Figure 4(b) shows the cases where wrong predictions are given by only feature distribution adaptation and only task semantic adaptation, and correct predictions are given by adapting both the feature distributions and task semantics.

From these results, we observe that in general cases where the samples are easy to classify, all adaptation methods give correct prediction. However, when the samples are hard such as the second one in Figure 4(b) which is sideways, it becomes harder for the model to classify. In this case, only adapting feature distributions or task semantics are not sufficient. Our proposed TAN can perform reasonably well in all situations.

4.6. Further analysis

Pivot data We empirically analyze size m of the pivot data \mathcal{P} . It is obvious that a larger m will bring more uncertainty, and a smaller m is likely to make the meta-network unstable. We record the performance of TAN using different values of m on COVID-DA dataset in Figure 5(a). The results indicate that TAN is robust to m and a small m can lead to competitive performance. Therefore, we set $m = 8$ in experiments for computational efficiency. We also compare different pivot data selection strategies in the appendix.

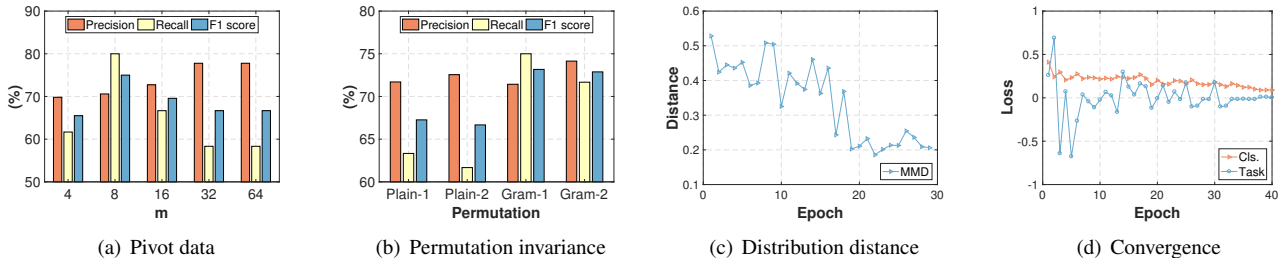


Figure 5. Detailed analysis of our method. (a) The number of pivot data in the target domain. (b) Evaluation of the permutation-invariant property of task semantic adaptor. (c) Distribution distance of feature adaptation. (d) Convergence analysis.

Permutation invariance We evaluate the permutation invariance property of M_θ . Figure 5(b) shows the two random results of using Gram matrix as inputs (i.e., permutation-invariant) and the raw features (i.e., the permutation-variant case, which we denote as ‘plain’ in the figure). While both networks outperform other baselines in Table 2, the Gram matrix gets the best performance, indicating that permutation invariance is important in the task semantic adaptor.

Feature distribution We compute the feature distribution distance using the maximum mean discrepancy (MMD) [3] as shown in Figure 5(c). The results indicate that our proposed TAN can gradually reduce the distribution distance.

Convergence and parameter sensitivity analysis We record the loss for classification and task semantic adaptation in Figure 5(d). The results show that although TAN involves both the classification and task adaptation networks, it can quickly reach a steady performance. This makes it easy to train in real applications. We also empirically analyze the sensitivity of the two trade-off parameters λ and μ and present the results in appendix, which shows that our method is relatively robust to these parameters.

5. Experiments on other datasets

Although our main focus is COVID-19, TAN is not limited to this problem. In fact, TAN can be applied to any dataset with similar setting and standard DA datasets.

Constructed dataset. We construct a new dataset from the VisDA-2017 DA challenge [28] by selecting its two random classes. In this dataset, the source domain is rendered 3D objects and the target domain is natural images. We call this constructed dataset *VisDA-binary*. Specifically, the source domain contains two classes: `train` (16,000) and `truck` (9,600) and the target domain contains `train` (4,236) and `bus` (4,690). The goal is to maximize the binary classification performance on the target domain, especially on class `bus`. This dataset is more balanced with more samples than COVID-19, which can be regarded as its

Table 6. Results on VisDA for binary classification (ResNet-50)

Method	P (%)	R (%)	F1 (%)
Pretrain-only	64.6	71.9	68.0
DAN [18]	78.4	54.1	64.0
DANN [11]	89.2	70.1	78.5
BNM [7]	79.6	60.5	68.8
TAN (Ours)	75.3	88.0	81.2

Table 7. Accuracy (%) on ImageCLEF-DA for UDA (ResNet-50).

Method	I→P	P→I	I→C	C→I	C→P	P→C	AVG
ResNet [15]	74.8	83.9	91.5	78.0	65.5	91.2	80.7
DAN [18]	75.0	86.2	93.3	84.1	69.8	91.3	83.3
DANN [11]	75.0	86.0	96.2	87.0	74.3	91.5	85.0
D-CORAL [35]	76.9	88.5	93.6	86.8	74.0	91.6	85.2
CAN [49]	78.2	87.5	94.2	89.5	75.8	89.2	85.7
JAN [20]	76.8	88.0	94.7	89.5	74.2	91.7	85.8
MADA [27]	75.0	87.9	96.0	88.8	75.2	92.2	85.8
CDAN [19]	77.7	90.7	97.7	91.3	74.2	94.3	87.7
TransNorm [43]	78.3	90.8	96.7	92.3	78.0	94.8	88.5
TAN (Ours)	78.7	91.0	97.0	92.0	79.7	96.0	89.1

complement. The results in Table 6 show that TAN outperforms all other comparison methods in recall and F1 score.

Standard DA benchmark. We also evaluate the performance of TAN on several standard domain adaptation benchmarks. Table 7 shows the results of TAN against other strong baselines on the ImageCLEF-DA [1] dataset. We see that although TAN is not specifically designed for traditional DA tasks, it still achieves competitive performance. We show the results on Office-Home [40] dataset in appendix where TAN also produces competitive performance.

6. Conclusions and Future Work

In this paper, we propose a Task Adaptation Network (TAN) for COVID-19 chest X-ray image classification by transferring knowledge from the typical pneumonia. TAN can adapt both of the cross-domain feature distributions and task semantics to produce accurate prediction on the target domain. Specifically for task semantic adaptation which is hard to model, we design a semantic adaptor that leverages the learning-to-learn strategy to learn the

adaptation ability from the domain-adversarial training. Experiments on several public datasets show that TAN significantly outperforms other comparison approaches. Moreover, TAN can also achieve competitive performance on several domain adaptation benchmarks.

In the future, we plan to apply TAN to more fine-grained COVID-19 diagnosis tasks such as detection and segmentation. TAN can also be applied to other COVID-19 data modalities like CT scans. In addition, we also plan to apply TAN to other similar transfer learning problems.

References

- [1] The imageclef-da challenge 2014. <https://www.imageclef.org/2014>. 8
- [2] Samy Bengio, Yoshua Bengio, Jocelyn Cloutier, and Jan Gecsei. On the optimization of a synaptic learning rule. In *Preprints Conf. Optimality in Artificial and Biological Neural Networks*, pages 6–8, 1992. 2, 3
- [3] Karsten M Borgwardt, Arthur Gretton, Malte J Rasch, Hans-Peter Kriegel, Bernhard Schölkopf, and Alex J Smola. Integrating structured biological data by kernel maximum mean discrepancy. *Bioinformatics*, 22(14):e49–e57, 2006. 8
- [4] Konstantinos Bousmalis, Nathan Silberman, David Dohan, Dumitru Erhan, and Dilip Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *CVPR*, pages 3722–3731, 2017. 2
- [5] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997. 2
- [6] Balázs Csanád Csáji. Approximation with artificial neural networks. *Faculty of Sciences, Eötvös Loránd University, Hungary*, 24:48, 2001. 4
- [7] Shuhao Cui, Shuhui Wang, Junbao Zhuo, Liang Li, Qingming Huang, and Qi Tian. Towards discriminability and diversity: Batch nuclear-norm maximization under label insufficient situations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3941–3950, 2020. 2, 6, 8
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 2
- [9] Jeff Donahue, Yangqing Jia, Oriol Vinyals, Judy Hoffman, Ning Zhang, Eric Tzeng, and Trevor Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In *International conference on machine learning*, pages 647–655, 2014. 2
- [10] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, pages 1126–1135, 2017. 2, 3
- [11] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*, 2015. 2, 3, 5, 6, 8, 12
- [12] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *JMLR*, 17(1):2096–2030, 2016. 3
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014. 2, 3
- [14] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *JMLR*, 13(Mar):723–773, 2012. 2
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 6, 8, 12
- [16] Yiyang Li, Yongxin Yang, Wei Zhou, and Timothy M Hospedales. Feature-critic networks for heterogeneous domain generalization. In *ICML*, 2019. 4
- [17] Shikun Liu, Edward Johns, and Andrew J Davison. End-to-end multi-task learning with attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1871–1880, 2019. 2
- [18] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I Jordan. Learning transferable features with deep adaptation networks. In *ICML*, 2015. 2, 6, 8, 12
- [19] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. Conditional adversarial domain adaptation. In *NeurIPS*, pages 1640–1650, 2018. 8
- [20] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Deep transfer learning with joint adaptation networks. In *ICML*, 2017. 8
- [21] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Deep transfer learning with joint adaptation networks. In *ICML*, pages 2208–2217, 2017. 12
- [22] Zelun Luo, Yuliang Zou, Judy Hoffman, and Li F Fei-Fei. Label efficient learning of transferable representations across domains and tasks. In *Advances in Neural Information Processing Systems*, pages 165–177, 2017. 2
- [23] Behnam Neyshabur, Hanie Sedghi, and Chiyuan Zhang. What is being transferred in transfer learning? In *Advances in neural information processing systems*, 2020. 1, 2
- [24] Arghya Pal and Vineeth N Balasubramanian. Zero-shot task transfer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2189–2198, 2019. 2
- [25] Sinno Jialin Pan, Qiang Yang, et al. A survey on transfer learning. *IEEE TKDE*, 22(10):1345–1359, 2010. 1, 2
- [26] Pau Panareda Busto and Juergen Gall. Open set domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 754–763, 2017. 2
- [27] Zhongyi Pei, Zhangjie Cao, Mingsheng Long, and Jianmin Wang. Multi-adversarial domain adaptation. In *AAAI*, 2018. 8
- [28] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. Visda: The visual domain adaptation challenge, 2017. 8
- [29] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*, pages 3723–3732, 2018. 2, 6
- [30] Kuniaki Saito, Shohei Yamamoto, Yoshitaka Ushiku, and Tatsuya Harada. Open set domain adaptation by backpropagation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 153–168, 2018. 2

- [31] Adam Santoro, Sergey Bartunov, Matthew Botvinick, et al. Meta-learning with memory-augmented neural networks. In *ICML*, pages 1842–1850, 2016. 2, 3
- [32] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. 7
- [33] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. In *Advances in Neural Information Processing Systems*, pages 527–538, 2018. 2
- [34] Richard Socher, Milind Ganjoo, Christopher D Manning, and Andrew Ng. Zero-shot learning through cross-modal transfer. In *Advances in neural information processing systems*, pages 935–943, 2013. 2
- [35] Baochen Sun and Kate Saenko. Deep coral: Correlation alignment for deep domain adaptation. In *ECCV*, pages 443–450, 2016. 2, 8
- [36] Qianru Sun, Yaoyao Liu, Tat-Seng Chua, and Bernt Schiele. Meta-transfer learning for few-shot learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 403–412, 2019. 2
- [37] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *CVPR*, pages 7167–7176, 2017. 2, 4
- [38] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint:1412.3474*, 2014. 2
- [39] Sait Unais, k v Gokul Lal, Prajapati Sunny, Bhaumik Rahul, Kumar Tarun, S Sanjana, and Bhalla Kriti. Curated dataset for covid-19 posterior-anterior chest radiography images (x-rays). 2020. 5
- [40] Hemant Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *CVPR*, pages 5018–5027, 2017. 8, 12
- [41] Jindong Wang, Wenjie Feng, Yiqiang Chen, Han Yu, Meiyu Huang, and Philip S Yu. Visual domain adaptation with manifold embedded distribution alignment. In *MM*, pages 402–410, 2018. 2
- [42] Linda Wang and Alexander Wong. Covid-net: A tailored deep convolutional neural network design for detection of covid-19 cases from chest radiography images. *arXiv preprint arXiv:2003.09871*, 2020. 1, 2
- [43] Ximei Wang, Ying Jin, Mingsheng Long, Jianmin Wang, and Michael I Jordan. Transferable normalization: Towards improving transferability of deep neural networks. In *NeurIPS*, pages 1951–1961, 2019. 6, 8
- [44] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *NIPS*, pages 3320–3328, 2014. 1, 2
- [45] Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3712–3722, 2018. 1, 2
- [46] Werner Zellinger, Thomas Grubinger, Edwin Lughofer, Thomas Natschläger, and Susanne Saminger-Platz. Central moment discrepancy (cmd) for domain-invariant representation learning. In *ICLR*, 2017. 2
- [47] Jing Zhang, Zewei Ding, Wanqing Li, and Philip Ogunbona. Importance weighted adversarial nets for partial domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8156–8164, 2018. 2
- [48] Jianpeng Zhang, Yutong Xie, Yi Li, Chunhua Shen, and Yong Xia. Covid-19 screening on chest x-ray images using deep learning based anomaly detection. *arXiv preprint arXiv:2003.12338*, 2020. 6
- [49] Weichen Zhang, Wanli Ouyang, Wen Li, and Dong Xu. Collaborative and adversarial network for unsupervised domain adaptation. In *CVPR*, pages 3801–3809, 2018. 2, 8
- [50] Yuchen Zhang, Tianle Liu, Mingsheng Long, and Michael I Jordan. Bridging theory and algorithm for domain adaptation. In *ICML*, 2019. 2, 4, 6
- [51] Yifan Zhang, Shuaicheng Niu, Zhen Qiu, Ying Wei, Peilin Zhao, Jianhua Yao, Junzhou Huang, Qingyao Wu, and Mingkui Tan. Covid-da: Deep domain adaptation from typical pneumonia to covid-19. *arXiv: 2005.01577*, 2020. 5, 6, 12

A. Evaluation of pivot data

In our main experiments, we set the $m = 8$ for the pivot data to select the top m samples belonging to one class with the highest probabilities. For pivot data construction, we further evaluate the performance of other two strategies: (1) select random m samples for each class and (2) select the bottom m samples for each class.

Here, ‘bottom m ’ is the opposite of top m in the main paper, which is selecting the m samples with the lowest probabilities. The results in Figure 6 show that the top m strategy achieves the best performance.

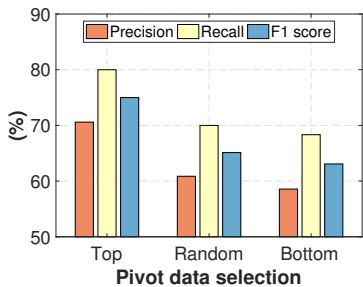


Figure 6. Different strategies on pivot data selection.

B. ROC curve

As for evaluation metrics, other than precision, recall, and F1 score, we further draw the ROC curve in Figure 7. The results show that our proposed TAN can achieve superior performance on this binary classification problem.

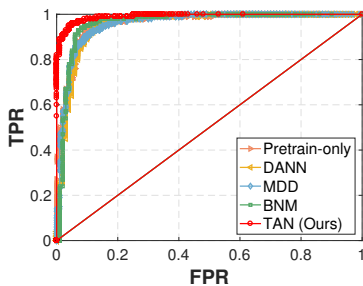


Figure 7. ROC curve.

We also compute the AUC (Area Under Curve) of our method and compare it with other baselines in Figure 8. The results show that other than target-train, which is the ideal state that uses the target domain labels for training, our method achieves the best AUC values compared to all other baselines.

C. Design criteria of task semantic adaptor

In this section, we pay special attention to the design criteria of the task semantic adaptation network M_θ . We ex-

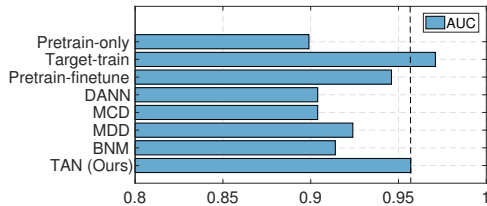


Figure 8. Area Under Curve (AUC) of different methods, indicating that the proposed TAN achieves the best AUC scores in all comparison methods other than the ideal state (target-train)

tensive analyze the following aspects: 1) network structure, 2) training criterion, and 3) activation function. Thorough analysis of these properties is valuable for designing better network and learning strategies in similar problems to ensure better performance.

C.1. Network structure

We design different structures of M_θ and record its performance in Table 9. As shown in the results, different structures produce different results and all results are better than comparison methods in Table 2 of the main paper. This means that the MLP structure of M_θ is effective. In general, complex structure is better at learning meaningful representations. In contrast, a simple structure may be worse in feature learning, but more difficult to overfit. Based on our experiments, we choose the structure in $-128 - 64 - 1$ for M_θ .

C.2. Training criterion

With regarding to the training criterion, we select different learning rates for the task semantic adaptor M_θ in $[0.0001, 0.0005, 0.001, 0.005, 0.01]$ and record the performance in Figure 9. The results indicate that M_θ can achieve similar performance with different learning rates. Therefore, to achieve the best performance we use the learning rate 0.0005.

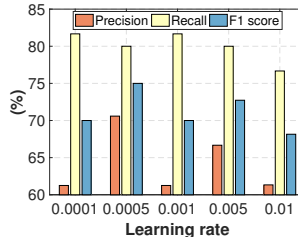


Figure 9. Different learning rates of M_θ

C.3. Activation function

We show the performance of different activation functions (i.e., σ in Eq. (7) of the main paper) in Table 10. The results show that tanh can generally lead to better per-

Table 8. Accuracy (%) on Office-Home dataset (ResNet-50)

Method	A→C	A→P	A→R	C→A	C→P	C→R	P→A	P→C	P→R	R→A	R→C	R→P	AVG
ResNet	34.9	50.0	58.0	37.4	41.9	46.2	38.5	31.2	60.4	53.9	41.2	59.9	46.1
DAN [18]	43.6	57.0	67.9	45.8	56.5	60.4	44.0	43.6	67.7	63.1	51.5	74.3	56.3
DANN [11]	45.6	59.3	70.1	47.0	58.5	60.9	46.1	43.7	68.5	63.2	51.8	76.8	57.6
JAN [21]	45.9	61.2	68.9	50.4	59.7	61.0	45.8	43.4	70.3	63.9	52.4	76.8	58.3
TAN (Ours)	46.9	61.3	71.4	49.0	62.8	63.9	53.6	48.8	73.4	66.9	45.4	75.2	59.9

Table 9. Different structures of M_θ

Structure	P (%)	R (%)	F1 (%)
in - 512 - 256 - 1	70.2	66.7	68.4
in - 256 - 128 - 1	69.0	66.7	67.8
in - 128 - 64 - 1	75.0	65.0	69.6

formance for our problem while sigmoid can also achieve good results. Therefore, we use tanh in this work.

Table 10. Different activation functions for M_θ

Activation func	P (%)	R (%)	F1 (%)
tanh	70.6	80.0	75.0
sigmoid	62.8	81.7	71.0
softplus	48.6	86.7	62.3
relu	67.2	71.7	69.4

D. Detailed information of experiments

For Bacterial and Viral datasets, the detailed information on source and target domain split in Table 11. The split information of COVID-DA is omitted which can be seen in [51] It is clear that the COVID-19 samples are not seen during training, which makes our unsupervised task transfer problem really challenging.

We also note that although the number of viral pneumonia samples is smaller than bacterial pneumonia samples, the transfer learning performance on viral dataset is almost the same as bacterial dataset (cf. Table 3 and 4 in the main paper, 95.2 vs. 95.4 F1 score). This maybe due to the high similarity between viral pneumonia and COVID-19 pneumonia as COVID-19 is also caused by a certain type of virus called ‘‘SARS-CoV-2’’. Again, this reflects the fact that similarity matters in transfer learning.

Table 11. Detailed information on domain split of Bacterial and Viral datasets.

Dataset	Domain	Symptom	#normal	#pneumonia	#COVID-19	#Total
Bacterial	source	Pneumonia	1,660	3,001	0	4,661
	target	COVID-19	1,610	0	1,281	2,891
Viral	source	Pneumonia	1,660	1,656	0	3,316
	target	COVID-19	1,610	0	1,281	2,891

E. Results on standard DA datasets

The results on Office-Home [40] dataset is shown in Table 8, and the results on VisDA-2017 dataset that uses all the classes are in Table 12. Note that TAN does not focus

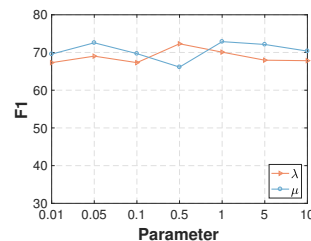
on the standard domain adaptation tasks, thus, we do not compare it with the latest DA methods. These results show that although the proposed TAN is not tailored for traditional domain adaptation tasks, it still achieves competitive results compared to several strong baselines. In the future, it is available to develop new TAN-based methods for traditional DA tasks to increase its performance.

Table 12. Accuracy (%) on VisDA-2017 dataset (ResNet-50)

Method	syth → real
ResNet [15]	52.4
DAN [18]	61.1
DANN [11]	57.4
JAN [21]	61.6
TAN (Ours)	64.0

F. Parameter sensitivity analysis

We show in Figure 10 that the two parameters λ, μ are relatively robust to different values, making our proposed method easy applicable to real applications.

Figure 10. Parameter sensitivity of λ and μ