MetaTeacher: Coordinating Multi-Model Domain Adaptation for Medical Image Classification

Anonymous Author(s) Affiliation Address email

Abstract

In medical image analysis, we often need to build an image recognition system for 1 a target scenario with the access to small labeled data and abundant unlabeled data, 2 as well as multiple related models pretrained on different source scenarios. This 3 presents the combined challenges of multi-source-free domain adaptation and semi-4 supervised learning *simultaneously*. However, both problems are typically studied 5 independently in the literature, and how to effectively combine existing methods is 6 non-trivial in design. In this work, we introduce a novel MetaTeacher framework 7 with three key components: (1) A learnable coordinating scheme for adaptive 8 domain adaptation of individual source models, (2) A mutual feedback mechanism 9 10 between the target model and source models for more coherent learning, and (3) A semi-supervised bilevel optimization algorithm for consistently organizing the 11 adaption of source models and the learning of target model. It aims to leverage 12 the knowledge of source models adaptively whilst maximize their complementary 13 benefits collectively to counter the challenge of limited supervision. Extensive 14 experiments on five chest x-ray image datasets show that our method outperforms 15 clearly all the state-of-the-art alternatives. 16

17 **1 Introduction**

Despite great stride made by existing deep learning methods on medical image classification re-18 sults [28, 46, 58], their performance often degrade drastically when applied to a new scenario. This 19 is mainly due to the domain shift challenge between the training and test data, caused by different 20 environments, different instruments, and different acquisition protocols. Unlike natural images, 21 annotating medical images requires special clinical expertise. It is hence more difficult to obtain 22 large-scale medical image datasets with high-quality labels at every single scenario. Domain adap-23 tation is a feasible solution, but comes with several limitations. Firstly, medical data is often under 24 strict privacy and license constraints. That means the source domain data is usually inaccessible 25 during domain adaptation. Secondly, medical data is typically multi-labeled which means that there 26 are multiple labels for a sample, and the multiple categories are not mutually exclusive. It has 27 more prominent different characteristics in different scenarios. Secondly, medical data is typically 28 multi-labeled which has more prominent different characteristics in different scenarios. Considering 29 these practical constraints, we propose a new Semi-supervised Multi-source-free Domain Adaptation 30 (SMDA) problem setting in the context of medical image classification. Our proposed setting has 31 three key conditions: (1) There are multiple source domain models trained on respective multi-label 32 33 medical image datasets, (2) All the source domain data is inaccessible for adaptation, and (3) The target domain data has only a small number of labelled samples along with abundant unlabeled data. 34

In medical image classification, there are limited domain adaptation works, with a need of accessing 35 the source domain data [4, 15, 20, 30, 38, 47, 49, 53]. Further, they usually consider a single source 36 37 domain. On the other hand, for employing multiple source domains, existing Multi-Source Domain Adaptation (MSDA) methods typically learn a common feature space for all source and target 38 39 domains [50] or use ensemble methods combined with source classifiers [6]. However, all of these MSDA methods require access to the source domain data. Regarding multi-label medical image 40 classification, there exists a solution which extends the standard classifier network by conditional 41 adversarial discriminator networks [39]. But it is still not source-free. Indeed, there have been 42 extensive study on Source-Free Domain Adaptation (SFDA) [31, 55]. However, they are not directly 43 applicable to our problem. Firstly, most of them assume a single source domain [31,55]. Using the 44 SFDA method to transfer each source domain model to the target domain separately and average 45 their predictions, this strategy cannot reveal the complementary information between different source 46 domains. Secondly, the source model is often domain biased. Different hospitals are featured with 47 different populations, leading to a situation that the source datasets focus on a specific set of class 48 labels. The existing SFDA methods can not assess the credibility of a source domain model with 49 different labels. 50

To address the above SMDA's limitations, employing knowledge distillation from multi-source 51 models to the target domain can be considered [14, 35, 56, 60, 61]. This forms a multi-teacher and 52 one-student scheme. In our problem setting, a few labels of the target domain are provided to judge 53 the credibility of multi-source models in different labels. In reality, it is common to exploit a few 54 labeled data in the target domain. Recent works [21, 25, 43, 44] have shown that a few labeled 55 data from the target domain can significantly improve the performance of the model. Inspired by 56 meta-learning approaches [33, 40, 42], we consider a bilevel optimization strategy to update both 57 the teachers and students. This is because different models vary in reliability and there is a need for 58 optimizing the update direction for each source model. This offers an opportunity of leveraging the 59 complementary and collaboration of different source models during model optimization, critical for 60 solving the low-supervision challenge. 61

Based on the above analysis, we propose a novel framework termed as MetaTeacher. It is based 62 on multi-teacher and one-student model. Each teacher model is pre-trained on a specific labeled 63 source data. The student model is initialized by a randomly chosen teacher. In order to provide 64 different update directions for multiple teachers, a coordinating weight learning method is proposed 65 66 to determine the contribution of each teacher for each target sample. In addition to knowledge transfer from multiple teachers, when adapting a specific teacher model, we also explore the feedback from 67 student and other teachers in a semi-supervised meta learning manner [13, 40]. Unlike the previous 68 MSDA approaches, MetaTeacher can adapt each teacher in different directions according to the 69 learned coordinating weight. This enable us to fully use different characteristics of source models, 70 71 whilst avoiding the problem of insufficient training samples for multi-label classification to some extent. 72

Our contributions are summarized as follows: (1) We propose a new problem setting, i.e., semi-73 supervised multi-source-free domain adaptation for multi-label medical image classification. To our 74 best knowledge, our work is the first exploration of multi-source-free and semi-supervised domain 75 adaptation in the field of transfer learning. (2) A novel framework MetaTeacher based on a multi-76 teacher and one-student scheme is introduced to solve the proposed SMDA problem. A mutual 77 feedback mechanism is designed based on meta-learning between the target model and the source 78 models for more coherent learning and adaptation. The knowledge from multiple source models 79 are sufficiently leveraged. (3) A coordinating weight learning method is derived for dynamically 80 revealing the performance differences of different source models over different classes. It is integrated 81 with the semi-supervised bilevel optimization algorithm for consistently updating the teacher and 82 student models. Extensive experiments on five well-known chest radiography datasets show that our 83 approach outperforms state-of-the-art alternatives clearly, along with in-detail ablation studies for 84 85 verifying the design choices of our model components.

86 2 Relate Works

Source-free domain adaptation. Source-free domain adaptation methods can be roughly divided 87 into two routes, i.e., generative approach [23, 24, 29, 54] and pseudo-label approach [3, 22, 31, 48]. 88 The generation approach generates target-style training samples to train the prediction model. Since 89 learning to generate features is difficult, this approach is extremely limited. The pseudo-label 90 approach generates pseudo-labels through the source domain model, which is simple and general and 91 has recently achieved good results in the machine learning community. The research of source-free 92 domain adaptation in the medical image analysis field mainly focuses on image segmentation. Bateson 93 et al. [3] maximized the mutual information between the target images and their label predictions 94 to perform spine, prostate and cardiac segmentation. Vibashan et al. [48] implemented source-free 95 domain adaptive image segmentation by generating pseudo-labels and applied self-training methods 96 for task-specific representation. These works are all conducted in the single-source domain case. 97 Currently, the research on multi-source-free domain adaptation is extremely limited, and most of the 98 works adapt the method of generating trusted pseudo-labels [1, 11]. 99

Multi-source domain adaptation for medical image classification. In machine learning commu-100 nity, MSDA works mainly have two strategies, i.e. distribution alignment [36, 64] and adversarial 101 learning [52, 62, 63]. The first strategy computes the statistical discrepancy between multi-source 102 domains and target domain, and then combines all predictions. The second strategy trains a domain 103 discriminator and forces the feature extraction network to learn domain-invariant features to confuse 104 the domain discriminator. For medical image classification, there only exist several shallow DA 105 models. Wang et al. [50] proposed to map multiple source and target data to a common latent space 106 for autism spectrum disorder classification. Cheng et al. [6] constructed a multi-domain transfer 107 classifier for the early diagnosis of Alzheimer's disease. All of these strategies require to access 108 source domain data and are not suitable for solving the proposed SMDA problem. To the best of our 109 knowledge, current teacher-student domain adaptation methods in the medical and machine learning 110 communities only consider the single-source domain case. When extending it to the multi-source 111 domain, it will face a challenging multi-objective optimization problem [8, 34]. 112

Semi-supervised domain adaptation (SSDA). Our problem is also related to SSDA which assumes 113 a small number of labeled samples in the target domain. Compared to UDA, using a few labeled 114 samples of the target domain allows to further achieve better domain alignment [27, 37, 57]. Due 115 to the shift of domain distribution, directly applying classical semi-supervised learning methods to 116 the SSDA problem will lead to sub-optimal performance. Representative SSDA works are based 117 on subspace learning [37, 57], entropy minimization [16, 43], label smoothing [10, 41] and active 118 119 learning [41,45]. However, all of these methods assume a single source domain with the source domain data accessible. Unlike these works, our method incorporates meta-learning and uses the 120 performance on the labeled target data as a feedback signal. 121

122 3 The Proposed Method

Problem statements. Suppose $D_T = \{(X_L^t, Y_L^t), X_U^t\}$ where Y_L^t denotes label annotations for a 123 small amount of target domain samples X_L^{t} and X_U^{t} for target domain samples without any label 124 annotations. The dimension of label vector is m. $D_{S_i} = \{(X_L^i, Y_L^i)\}$ where Y_L^i denotes label 125 annotations for *i*-th source domain samples X_L^i . For semi-supervised multi-source-free domain 126 adaptation problem, when the pretrained source classifiers f_{T_i} is applied to the target domain, the 127 source dataset D_{S_i} is not accessible for $i = 1, \dots, n$. Given source classifiers f_{T_i} for $i = 1, \dots, n$ 128 and the target data D_T , the task is to find a mapping $f_S: X_U^t \to Y_U^t$ where Y_U^t denotes the predicted 129 labels for target domain sample X_{U}^{t} that can work well in the target domain. 130

Overview. As shown in Fig.1, our framework is based on multi-teacher and one-student scheme. First,
multiple teacher models are pretrained according to each source domain, and then the student model
is initialed using a randomly chosen teacher model. They are all composed of a feature extractor
based on Resnet50 [17] and a multi-label classifier. The classifier consists of a fully connected layer,



Figure 1: Overview of MetaTeacher. (a) learns coordinating weight mapping which will be used later to provide guidance for updating teacher model. (b) alternately updates the teacher and student models. Each teacher is updated with feedback signals from student and other teachers.

where the input is an one-dimensional expanded feature, and the output is the probability of each
label. The objective function is the error loss between the predicted output and the ground truth.

Compared with traditional teacher-student model, our method has two differences: (1) coordinating 137 weight learning; (2) bilevel optimization. For the first part, a mapping is trained based on labeled 138 target domain samples, which fuses the multi-teacher predictions adaptively for each target sample. 139 This mapping will be used in the another part. In the initial iteration, the mapping and student model 140 are trained based on labeled target samples. In subsequent iterations, this part will only optimize the 141 mapping while the student model will be updated based on bilevel optimization. Then, in the bilevel 142 optimization part, the student and teacher models are updated alternately based on meta-learning. 143 Specifically, for an unlabeled target sample, a coordinating weight is generated, which provides 144 optimization direction for each teacher model. Finally, these two parts will be iteratively undated 145 until convergence. 146

147 3.1 Coordinating Weight Learning

As mentioned earlier, the teacher models are trained on different source domain data. Due to different distributions, they have different characteristics. Therefore, for a target domain sample, the classification probability of each teacher model is inconsistent. When we want to optimize a teacher model based on the target domain samples, the optimization direction of each teacher model should be different. So it is necessary to obtain the contribution weight of each teacher model to the final classification results. We call it coordinating weight. Fortunately, we can obtain the weight mapping with the labeled samples in the target domain.

As shown in Fig.1(a), for obtaining the coordinating weight, we first input the labeled target sample x_l^t into the student network, and get the output $B = g(x_l^t)$ from feature extraction network g, where $B \in R^{c \times h \times w}$. c, h, and w are the number of channels, height, and width respectively. Then, we perform a maximum pool operation on the feature map B to get $\psi \in R^{1 \times c}$ which retains the most important information of each channel. Our mapping consists of two learnable variables μ and ν , where $\mu \in R^{n \times 1}, \nu \in R^{c \times m}$. Then, we define a mapping $\phi = \mu \psi \nu \in R^{n \times m}$ for the target sample x_l^t . After normalizing, we get the coordinating weight matrix W where

$$V_{j,k} = \frac{exp(\phi_{j,k})}{\sum_{z=1}^{n} exp(\phi_{z,k})}.$$
(1)

Suppose for the sample x_l^t , the predictions of all teachers are formed as a matrix $P \in \mathbb{R}^{n \times m}$, by taking the Hadamard product between the teacher predictions and the coordinating weight matrix, we can get the fused prediction as the following,

ι

$$\bar{y}_l^t = Sum(P \circ W) \tag{2}$$

where $Sum(\cdot)$ means adding by rows. Denoting $\bar{y}_l^s = f_S(x_l^t; \theta_S)$ as the student prediction on the target sample x_l^t , we train the weight mapping and initialize student network using the following loss,

$$\mathcal{L}_W = \mathcal{L}\left(\bar{y}_l^s, y_l\right) + \alpha \mathcal{L}_{KL}\left(\bar{y}_l^t, \bar{y}_l^s\right) + \beta \left(\|\mu\| + \|\nu\|\right) \tag{3}$$

where $\mathcal{L}(\bar{y}_{l}^{s}, y_{l}) = \frac{1}{m} \sum_{i=1}^{m} [y_{l,i} log(\bar{y}_{l,i}^{s}) + (1 - y_{l,i}) log(1 - \bar{y}_{l,i}^{s})]$ represent the BCE (Binary Cross Entropy) loss, y_{l} is the groud truth, θ_{S} is the parameter of student network. $\mathcal{L}_{KL}(\bar{y}_{l}^{t}, \bar{y}_{l}^{s}) = \sum_{i=1}^{m} \bar{y}_{l,i}^{t} \log(\bar{y}_{l,i}^{t}/\bar{y}_{l,i}^{s})$ represents the KL (Kullback-Leibler divergence) loss which measures the distribution difference between the fused teacher prediction and student prediction. α and β are two balance parameters.

Remark. The mapping generates coordinating weight by Eq.(1), which not only reveals the complementarity of different teachers on different instances, but also, more interestingly, participates in the derivation of the update formula of teacher in the bilevel optimization process (see Appendix), providing a reference for the update direction of different teachers.

176 3.2 Bilevel Optimization

The bilevel optimization problem [5,7] was first proposed in the field of game theory. It includes an upper-level optimization task and a lower-level optimization task, where upper-level optimization task contains lower-level optimization task as a constraint. Here, the upper-level optimization task (student) provides feedback signals to the lower-level optimization tasks (teachers) through the performance on labeled data and the coordinating weight mapping. For an unlabeled target sample x_u^t , suppose the pseudo-label based on the learned coordinating weight mapping ϕ from multi-teachers Eq.(2) is \bar{y}_u^t and the corresponding coordinating weight matrix is W_u , we can define a loss function Γ_u as follows,

$$\Gamma_u(\theta_{T_1},\cdots,\theta_{T_n},\theta_S) = \mathcal{L}(\bar{y}_u^t,\bar{y}_u^s) \tag{4}$$

where $\bar{y}_{u}^{s} = f_{T}(x_{u}^{t}; \theta_{S}), \theta_{T_{i}}$ is the parameter of the *i*-th teacher network. Similarly, a loss function $\Gamma_{l}(\theta_{T_{1}}, \cdots, \theta_{T_{n}}, \theta_{S}) = \mathcal{L}(y_{l}, \bar{y}_{l}^{s})$ is defined for a labeled target samples x_{l}^{t} . In the bilevel optimization task, updating θ_{S} is the upper-level optimization task objective, while updating $\theta_{T_{1}}, \cdots, \theta_{T_{n}}$ is the lower-level optimization task objective. The upper-level optimization task and the lower-level optimization task are mutually constrained. To reach the lower-level optimization task objective, the performance of the upper-level optimization task objective on the labeled target data is utilized as feedback signal. So we get the objective function in lower-level optimization task as the following,

$$\underset{\theta_{T_1},\cdots,\theta_{T_n}}{\operatorname{argmin}} \Gamma_l\left(\theta_{T_1},\cdots,\theta_{T_n},\theta_S^{OP}\right) \quad \text{s.t.} \quad \theta_S^{OP} = \underset{\theta_S}{\operatorname{argmin}} \Gamma_u\left(\theta_{T_1},\cdots,\theta_{T_n},\theta_S\right). \tag{5}$$

Obviously, Eq.(5) cannot be optimized simply by gradient descent method, because the teacher models parameters can not be updated until θ_S reaches the optimum. We refer to the idea of meta-learning [13, 32, 40] and make a one-step approximation of the problem,

$$\theta_S^{OP} \approx \theta_S - \eta_S \cdot \nabla_{\theta_s} \Gamma_u \left(\theta_{T_1}, \theta_{T_2}, \cdots, \theta_{T_n}, \theta_S \right) \tag{6}$$

where η_S is the learning rate of the student network. Substitute Eq. (6) into Eq. (5) to obtain a new optimization objective function

$$\Gamma_l\left(\theta_{T_1},\cdots,\theta_{T_n},\theta_S-\eta_S\cdot\nabla_{\theta_s}\Gamma_u\left(\theta_{T_1},\theta_{T_2},\cdots,\theta_{T_n},\theta_S\right)\right).$$
(7)

¹⁹⁶ By optimizing Eq. (7) (see Appendix), we get the following update rules,

$$\theta_S' = \theta_S - \eta_S \cdot \nabla_{\theta_s} \Gamma_u, \tag{8}$$

$$\theta_{T_i}' = \theta_{T_i} - \eta_{T_i} \cdot \left[\left(\nabla_{\theta_S'} \Gamma_l \right)^T \cdot \nabla_{\theta_S} \Gamma_u \right]^T \cdot \nabla_{\theta_{T_i}} \mathcal{L} \left(\bar{y}_u^i, \tilde{y}_u^i \right)$$
(9)

for $i = 1, \dots, n$, where θ'_{S} and $\theta'_{T_{i}}$ are the updated parameters corresponding to the student and teachers respectively. $\bar{y}_{u}^{i} = f_{T_{i}} (x_{u}^{t}; \theta_{T_{i}}) \cdot W_{u}^{i}$ and W_{u}^{i} is the *i*th-row coordinating weight vector of W_{u} respect to the *i*-th teacher. \tilde{y}_{u}^{i} is the pseudo labels after normalizing the values of \bar{y}_{u}^{i} to 0 or 1, i.e., $\tilde{y}_{u,j}^{i} = 0$ when $\bar{y}_{u,j}^{i} < 0.5$ and $\tilde{y}_{u,j}^{i} = 1$ for other cases. Additionally, in order to prevent optimizing teachers in the same direction, the predictions of the updated multiple teachers should be as far away from each other as possible. So, we define a divergence loss as follows,

$$\mathcal{L}_D = -ln \sum_{j=1, j \neq i}^n \mathcal{L}_2\left(B_{T_i}\left(x_u^t; \theta_{T_i}\right), B_{T_j}\left(x_u^t; \theta_{T_j}\right)\right)$$
(10)

where $B_{T_i}(x_u; \theta_{T_i})$ represents the max-pooled results of the output feature map of the *i*-th teacher network. Here, we apply a max-pooling operation to the output features of multiple teachers and calculate the distance with L_2 norm. By requiring these feature maps to be far away each other, the optimization direction of teachers will be effectively adjusted. Finally, we update the *i*-th teacher network by the following rule,

$$\theta_{T_i}' = \theta_{T_i} - \eta_{T_i} \cdot \left(\left[(\nabla_{\theta_{S'}} \Gamma_l)^T \cdot \nabla_{\theta_S} \Gamma_u \right]^T \cdot \nabla_{\theta_{T_i}} \mathcal{L} \left(\bar{y}_u^i, \tilde{y}_u^i \right) + \gamma \nabla_{\theta_{T_i}} \mathcal{L}_D \right)$$
(11)

where γ is a hyperparameter.

Remark. Eq.(11) reveals that the update direction of θ_{T_i} is determined by three factors: (1) coordinating weight confuses feedback signals from different teachers; (2) student network parameters provide feedback signals and generate coordinating weight; (3) diversity constraint emphasizes the characteristic of different teacher networks. Interestingly, these three factors change over time during the meta-learning process. In addition to alternating updates of the student and teacher models, we also update the mapping periodically.

216 4 Experiments

Datasets. Five publicly available chest x-ray datasets are used to construct our multi-domain 217 adaptation scenarios. NIH-CXR14 [51] is a large public dataset of chest x-ray, which contains 108,948 218 front view x-ray images of 32,717 patients collected from NIH Clinical Center, with a total of 14 219 disease labels. MIMIC-CXR [19] contains 377,110 images and text reports, corresponding to 227,835 220 radiological studies conducted by Beth Israel Deaconess Medical Center in Boston, Massachusetts. 221 CheXpert [18] consists of 224,316 chest x-ray of 65,240 patients. The dataset collected chest x-ray 222 223 examinations and related radiology reports performed at inpatient and outpatient centers at Stanford Hospital from October 2002 to July 2017. Open-i [9] is collected by Indiana University Hospital 224 through the network from open source literature and biomedical image collection. It contains 3955 225 radiology reports, corresponding to 7470 frontal and lateral chest films. To be consistent with other 226 datasets, we filter out the side chest x-ray in Open-I, leaving only 3955 frontal images. Google-Health-227 CXR [2] is manually labeled by medical experts for CXR images with high accuracy and contains 228 about 4000 images. We follow the traditional UDA setting, and choose the disease closed set in these 229 five datasets as multi classification labels, i.e., Atelectasis, Cardiomegaly, Effusion, Consolidation, 230 Edema and Pneumonia. Four transfer scenarios are constructed, which are NIH-CXR14, CheXpert, 231 MIMIC-CXR to Open-i; NIH-CXR14, CheXpert, MIMIC-CXR to Google-Health-CXR; CheXpert, 232 MIMIC-CXR to NIH-CXR14 and NIH-CXR14, CheXpert to Open-i. 233

Implementation details. In order to make a compromise between images in different datasets, we 234 scale the images to 128*128 before feeding them into the network. To expand the training set, several 235 data augmentation techniques are used, including random cropping and horizontal flipping. SGD with 236 momentum of 0.9 is used as the optimizer. For the student model, the initial learning rate is 0.01 and 237 the weight decay is 5e-4. The learning rate for coordinating weight mapping is 0.001; For the teacher 238 models, the initial learning rate is 0.001 and the weight decay is 5e-6. The values of α , β and γ are 239 set as 0.5, 0.01 and 0.01 respectively. For the case when the target domains datasets are small-scale. 240 such as Open-i and Google-Health-CXR, we assume that there are 200 labeled data in the target 241 domains, and in order to give a good initial condition for training, we randomly select a source model 242 243 to initialize the target model. For the case when the target domains datasets are large-scale, such as NIH-CXR14, we assume that there are 500 labeled data in the target domains. Unless otherwise 244 specified, the interval for updating coordinating weight mapping is set as 100 iterations. Following 245

Method	Atelectasis	Cardiomegaly	Effusion	Consolidation	Edema	Pneumonia	Average
DECISION [1]	83.27	91.55	96.18	97.02	92.74	89.24	91.67
CAiDA [11]	82.45	92.16	95.12	95.92	89.89	90.37	90.99
SHOT-best [31]	81.48	91.22	94.19	95.10	88.96	89.58	90.09
MME [43]	82.44	90.82	95.46	96.07	90.26	87.20	90.38
ECACL [26]	82.60	92.18	96.32	95.97	90.70	89.61	91.23
Source Only(N)	83.09	87.20	96.11	95.10	86.87	77.40	87.63
Source Only(C)	82.26	87.64	94.71	96.61	90.22	75.12	87.76
Source Only(M)	80.63	91.31	94.87	94.53	84.91	82.78	88.05
Fine-tune(average)	82.14	88.71	95.32	95.52	88.77	78.48	88.16
Ours(w/o mapping)	79.99	92.64	98.22	93.64	95.50	84.54	90.76
Ours(w/o update)	81.98	90.72	95.76	95.51	89.40	82.53	89.32
Ours(all)	81.72	92.59	96.25	97.64	94.52	94.33	92.84

Table 1: Comparing the state-of-the-art methods on the transfer from *NIH-CXR14*, *CheXpert*, *MIMIC-CXR* to *Open-i*. Metric: AUROC.

Table 2: Comparing the state-of-the-art methods on the transfer from *NIH-CXR14*, *CheXpert*, *MIMIC-CXR* to *Google-Health-CXR*. Metric: AUROC.

Method	Atelectasis	Cardiomegaly	Effusion	Consolidation	Edema	Pneumonia	Average
DECISION [1]	77.24	81.71	85.94	79.03	83.48	83.68	81.85
CAiDA [11]	76.90	81.82	87.55	79.62	85.10	82.72	82.29
SHOT-best [31]	75.43	80.28	86.63	77.88	82.37	81.22	80.64
MME [43]	77.34	84.93	86.17	78.65	85.33	71.28	80.62
ECACL [26]	76.27	84.54	87.06	79.95	85.82	72.66	81.05
Source Only(N)	76.54	84.48	86.36	75.66	83.94	62.59	78.26
Source Only(C)	72.09	76.45	84.55	79.07	68.25	58.39	73.13
Source Only(M)	68.04	79.38	84.17	72.41	68.71	52.60	70.88
Fine-tune(average)	73.48	80.14	85.96	74.17	74.74	60.20	74.78
Ours(w/o <i>mapping</i>)	75.62	83.91	85.40	80.27	75.13	81.77	80.35
Ours(w/o update)	76.75	84.30	86.67	78.59	82.31	65.84	79.08
Ours(all)	77.65	79.52	88.73	78.74	86.73	84.78	82.69

the setting of multi-label medical image classification problems, the evaluation criterion is Area
Under the Receiver Operating Characteristic (AUROC) [12] curve score.

248 4.1 Comparisons to State-of-the-Art

At present, there does not exist any experimental report on our problem setting. So we choose 249 four category of methods for compare. The first category is Source only which means directly 250 applying a teacher model to the target domain. The second category is Fine-tune(average) which 251 fine-tune each teacher network using labeled target domain data, then average their predicted values. 252 The third category is the state-of-the-art multi-source-free domain adaptation methods, which are 253 DECISION [1], CAiDA [11], and SHOT-best. The SHOT-best refers to adapting each source domain 254 separately through the SHOT [31] method. The model with the best performance on the validation set 255 is selected. The final category is semi-supervised domain adaptation methods, which are MME [43] 256 and ECACL [26]. For the semi-supervised domain adaptation methods, we assume that the labeled 257 target data are the same as our method. Since they are single-source based methods, we perform 258 domain adaptation for each source model and take the best result. 259

Tables 1-4 show the comparison results on four transfer scenarios. Ours(all) is our proposed method.
 Source Only(N), Source Only(C) and Source Only(M) are the teacher models respect to the *NIH-CXR14*, *CheXpert* and *MIMIC-CXR* datasets respectively. For the scenario from *CheXpert*, *MIMIC-CXR* to *NIH-CXR14*, since the dataset *NIH-CXR14* contains 108,948 x-ray images, different from

Method	Atelectasis	Cardiomegaly	Effusion	Consolidation	Edema	Pneumonia	Average
DECISION [1]	72.99	80.73	79.37	75.52	82.30	71.38	77.05
CAiDA [11]	72.64	81.12	80.25	74.73	81.02	70.44	76.70
SHOT-best [31]	70.79	79.62	79.24	72.25	80.79	69.65	75.39
MME [43]	72.90	81.73	81.01	73.11	81.03	71.52	76.88
ECACL [26]	72.41	81.98	82.07	72.92	80.82	71.65	76.98
Source Only(N)	72.31	80.52	79.42	69.66	77.95	67.37	74.54
Source Only(C)	70.45	79.66	79.98	68.26	78.01	70.82	73.86
Fine-tune(average)	71.52	80.29	80.08	68.97	78.02	69.05	74.66
Ours(w/o <i>mapping</i>)	72.05	81.58	78.36	72.94	82.19	69.82	76.16
Ours(w/o update)	72.24	80.69	79.56	69.80	78.13	70.55	75.16
Ours(all)	73.63	86.64	80.86	72.24	86.68	66.37	77.74

Table 3: Comparing the state-of-the-art methods on the transfer from *CheXpert, MIMIC-CXR* to *NIH-CXR14*. Metric: AUROC.

Table 4: Comparing the state-of-the-art methods on the transfer from *NIH-CXR14*, *CheXpert* to *Open-i*. Metric: AUROC.

Method	Atelectasis	Cardiomegaly	Effusion	Consolidation	Edema	Pneumonia	Average
DECISION [1]	83.15	90.86	96.12	96.32	92.33	88.79	91.26
CAiDA [11]	82.38	91.97	94.89	95.30	89.81	90.44	90.80
SHOT-best [31]	81.48	91.22	94.19	95.10	88.96	89.58	90.09
MME [43]	81.46	90.40	94.86	97.73	89.79	87.31	90.26
ECACL [26]	82.22	88.76	96.04	96.85	92.43	87.90	90.70
Source Only(N)	83.09	87.20	96.11	95.10	86.87	77.40	87.63
Source Only(C)	82.26	87.64	94.71	96.61	90.22	75.12	87.76
Fine-tune(average)	82.66	87.98	95.85	95.67	88.58	77.02	87.96
Ours(w/o mapping)	83.73	93.37	96.04	97.30	91.51	82.34	90.72
Ours(w/o update)	82.70	88.91	95.47	95.48	88.96	78.85	88.40
Ours(all)	82.11	92.42	96.80	97.07	92.20	91.27	91.98

other scenarios, this time we do not need to initialize the target model with the source models. It can be observed that our method achieves the best performance. The extensive experiments on four different transfer scenarios verify the adaptability of our method under multi-label chest x-ray dataset transfer cases. For the scenario from *NIH-CXR14*, *CheXpert* to *Open-i*, as show in Table 4, the performance of two source domains is 0.86% lower than that of three source domains. Furthermore, MetaTeacher also has moderate training time and more clearer background (see Appendix).

270 4.2 Ablation Analysis and Discussion

Component analysis. In Tables 1-4, Ours(w/o mapping) represents that our proposed method 271 removes the part of coordinating weight learning and optimization substituted by average. Ours(w/o 272 update) means to remove the bilevel optimization process. In this situation, the weighted output 273 of teachers is used to supervise the learning of student network. The results in the last three rows 274 of Tables 1-4 show that these two parts are indispensable. It is worth mentioning that Ours(w/o 275 mapping) still obtains promising performance due to the following reasons. First, for student updating, 276 averaging predictions from multiple teachers is beneficial for student performance, consistent with 277 the finding by [59]. Second, the fixed W is also involved in the teacher optimization. It means bilevel 278 optimization contributes more gain to the overall performance than the coordinating weight learning. 279 However, the coordinating weight learning can judge which disease category the teacher is good at 280 by weight, knowledge with different weights can be learned from different teachers. Therefore, the 281 results in each disease category are close to the predictions of the best teachers, such as Pneumonia in 282 Table 1 and Atelectasis in Table 2 (also see Appendix). 283

Number(propotion)	Atelectasis	Cardiomegaly	Effusion	Consolidation	Edema	Pneumonia	Average		
50(1.4%)	82.48	92.22	95.19	96.10	89.96	90.58	91.09		
100(2.8%)	82.19	92.50	96.83	97.02	92.43	91.20	92.03		
200(5.6%)	81.72	92.59	96.25	97.64	94.52	94.33	92.84		
300(8.4%)	82.21	92.97	96.83	97.42	94.07	94.33	92.97		
Sensitivity analy 91 92 93 94 95 90 94 95 97 85 0 005 0.1 0.25 0.5 (2)	vsis of <i>a</i>	Sensitiv 92 92 91 92 92 91 92 91 92 91 90 90 90 88 88 88 88 90 88 88 90 88 89 88 90 90 90 90 90 90 90 90 90 90 90 90 90	ity analysis o	f β β β β β β β β	Sensiti	vity analysis of MetaTeach Baseline	Y er 1 lovalue		
(a)			(b)			(C)			

Table 5: Effect of the size of labeled target data on the transfer from *NIH-CXR14*, *CheXpert*, *MIMIC-CXR* to *Open-i*. Metric: AUROC.

Figure 2: Effect of different hyperparameters on the transfer from *NIH-CXR14*, *CheXpert*, *MIMIC-CXR* to *Open-i*. Baseline: source only(M).

Effects of proportion of labeled target data. Table 5 shows the influence of the amount of labeled
 data in the target domain on the transfer scenario of *NIH-CXR14*, *CheXpert*, *MIMIC-CXR* to *Open-i*.
 The experimental results show that the performance slowly improves as the amount of labeled data
 increases; a small number of labeled target domain samples can achieve good results.

Parameter analysis. We conduct parameter analysis experiment on the transfer scenario of NIH-288 CXR14, CheXpert, MIMIC-CXR to Open-i. The basic strategy is to change a parameter while other 289 parameters are fixed. Our method MetaTeacher has three hyperparameters, i.e., α and β in Eq. (3), 290 and γ in Eq.(11). Fig.(2)(a) shows performance changing with the parameter α . When $\alpha = 0$, the 291 coordinating weight mapping is not trained effectively resulting in the inability to determine the 292 optimization direction of each teacher. When α gradually increases to around 0.5, the result achieve 293 optimal performance. Fig.(2)(b) shows the influences of the parameters β . When the β is too large, it 294 means that the coordinating weight learning part is ineffective and cannot express the relationship 295 between the source domains. When β is set to 0, coordinating weight learning may overfit, which 296 may cause coordinating weights to work well on some instances but poorly on other instances; for 297 this case, the performance is 92.49% about 0.35% lower than the result 92.84% in Table 1. Fig.(2)(c) 298 shows the influences of the parameter γ on divergence loss. When γ is set to 0.01, the performance 299 reaches the best, but with the continuous increase of γ , the performance decreases obviously. When 300 $\gamma = 0$, the result is 92.29%, which is 0.55% lower. We can also see that our method is also quite 301 stable for the parameters α , β and γ in a large interval. 302

303 5 Conclusion

In this paper, we proposed a new framework, termed as MetaTeacher, for semi-supervised multi-304 source-free domain adaptation of medical image classification. The transfer learning process is 305 modeled as a multi-teacher and one-student scheme. We not only optimize student, but also optimize 306 teachers through student's feedback in the target domain. The optimization is based on meta-learning, 307 which consists of two main part: coordinating weight learning, and bilevel optimization. The first 308 part obtains the coordinating weight mapping which is used to coordinate the teacher outputs and 309 updates. Bilevel optimization updates the student base on the pseudo-labeled data produced by 310 teachers and updates each teacher base on the feedback signal generated by student and other teachers. 311 Extensive experiments on multi-label chest x-ray datasets empirically demonstrated the superiority of 312 313 our method over many state-of-the-art approaches.

314 **References**

- [1] AHMED, S. M., RAYCHAUDHURI, D. S., PAUL, S., OYMAK, S., AND ROY-CHOWDHURY, A. K.
 Unsupervised multi-source domain adaptation without access to source data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 10103–10112.
- [2] BALTRUSCHAT, I. M., NICKISCH, H., GRASS, M., KNOPP, T., AND SAALBACH, A. Comparison of deep learning approaches for multi-label chest x-ray classification. *Scientific reports* 9, 1 (2019), 1–10.
- [3] BATESON, M., DOLZ, J., KERVADEC, H., LOMBAERT, H., AND AYED, I. B. Source-free domain adaptation for image segmentation. *arXiv preprint arXiv:2108.03152* (2021).
- [4] BERMÚDEZ-CHACÓN, R., MÁRQUEZ-NEILA, P., SALZMANN, M., AND FUA, P. A domain-adaptive two-stream u-net for electron microscopy image segmentation. In 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018) (2018), IEEE, pp. 400–404.
- [5] BRACKEN, J., AND MCGILL, J. T. Mathematical programs with optimization problems in the constraints.
 Operations Research 21, 1 (1973), 37–44.
- [6] CHENG, B., LIU, M., SHEN, D., LI, Z., AND ZHANG, D. Multi-domain transfer learning for early diagnosis of alzheimer's disease. *Neuroinformatics* 15, 2 (2017), 115–132.
- [7] COLSON, B., MARCOTTE, P., AND SAVARD, G. An overview of bilevel optimization. *Annals of operations research 153*, 1 (2007), 235–256.
- [8] DEB, K. Multi-objective optimization. In Search methodologies. Springer, 2014, pp. 403–449.
- [9] DEMNER-FUSHMAN, D., KOHLI, M. D., ROSENMAN, M. B., SHOOSHAN, S. E., RODRIGUEZ, L.,
 ANTANI, S., THOMA, G. R., AND MCDONALD, C. J. Preparing a collection of radiology examinations
 for distribution and retrieval. *Journal of the American Medical Informatics Association 23*, 2 (2016),
 304–310.
- [10] DONAHUE, J., HOFFMAN, J., RODNER, E., SAENKO, K., AND DARRELL, T. Semi-supervised domain
 adaptation with instance constraints. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2013), pp. 668–675.
- [11] DONG, J., FANG, Z., LIU, A., SUN, G., AND LIU, T. Confident anchor-induced multi-source free domain
 adaptation. *Advances in Neural Information Processing Systems 34* (2021).
- 12] FAWCETT, T. An introduction to roc analysis. Pattern recognition letters 27, 8 (2006), 861–874.
- [13] FINN, C., ABBEEL, P., AND LEVINE, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning* (2017), PMLR, pp. 1126–1135.
- FURLANELLO, T., LIPTON, Z., TSCHANNEN, M., ITTI, L., AND ANANDKUMAR, A. Born again neural
 networks. In *International Conference on Machine Learning* (2018), PMLR, pp. 1607–1616.
- [15] GAO, Y., ZHANG, Y., CAO, Z., GUO, X., AND ZHANG, J. Decoding brain states from fmri signals by
 using unsupervised domain adaptation. *IEEE Journal of Biomedical and Health Informatics 24*, 6 (2019),
 1677–1685.
- [16] GRANDVALET, Y., AND BENGIO, Y. Semi-supervised learning by entropy minimization. Advances in neural information processing systems 17 (2004).
- [17] HE, K., ZHANG, X., REN, S., AND SUN, J. Deep residual learning for image recognition. In *Proceedings* of the IEEE conference on computer vision and pattern recognition (2016), pp. 770–778.
- IRVIN, J., RAJPURKAR, P., KO, M., YU, Y., CIUREA-ILCUS, S., CHUTE, C., MARKLUND, H.,
 HAGHGOO, B., BALL, R., SHPANSKAYA, K., ET AL. Chexpert: A large chest radiograph dataset with
 uncertainty labels and expert comparison. In *Proceedings of the AAAI conference on artificial intelligence* (2019), vol. 33, pp. 590–597.
- JOHNSON, A. E., POLLARD, T. J., GREENBAUM, N. R., LUNGREN, M. P., DENG, C.-Y., PENG, Y.,
 LU, Z., MARK, R. G., BERKOWITZ, S. J., AND HORNG, S. Mimic-cxr-jpg, a large publicly available
 database of labeled chest radiographs. *arXiv preprint arXiv:1901.07042* (2019).
- [20] KAMPHENKEL, J., JÄGER, P. F., BICKELHAUPT, S., LAUN, F. B., LEDERER, W., DANIEL, H.,
 KUDER, T. A., DELORME, S., SCHLEMMER, H.-P., KÖNIG, F., ET AL. Domain adaptation for deviating
 acquisition protocols in cnn-based lesion classification on diffusion-weighted mr images. In *Image Analysis for Moving Organ, Breast, and Thoracic Images.* Springer, 2018, pp. 73–80.
- [21] KIM, T., AND KIM, C. Attract, perturb, and explore: Learning a feature alignment network for semisupervised domain adaptation. In *European conference on computer vision* (2020), Springer, pp. 591–607.
- [22] KIM, Y., CHO, D., HAN, K., PANDA, P., AND HONG, S. Domain adaptation without source data. *arXiv preprint arXiv:2007.01524* (2020).

- [23] KUNDU, J. N., VENKAT, N., BABU, R. V., ET AL. Universal source-free domain adaptation. In
 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (2020), pp. 4544–4553.
- [24] KURMI, V. K., SUBRAMANIAN, V. K., AND NAMBOODIRI, V. P. Domain impression: A source data
 free domain adaptation method. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (2021), pp. 615–625.
- [25] LI, B., WANG, Y., ZHANG, S., LI, D., KEUTZER, K., DARRELL, T., AND ZHAO, H. Learning invariant
 representations and risks for semi-supervised domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 1104–1113.
- [26] LI, K., LIU, C., ZHAO, H., ZHANG, Y., AND FU, Y. Ecacl: A holistic framework for semi-supervised domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 8578–8587.
- [27] LI, L., AND ZHANG, Z. Semi-supervised domain adaptation by covariance matching. *IEEE transactions* on pattern analysis and machine intelligence 41, 11 (2018), 2724–2739.
- [28] LI, Q., CAI, W., WANG, X., ZHOU, Y., FENG, D. D., AND CHEN, M. Medical image classification with
 convolutional neural network. In 2014 13th international conference on control automation robotics &
 vision (ICARCV) (2014), IEEE, pp. 844–848.
- [29] LI, R., JIAO, Q., CAO, W., WONG, H.-S., AND WU, S. Model adaptation: Unsupervised domain
 adaptation without source data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020), pp. 9641–9650.
- [30] LI, W., ZHAO, Y., CHEN, X., XIAO, Y., AND QIN, Y. Detecting alzheimer's disease on small dataset:
 A knowledge transfer perspective. *IEEE journal of biomedical and health informatics 23*, 3 (2018),
 1234–1242.
- [31] LIANG, J., HU, D., AND FENG, J. Do we really need to access the source data? source hypothesis transfer
 for unsupervised domain adaptation. In *International Conference on Machine Learning* (2020), PMLR,
 pp. 6028–6039.
- [32] LIU, H., SIMONYAN, K., AND YANG, Y. Darts: Differentiable architecture search. arXiv preprint
 arXiv:1806.09055 (2018).
- [33] MADHAWA, K., AND MURATA, T. Metal: Active semi-supervised learning on graphs via meta-learning.
 In Asian Conference on Machine Learning (2020), PMLR, pp. 561–576.
- [34] MARLER, R. T., AND ARORA, J. S. Survey of multi-objective optimization methods for engineering.
 Structural and multidisciplinary optimization 26, 6 (2004), 369–395.
- [35] PARK, S., AND KWAK, N. Feature-level ensemble knowledge distillation for aggregating knowledge from multiple networks. In *ECAI 2020*. IOS Press, 2020, pp. 1411–1418.
- PENG, X., BAI, Q., XIA, X., HUANG, Z., SAENKO, K., AND WANG, B. Moment matching for multi source domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision* (2019), pp. 1406–1415.
- [37] PEREIRA, L. A., AND DA SILVA TORRES, R. Semi-supervised transfer subspace for domain adaptation.
 Pattern Recognition 75 (2018), 235–249.
- [38] PERONE, C. S., BALLESTER, P., BARROS, R. C., AND COHEN-ADAD, J. Unsupervised domain
 adaptation for medical imaging segmentation with self-ensembling. *NeuroImage 194* (2019), 1–11.
- [39] PHAM, D., KOESNADI, S., DOVLETOV, G., AND PAULI, J. Unsupervised adversarial domain adaptation
 for multi-label classification of chest x-ray. In 2021 IEEE 18th International Symposium on Biomedical
 Imaging (ISBI) (2021), IEEE, pp. 1236–1240.
- [40] PHAM, H., DAI, Z., XIE, Q., AND LE, Q. V. Meta pseudo labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 11557–11568.
- [41] PRABHU, V., CHANDRASEKARAN, A., SAENKO, K., AND HOFFMAN, J. Active domain adaptation via clustering uncertainty-weighted embeddings. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 8505–8514.
- [42] REN, M., TRIANTAFILLOU, E., RAVI, S., SNELL, J., SWERSKY, K., TENENBAUM, J. B., LAROCHELLE,
 H., AND ZEMEL, R. S. Meta-learning for semi-supervised few-shot classification. *arXiv preprint arXiv:1803.00676* (2018).
- [43] SAITO, K., KIM, D., SCLAROFF, S., DARRELL, T., AND SAENKO, K. Semi-supervised domain
 adaptation via minimax entropy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2019), pp. 8050–8058.
- 423 [44] SINGH, A. Clda: Contrastive learning for semi-supervised domain adaptation. *Advances in Neural* 424 *Information Processing Systems 34* (2021).

- [45] SU, J.-C., TSAI, Y.-H., SOHN, K., LIU, B., MAJI, S., AND CHANDRAKER, M. Active adversarial
 domain adaptation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (2020), pp. 739–748.
- [46] TALEB, A., LOETZSCH, W., DANZ, N., SEVERIN, J., GAERTNER, T., BERGNER, B., AND LIPPERT, C.
 3d self-supervised methods for medical imaging. *Advances in Neural Information Processing Systems 33* (2020), 18158–18172.
- [47] VAN OPBROEK, A., VERNOOIJ, M. W., IKRAM, M. A., AND DE BRUIJNE, M. Weighting training
 images by maximizing distribution similarity for supervised segmentation across scanners. *Medical image analysis* 24, 1 (2015), 245–254.
- [48] VS, V., VALANARASU, J. M. J., AND PATEL, V. M. Target and task specific source-free domain adaptive
 image segmentation. *arXiv preprint arXiv:2203.15792* (2022).
- [49] WACHINGER, C., REUTER, M., INITIATIVE, A. D. N., ET AL. Domain adaptation for alzheimer's disease
 diagnostics. *Neuroimage 139* (2016), 470–479.
- [50] WANG, J., ZHANG, L., WANG, Q., CHEN, L., SHI, J., CHEN, X., LI, Z., AND SHEN, D. Multi-class asd
 classification based on functional connectivity and functional correlation tensor via multi-source domain
 adaptation and multi-view sparse representation. *IEEE transactions on medical imaging 39*, 10 (2020),
 3137–3147.
- WANG, X., PENG, Y., LU, L., LU, Z., BAGHERI, M., AND SUMMERS, R. M. Chestx-ray8: Hospital scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common
 thorax diseases. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017),
 pp. 2097–2106.
- [52] XU, R., CHEN, Z., ZUO, W., YAN, J., AND LIN, L. Deep cocktail network: Multi-source unsupervised
 domain adaptation with category shift. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2018), pp. 3964–3973.
- YAN, W., WANG, Y., GU, S., HUANG, L., YAN, F., XIA, L., AND TAO, Q. The domain shift problem of medical image segmentation and vendor-adaptation by unet-gan. In *International Conference on Medical Image Computing and Computer-Assisted Intervention* (2019), Springer, pp. 623–631.
- 452 [54] YANG, C., GUO, X., CHEN, Z., AND YUAN, Y. Source free domain adaptation for medical image 453 segmentation with fourier style mining. *Medical Image Analysis* (2022), 102457.
- 454 [55] YANG, S., WANG, Y., VAN DE WEIJER, J., HERRANZ, L., AND JUI, S. Unsupervised domain adaptation 455 without source data by casting a bait. *arXiv e-prints* (2020), arXiv–2010.
- YANG, Z., SHOU, L., GONG, M., LIN, W., AND JIANG, D. Model compression with two-stage multi teacher knowledge distillation for web question answering system. In *Proceedings of the 13th International Conference on Web Search and Data Mining* (2020), pp. 690–698.
- [57] YAO, T., PAN, Y., NGO, C.-W., LI, H., AND MEI, T. Semi-supervised domain adaptation with subspace
 learning for visual recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition* (2015), pp. 2142–2150.
- 462 [58] YOSINSKI, J., CLUNE, J., BENGIO, Y., AND LIPSON, H. How transferable are features in deep neural 463 networks? *Advances in neural information processing systems* 27 (2014).
- 464 [59] YOU, S., XU, C., XU, C., AND TAO, D. Learning from multiple teacher networks. In *Proceedings* 465 of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2017),
 466 pp. 1285–1294.
- YUAN, F., SHOU, L., PEI, J., LIN, W., GONG, M., FU, Y., AND JIANG, D. Reinforced multi-teacher
 selection for knowledge distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence* (AAAI'21) (2021).
- [61] ZHAO, H., SUN, X., DONG, J., CHEN, C., AND DONG, Z. Highlight every step: Knowledge distillation
 via collaborative teaching. *IEEE Transactions on Cybernetics* (2020).
- [62] ZHAO, H., ZHANG, S., WU, G., MOURA, J. M., COSTEIRA, J. P., AND GORDON, G. J. Adversarial
 multiple source domain adaptation. *Advances in neural information processing systems 31* (2018).
- [63] ZHAO, S., WANG, G., ZHANG, S., GU, Y., LI, Y., SONG, Z., XU, P., HU, R., CHAI, H., AND
 KEUTZER, K. Multi-source distilling domain adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2020), vol. 34, pp. 12975–12983.
- [64] ZHU, Y., ZHUANG, F., AND WANG, D. Aligning domain-specific distribution and classifier for cross domain classification from multiple sources. In *Proceedings of the AAAI Conference on Artificial Intelli- gence* (2019), vol. 33, pp. 5989–5996.

480 Checklist

 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes] (b) Did you describe the limitations of your work? [No] (c) Did you discuss any potential negative societal impacts of your work? [No] (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes] (a) Did you are including theoretical results (a) Did you are including theoretical results (a) Did you are include the ethics review guidelines and ensured that your paper conforms to them? [Yes] (b) Did you include complete proofs of all theoretical results? [Yes] (c) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] (e) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details. (c) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) Did you include any new assets either in the supplemental material or as a URL? [No] (e) Did you usics whether and how consent was obtained from people whose data you're using/curating? [No] (b) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (c) Did you include the stimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	481	1. For all authors
 (b) Did you describe the limitations of your work? [No] (c) Did you discuss any potential negative societal impacts of your work? [No] (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes] (a) Did you are including theoretical results (a) Did you state the full set of assumptions of all theoretical results? [Yes] (b) Did you include complete proofs of all theoretical results? [Yes] (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details. (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If you montion the license of the assets? [No] (b) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you include the stimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	482 483	 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 (c) Did you discuss any potential negative societal impacts of your work? [No] (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes] 2. If you are including theoretical results (a) Did you state the full set of assumptions of all theoretical results? [Yes] (b) Did you include complete proofs of all theoretical results? [Yes] 3. If you ran experiments (a) Did you upccify all the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] (c) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details. (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you include any new assets either in the supplemental material or as a URL? [No] (c) Did you include any new assets either in the supplemental material or as a URL? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you include the estimated hourly wage paid to participants and the total amount spent on participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did yo	484	(b) Did you describe the limitations of your work? [No]
 (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes] If you are including theoretical results (a) Did you state the full set of assumptions of all theoretical results? [Yes] (b) Did you include complete proofs of all theoretical results? [Yes] If you ran experiments (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details. (c) Did you erport error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (b) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] 	485	(c) Did you discuss any potential negative societal impacts of your work? [No]
 If you are including theoretical results (a) Did you state the full set of assumptions of all theoretical results? [Yes] (b) Did you include complete proofs of all theoretical results? [Yes] If you ran experiments (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details. (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you include any new assets either in the supplemental material or as a URL? [No] (c) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	486 487	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
 (a) Did you state the full set of assumptions of all theoretical results? [Yes] (b) Did you include complete proofs of all theoretical results? [Yes] 3. If you ran experiments (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details. (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you include any new assets either in the supplemental material or as a URL? [No] (c) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	488	2. If you are including theoretical results
 (b) Did you include complete proofs of all theoretical results? [Yes] 3. If you ran experiments (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details. (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you include any new assets either in the supplemental material or as a URL? [No] (c) Did you include any new assets either in the supplemental material or as a URL? [No] (c) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (c) Did you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	489	 (a) Did you state the full set of assumptions of all theoretical results? [Yes] (b) Did you state the full set of assumptions of all theoretical results?
 3. If you ran experiments (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details. (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you include any new assets either in the supplemental material or as a URL? [No] (c) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	490	(b) Did you include complete proofs of all theoretical results? [Yes]
 (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details. (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you include any new assets either in the supplemental material or as a URL? [No] (c) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount specific participant compensation? [N/A] 	491	3. If you ran experiments
 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details. (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you include any new assets either in the supplemental material or as a URL? [No] (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you discribe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	492 493	 (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]
 (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you mention the license of the assets? [No] (c) Did you include any new assets either in the supplemental material or as a URL? [No] (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	494 495	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See the Implementation details.
 (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details. 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you mention the license of the assets? [No] (c) Did you include any new assets either in the supplemental material or as a URL? [No] (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	496 497	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you mention the license of the assets? [No] (c) Did you include any new assets either in the supplemental material or as a URL? [No] (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	498 499	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the Implementation details.
 (a) If your work uses existing assets, did you cite the creators? [Yes] (b) Did you mention the license of the assets? [No] (c) Did you include any new assets either in the supplemental material or as a URL? [No] (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	500	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
 (b) Did you mention the license of the assets? [No] (c) Did you include any new assets either in the supplemental material or as a URL? [No] (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	501	(a) If your work uses existing assets, did you cite the creators? [Yes]
 (c) Did you include any new assets either in the supplemental material or as a URL? [No] (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	502	(b) Did you mention the license of the assets? [No]
 (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	503	(c) Did you include any new assets either in the supplemental material or as a URL? [No]
 (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	504 505	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No]
 5. If you used crowdsourcing or conducted research with human subjects (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	506 507	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No]
 (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	508	5. If you used crowdsourcing or conducted research with human subjects
 (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	509 510	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] 	511 512	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
	513 514	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]