# Transferability-Guided Cross-Domain Cross-Task Transfer Learning

Yang Tan, Enming Zhang, Yang Li, Shao-Lun Huang, Xiao-Ping Zhang, *Fellow, IEEE*

*Abstract*—We propose two novel transferability metrics F-OTCE (Fast Optimal Transport based Conditional Entropy) and JC-OTCE (Joint Correspondence OTCE) to evaluate how much the source model (task) can benefit the learning of the target task and to learn more generalizable representations for cross-domain cross-task transfer learning. Unlike the original OTCE metric that requires evaluating the empirical transferability on auxiliary tasks, our metrics are auxiliary-free such that they can be computed much more efficiently. Specifically, F-OTCE estimates transferability by first solving an Optimal Transport (OT) problem between source and target distributions, and then uses the optimal coupling to compute the Negative Conditional Entropy between source and target labels. It can also serve as an objective function to enhance downstream transfer learning tasks including model finetuning and domain generalization. Meanwhile, JC-OTCE improves the transferability accuracy of F-OTCE by including label distances in the OT problem, though it incurs additional computation costs. Extensive experiments demonstrate that F-OTCE and JC-OTCE outperform state-of-the-art auxiliary-free metrics by 21.1% and 25.8%, respectively in correlation coefficient with the ground-truth transfer accuracy. By eliminating the training cost of auxiliary tasks, the two metrics reduce the total computation time of the previous method from 43 minutes to 9.32s and 10.78s, respectively, for a pair of tasks. When applied in the model finetuning and domain generalization tasks, F-OTCE shows significant improvements in the transfer accuracy in few-shot classification experiments, with up to 4.41% and 2.34% accuracy gains, respectively.

*Index Terms*—Transfer learning, few-shot learning, transferability estimation, task relatedness, cross-domain, cross-task, source selection.

## I. INTRODUCTION

**T**RANSFER learning is an effective learning paradigm to enhance the performance on target tasks via leveraging prior knowledge from the related source tasks (or source models), especially when there are only few labeled data for supervision [1]–[4]. However, the success of transfer learning is not always guaranteed. If the source and target tasks are unrelated, or if the transferred representation does not carry sufficient information about the target task, transfer learning will not obtain a notable gain on the target task performance, and may even experience negative transfer, i.e., the performance becomes worse than that of training from scratch on the target task [5]. Therefore, understanding when and what to transfer between tasks is crucial to the success of transfer learning.

Yang Tan, Enming Zhang, Yang Li, Shao-Lun Huang and Xiao-Ping Zhang are with Shenzhen Key Laboratory of Ubiquitous Data Enabling, Shenzhen International Graduate School, Tsinghua University. The corresponding author is Yang Li (e-mail: yangli@sz.tsinghua.edu.cn).
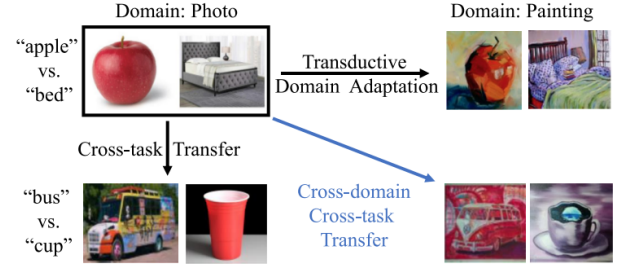


Fig. 1. Illustration of three different transfer learning settings, i.e., transductive domain adaptation [6], cross-task transfer [7] and the cross-domain cross-task transfer we investigating.

The "when to transfer" problem was traditionally studied theoretically through the derivation of generalization bounds of transfer learning across tasks [8], [9] and across domains (also known as the domain adaptation problem) [10]–[15]. Such studies bound the target task generalization error by a function that depends on certain divergence between the source and target domain, or the complexity of the hypothesis class for the source and target tasks. In practice, however, these bounds are difficult to compute from data and they tend to rely on strict assumptions that can not be verified. In recent years, the notion of task transferability was proposed to address the "when to transfer" problem in the context of deep transfer learning [7], [16]–[23]. The transferability problem aims to quantitatively evaluate how much the source task or source model could benefit the learning of the target task. It can be used to directly select the most "transferable" source model from a model zoo for a target task, rather than exhaustively trying each source model on the target data. In addition, transferability can help prioritize different tasks for joint training [16] and multi-source feature fusion [22].

As empirical transferability studies [16]–[19], [24] incur heavy computation burdens in retraining the transfer learning model on the target training data, a new trend of transferability research aims to efficiently estimate the transfer performance a-priori with little or no training of the transfer model. Several efficient transferability metrics are proposed, including NCE [20], H-score [7], LEEP [21], and LogME [23]. Despite being evidently more efficient to compute from practical data than empirical methods, they are also prone to strict data assumptions [7], [20] and insufficient performance [21], [23] while task complexities are similar. Moreover, the aforementioned metrics are solely used for determining when to transfer between a pair of source and target tasks, but they do not contribute to solving the "what to transfer" problem, i.e. how

to obtain more generalizable feature representations across domains and tasks.

Recently, a novel transferability metric **OTCE** (Optimal Transport based Conditional Entropy) [22] is proposed to effectively estimate the transferability under the challenging cross-domain cross-task transfer setting, as shown in Fig. 1. Unlike the transferability metrics mentioned earlier, OTCE adopts a more analytical disentanglement approach. It explicitly assesses the domain difference (measured by Wasserstein distance) and the task difference (determined by conditional entropy) between tasks, and then integrates them via a linear model to quantify transferability. This technically sound design yields substantial accuracy improvement over the aforementioned metrics. Nevertheless, a major limitation of OTCE is its dependency on auxiliary tasks with known transfer performance to determine the intrinsic parameters of the linear model. On one hand, the availability of sufficient labeled data for creating auxiliary tasks is not always guaranteed. On the other hand, assessing the transfer performance of such auxiliary tasks necessitates retraining the source model, incurring additional computation costs. As a result, the reliance on auxiliary tasks makes OTCE relatively inefficient and less applicable in general practical scenarios.

In this paper, we aim to broaden the applicability of the OTCE framework and investigate the potential uses of transferability in downstream transfer learning tasks. We propose two auxiliary-free transferability metrics, namely **F-OTCE** (Fast OTCE) and **JC-OTCE** (Joint Correspondence OTCE), which eliminate the need for auxiliary tasks and substantially enhance the efficiency without compromising accuracy. For classification problems, the F-OTCE metric solves the Optimal Transport (OT) problem [25], [26] to estimate a probabilistic coupling between the unpaired samples from the source and target datasets. Then the optimal coupling enables us to derive the negative conditional entropy between the source and target task labels for representing transferability, which measures the label uncertainty of a target sample given the labels of corresponding source samples. While the F-OTCE metric does not explicitly evaluate the domain difference, the estimated probabilistic coupling between the source and target data implicitly captures the domain difference to some extent in this unified framework.

Then we propose the JC-OTCE metric to further improve the accuracy of the F-OTCE metric in diverse transfer configurations. Our motivation is that F-OTCE only considers the joint probability distribution of input samples when determining data correspondences between the source and target domains. But this approach is limited because the definition (label annotations) of the source task can also affect model generalization. To address this limitation, we incorporate label distance into the ground cost of the OT problem, allowing for the computation of correspondences in both sample and label spaces. By including additional label information, JC-OTCE produces improved data correspondences that partially compensate for the lack of explicit domain difference consideration. JC-OTCE achieves comparable transferability accuracy to the original OTCE metric but requires additional computation compared to F-OTCE, which remains preferable for efficiency purposes.

Moreover, we investigate the application of our transferability metric in two downstream transfer learning tasks including *model finetuning* and *domain generalization*, offering a solution to the "what to transfer" problem. Specifically, to enhance the model finetuning performance, we propose an OTCE-based finetune algorithm that optimizes the pretrained source model to learn more transferable feature representation via maximizing the F-OTCE score between the source and target tasks. The optimized model is then finetuned on target training data using the classification loss function.

We also demonstrate that incorporating the F-OTCE metric into a novel domain generalization method URL [27] can further improve its generalizability on unseen domains. Our motivation is to view distilling knowledge from domain-specific models to the universal model as maximizing the transferability between them. Therefore, we replace the knowledge distillation function in URL with our F-OTCE score, resulting in significant accuracy improvements in few-shot classification tasks on unseen domains.

This work is an extension of our previous conference paper [22], and the additional contributions are summarized as follows:

1) *Expanding the applicability of OTCE framework.* Our proposed F-OTCE and JC-OTCE metrics eliminate the need for auxiliary tasks and achieve comparable transferability accuracy to OTCE. They also outperform previous auxiliary-free transferability metrics in terms of accuracy while maintaining comparable efficiency.

2) *Investigating the potential uses of transferability.* We illustrate the effectiveness of using F-OTCE as an optimization objective in improving the performance of downstream tasks, such as model finetuning and domain generalization. We consider F-OTCE to be a general tool that can be easily integrated into various algorithms for transfer learning and other related applications.

In our experiments using several multi-domain classification datasets, we show that our proposed two metrics significantly outperform existing auxiliary-free metrics with $25.8\%$ correlation gain on average, while cutting more than $99\%$ of the computation time in the original OTCE. We also show that, when served as a loss function, F-OTCE leads to notable classification accuracy gains on the model finetuning and domain generalization tasks, with up to $4.41\%$ and $2.34\%$. The rest of this paper is organized as follows. Section II introduces the formulation of transferability. Section III provides a preliminary analysis on OTCE. Section IV presents our two auxiliary-free transferability metrics. Section V illustrates our proposed transferability-guided model finetuning and domain generalization algorithms. Section VI provides all the experimental results and analyses. Finally, we draw the conclusion in Section VII.

## II. TRANSFERABILITY FORMULATION

Here we introduce the formal definition of transferability. Suppose we have source data $D_s = \{(x_s^i, y_s^i)\}_{i=1}^m \sim P_s(x, y)$ and target data $D_t = \{(x_t^i, y_t^i)\}_{i=1}^n \sim P_t(x, y)$, where $x$ represents the input instance and $y$ denotes the label. We have
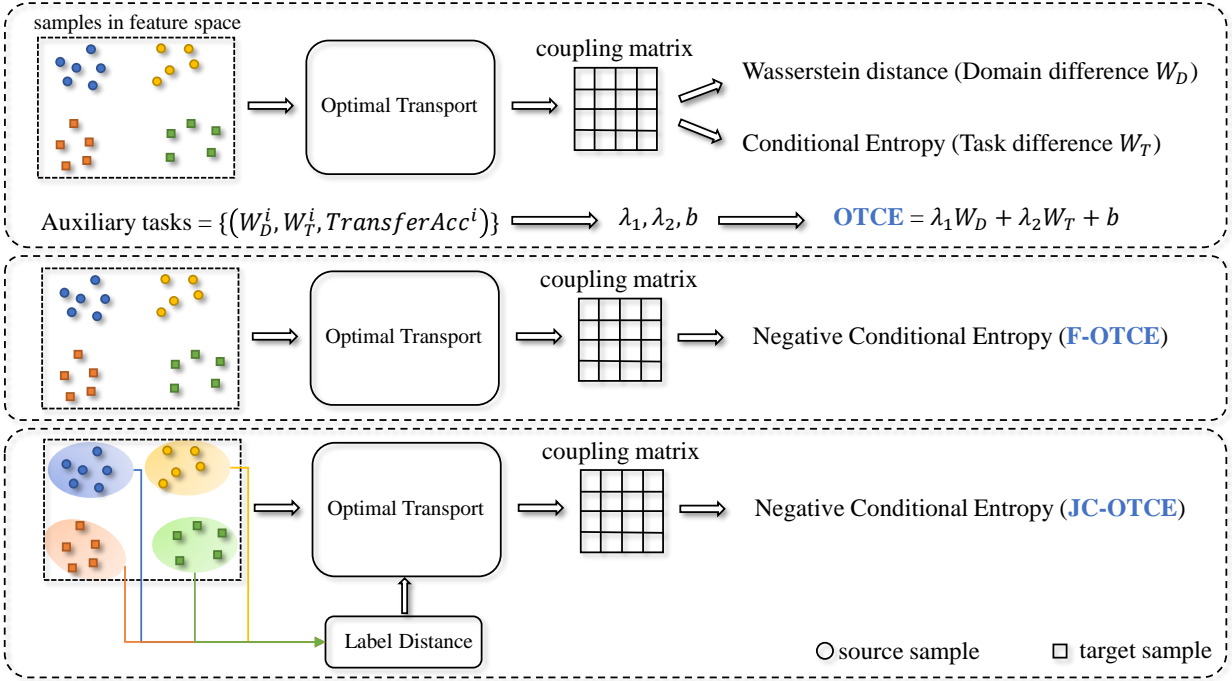
Fig. 2. Illustration of the auxiliary-based OTCE metric [22] (top), and our proposed F-OTCE (middle) and JC-OTCE (bottom) metrics which do not require auxiliary tasks with known transfer accuracy to learn the weighting coefficients. For OTCE (top), $W_D$ and $W_T$ represent the domain difference and task difference between two tasks, respectively. To estimate the coefficients $\lambda_1, \lambda_2, b$ of the linear model, we need to sample at least three auxiliary tasks from the target dataset and calculate $W_D^i$, $W_T^i$ and transfer accuracy $TransferAcc^i$ between the source task and each auxiliary task as training data.

$x_s^i, x_t^i$ from the input space $\mathcal{X}$, and $y_s^i$ from the source label space $\mathcal{Y}_s$, and $y_t^i$ from the target label space $\mathcal{Y}_t$. Meanwhile, $P(x_s) \neq P(x_t)$ and $\mathcal{Y}_s \neq \mathcal{Y}_t$ indicate different domains and tasks respectively. In addition, we are given a source model $(\theta_s, h_s)$ pretrained on source data $D_s$, in which $\theta_s : \mathcal{X} \to \mathbb{R}^d$ represents a feature extractor producing $d$-dimensional features and $h_s : \mathbb{R}^d \to \mathcal{P}(\mathcal{Y}_s)$ is the head classifier predicting the final probability distribution of labels, where $\mathcal{P}(\mathcal{Y}_s)$ is the space of all probability distributions over $\mathcal{Y}_s$. Note that the notation $\theta$ and $h$ can also represent model parameters.

In this paper, we mainly investigate the transferability estimation problem with two representative transfer paradigms for neural networks [5], i.e., *Retrain head* and *Finetune*. The *Retrain head* method keeps the parameters of the source feature extractor $\theta_s$ frozen and retrains a new head classifier $h_t$. But the *Finetune* method updates the source feature extractor and the head classifier simultaneously to obtain new $(\theta_t, h_t)$. Compared to *Retrain head*, *Finetune* trade-offs transfer efficiency for better transfer accuracy and it requires more target data to avoid overfitting [22].

To obtain the empirical transferability, we need to retrain the source model via Retrain head or Finetune on target data and then evaluate the expected log-likelihood on its testing set. Formally, the empirical transferability is defined as:

*Definition 1:* The empirical transferability from the source task $S$ to the target task $T$ is measured by the expected log-likelihood of the retrained $(\theta_s, h_t)$ or $(\theta_t, h_t)$ on the testing

set of target task:

$$\text{Trf}(S \to T) = \begin{cases} \mathbb{E}\left[\log P(y_t|x_t; \theta_s, h_t)\right] & \text{(Retrain head)} \\ \mathbb{E}\left[\log P(y_t|x_t; \theta_t, h_t)\right] & \text{(Finetune)} \end{cases}, \quad (1)$$

which indicates how good the transfer performance is on the target task. In practice, we usually take the testing accuracy as an approximation of the log-likelihood [20], [22].

Although the empirical transferability can be the golden standard of describing how easy it is to transfer the knowledge learned from a source task to a target task, it is computationally expensive to obtain. Efficient transferability metric is a function of the source and target data that approximates the empirical transferability, i.e., the *ground-truth* of the transfer performance on target tasks. It is therefore imperative to find efficient transferability metrics that can accurately estimate empirical transferability.

## III. PRELIMINARY ANALYSIS OF OTCE

OTCE (Optimal Transport based Conditional Entropy) is an analytical transferability metric proposed for the cross-domain cross-task transfer learning setting. As illustrated in Fig. 2 (upper part), OTCE quantifies transferability as a linear model of the domain difference $W_D$ (measured by Wasserstein distance) and task difference $W_T$ (determined by conditional entropy), which is denoted as:

$$\text{OTCE} = \lambda_1 W_D + \lambda_2 W_T + b. \quad (2)$$

However, a major limitation of OTCE is its dependency on auxiliary tasks with known transfer accuracy to determine the
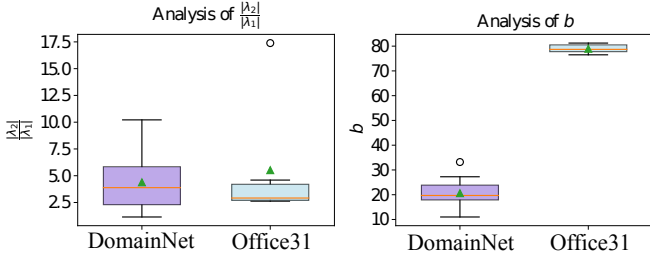
Fig. 3. Statistic of the learned weighting coefficients $\lambda_1, \lambda_2$ and the bias term $b$ of OTCE under diverse transfer configurations.

intrinsic parameters of the linear model. In practice, we are not always able to access sufficient labeled data from target domain for constructing auxiliary tasks. Meanwhile, obtaining the transfer performance of auxiliary tasks needs retraining the source model, which incurs additional computation costs. As a result, the reliance on auxiliary tasks makes OTCE relatively inefficient and less applicable in general scenarios.

The statistic of the learned parameters $\lambda_1, \lambda_2, b$ (shown in Fig. 3) reveals that $\frac{|\lambda_2|}{|\lambda_1|}$ among different transfer configurations varied irregularly, suggesting that the importance of domain difference and task difference varies for different cross-domain transfer learning settings. It is therefore incapable of using the pre-defined coefficients for computing OTCE scores. Additionally, we notice that the task difference $W_T$ plays a more important role ($\frac{|\lambda_2|}{|\lambda_1|} > 1$) in evaluating transferability. Therefore, our proposed auxiliary-free transferability metrics mainly utilize the task difference for describing transferability.

## IV. Auxiliary-free Transferability Metrics

Our proposed auxiliary-free transferability metrics **F-OTCE (Fast OTCE)** and **JC-OTCE (Joint Correspondence OTCE)** can be viewed as the efficient versions of the auxiliary-based OTCE metric, which only consider the Negative Conditional Entropy to describe transferability, as depicted in Fig. 2. Although we do not explicitly evaluate the domain difference, the estimated probabilistic coupling between the source and target data implicitly captures the domain difference to some extent in this unified framework.

Specifically, F-OTCE achieves higher efficiency, while JC-OTCE performs better in terms of accuracy across diverse scenarios. The main difference between the two metrics is that the ground cost of JC-OTCE considers both sample distance and label distance when calculating the optimal coupling between the source and target data, which approximates computing ground cost in the joint space $\mathcal{X} \times \mathcal{Y}$, resulting in more precise data correspondences.

### A. F-OTCE Metric

Formally, we first use the source feature extractor $\theta_s$ to embed the source and target input instances as latent features, denoted as $\hat{x}_s^i = \theta_s(x_s^i)$ and $\hat{x}_t^i = \theta_s(x_t^i)$ respectively. Then the computation process contains two steps as described below.

**Step1: Compute optimal coupling.** First, for the F-OTCE metric, we define the ground cost between samples as:

$$c_1(\hat{x}_s^i, \hat{x}_t^j) \triangleq \|\hat{x}_s^i - \hat{x}_t^j\|_2^2, \qquad (3)$$

so the OT problem with the entropic regularization [28] can be defined as:

$$OT(X_s, X_t) \triangleq \min_{\pi \in \mathcal{P}(X_s, X_t)} \sum_{i,j=1}^{m,n} c_1(\hat{x}_s^i, \hat{x}_t^j)\pi_{ij} - \lambda H(\pi), \quad (4)$$

where $\pi$ is the coupling matrix of size $m \times n$, and $H(\pi) = -\sum_{i=1}^{m} \sum_{j=1}^{n} \pi_{ij} \log \pi_{ij}$ is the entropic regularizer with $\lambda = 0.1$. The OT problem above can be solved efficiently by the Sinkhorn algorithm [28] to produce an optimal coupling matrix $\pi^*$.

From a probabilistic point of view, the coupling matrix $\pi^*$ is a non-parametric estimation of the joint probability distribution of the source and target latent features $P(X_s, X_t)$. We model the relationship between the source and the target data according to the following simple Markov random field: $Y_s - X_s - X_t - Y_t$, where label random variables $Y_s$ and $Y_t$ are only dependent on $X_s$ and $X_t$, respectively, i.e., $P(Y_s, Y_t|X_s, X_t) = P(Y_s|X_s)P(Y_t|X_t)$. Furthermore, we can derive the empirical joint probability distribution of the source and target labels,

$$P(Y_s, Y_t) = \mathbb{E}_{X_s, X_t}[P(Y_s|X_s)P(Y_t|X_t)]. \qquad (5)$$

This joint probability distribution can reveal the transfer performance since the goodness of class-to-class matching intuitively reveals the hardness of transfer.

**Step2: Compute negative conditional entropy.** We are inspired by Tran *et al.* [20] who use Conditional Entropy (CE) $H(Y_t|Y_s)$ to describe class-to-class matching quality over the same input instances. They have shown that the empirical transferability is lower bounded by the negative conditional entropy,

$$\widetilde{\mathrm{Trf}}(S \to T) \geq l_S(\theta_s, h_s) - H(Y_t|Y_s), \qquad (6)$$

where the training log-likelihood $\widetilde{\mathrm{Trf}}(S \to T) = l_T(\theta_s, h_t) = \frac{1}{n} \sum_{i=1}^{n} \log P(y_t^i|x_t^i; \theta_s, h_t)$ is an approximation of the empirical transferability when the retrained model is not overfitted. And $l_S(\theta_s, h_s)$ is a constant, so the empirical transferability can be attributed to the conditional entropy.

We consider it as a reasonable metric to evaluate transferability under the cross-domain cross-task transfer setting once we learn the soft correspondence $\pi^*$ between source and target features via optimal transport. We can also compute the empirical joint probability distribution of the source and target labels, and the marginal probability distribution of the source label, denoted as:

$$\hat{P}(y_s, y_t) = \sum_{i,j:y_s^i=y_s, y_t^j=y_t} \pi_{ij}^*, \qquad (7)$$

$$\hat{P}(y_s) = \sum_{y_t \in \mathcal{Y}_t} \hat{P}(y_s, y_t). \qquad (8)$$

Then we can compute the negative conditional entropy as the F-OTCE score,

$$\begin{aligned} \text{F-OTCE} &= -H_{\pi^*}(Y_t|Y_s) \\ &= \sum_{y_t \in \mathcal{Y}_t} \sum_{y_s \in \mathcal{Y}_s} \hat{P}(y_s, y_t) \log \frac{\hat{P}(y_s, y_t)}{\hat{P}(y_s)}. \end{aligned} \qquad (9)$$
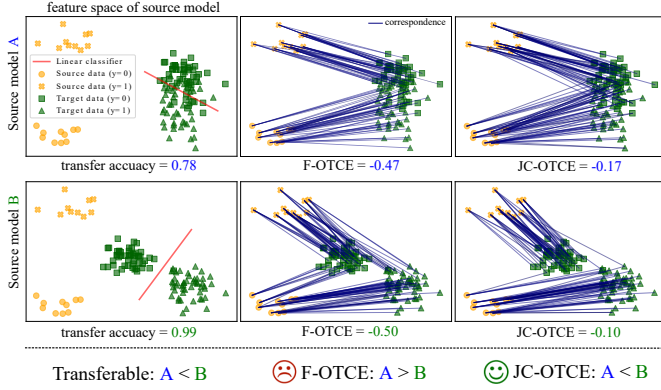
Fig. 4. A toy example shows that the F-OTCE metric fails to distinguish the more transferable source model, while the JC-OTCE predicts correctly by involving the label distance in computing the correspondences.

TABLE I
DIFFERENCES BETWEEN MODEL FINETUNING AND DOMAIN GENERALIZATION.

| | Model finetuning | Domain generalization |
|---|---|---|
| Source data | Single | Multiple |
| Target task | Known | Unknown |
| Goal | Achieving higher accuracy on the target task | Learning generalizable feature representations |

previous Section IV-A (see Equation (7), (8)), the JC-OTCE score is computed as the negative conditional entropy as well.

$$\text{JC-OTCE} = -H_{\pi^*}(Y_t|Y_s)$$
$$= \sum_{y_t \in \mathcal{Y}_t} \sum_{y_s \in \mathcal{Y}_s} \hat{P}(y_s, y_t) \log \frac{\hat{P}(y_s, y_t)}{\hat{P}(y_s)}. \quad (12)$$

## V. TRANSFERABILITY-GUIDED TRANSFER LEARNING

In this section, we present two examples of utilizing our transferability metric to boost the performance of downstream transfer learning tasks, including model finetuning and domain generalization. The differences between these two transfer learning tasks are described in Table I.

To facilitate the training process, we adopt the F-OTCE metric as the optimization objective since using the JC-OTCE metric needs solving multiple OT problems to compute pairwise label distances, which incurs significant computational costs. Additionally, due to GPU memory constraints, we typically perform mini-batch training which only loads a subset of the dataset in the current training iteration, while computing label distance requires loading the entire dataset.

### A. OTCE-based Model Finetuning

The vanilla finetune algorithm follows the "pretraining + finetuning" pipeline that is commonly used in transfer learning. However, this scheme does not consider the relatedness between the source and target tasks. To address this issue, our proposed OTCE-based finetune algorithm introduces an intermediate step into the conventional pipeline, i.e., maximize the transferability of transferring from the source task to the target task, resulting in a "pretraining + adaptation (maximizing transferability) + finetuning" framework. The moderate optimization during the adaptation step utilizes the task relationship characterized by our F-OTCE score to enable the source feature representation to become more transferable to the target task. This facilitates easier learning of the head classifier during the finetuning step and ultimately leads to higher transfer accuracy.

Suppose we have obtained the pretrained model on the source task, the OTCE-based finetune algorithm is a two-step framework, as depicted in Fig. 5 and Algorithm 1. First, we optimize the source feature extractor $\hat{\theta}_s$ by minimizing the conditional entropy within one epoch. Formally,

$$\hat{\theta}_s^* = \arg\min_{\hat{\theta}_s} H_{\pi^*}(Y_t|Y_s)$$
$$= -\arg\min_{\hat{\theta}_s} \sum_{y_t \in \mathcal{Y}_t} \sum_{y_s \in \mathcal{Y}_s} \hat{P}(y_s, y_t) \log \frac{\hat{P}(y_s, y_t)}{\hat{P}(y_s)}, \quad (13)$$

Compared to the auxiliary-based OTCE, we directly use the negative conditional entropy to characterize transferability, which avoids the cumbersome parameter fitting process on auxiliary tasks, resulting in a drastic efficiency improvement.

### B. JC-OTCE Metric

F-OTCE is an efficient transferability metric in practical scenarios, but its accuracy can be further improved. Take a toy example shown in Fig. 4 for illustration, where the F-OTCE metric fails to distinguish the more transferable source model. This observation suggests that computing data correspondences solely based on sample distance (in space $\mathcal{X}$) may not always accurately capture the class-to-class matching quality (or the label uncertainty of the target task) as expected. Therefore, to further improve the accuracy of F-OTCE, we propose the JC-OTCE metric which involves the additional label distance in computing the joint correspondences between data in the joint space $\mathcal{X} \times \mathcal{Y}$.

Formally, we first define the data instances of the source and target tasks as $z_s = (\hat{x}_s, y_s)$ and $z_t = (\hat{x}_t, y_t)$ respectively, where $z_s \in \mathcal{Z}_s = \mathcal{X} \times \mathcal{Y}_s$ and $z_t \in \mathcal{Z}_t = \mathcal{X} \times \mathcal{Y}_t$. And we define the $\alpha_y \triangleq P(X|Y=y)$, which can be estimated from a collection of finite samples with label $y$. Inspired by recent work [29], we compute the label distance as the Wasserstein distance $Wass(\alpha_{y_s}, \alpha_{y_t})$. Then the ground cost for JC-OTCE can be defined as:

$$c_2(z_s^i, z_t^j) \triangleq \gamma\|\hat{x}_s^i - \hat{x}_t^j\|_2^2 + (1-\gamma)Wass(\alpha_{y_s^i}, \alpha_{y_t^j}), \quad (10)$$

where $\gamma \in [0, 1]$ is a weighting coefficient to combine the sample distance and the label distance, and here we let $\gamma = 0.5$. More discussion about $\gamma$ is described in Section VI-C. Similarly, the OT problem for $Z_s$ and $Z_t$ is defined as:

$$OT(Z_s, Z_t) \triangleq \min_{\pi \in \mathcal{P}(Z_s, Z_t)} \sum_{i,j=1}^{m,n} c_2(z_s^i, z_t^j)\pi_{ij} - \lambda H(\pi). \quad (11)$$

By solving this OT problem, we also obtain the optimal coupling matrix $\pi^*$. Then following the **Step2** described in
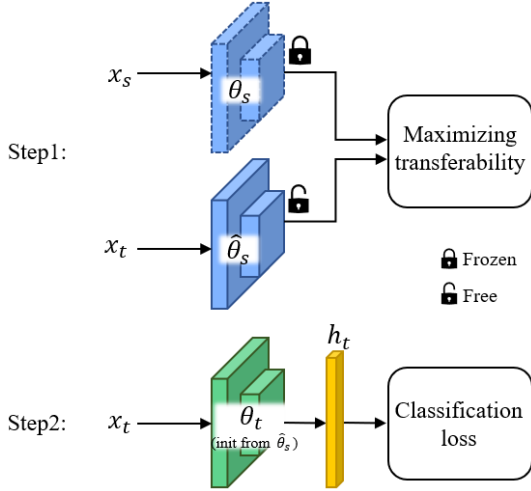
Fig. 5. The pipeline of our OTCE-based finetune algorithm.

---

**Algorithm 1** OTCE-based finetune

**Require:** source dataset $D_s = \{(x_s^i, y_s^i)\}_{i=1}^m$
     target dataset $D_t = \{(x_t^i, y_t^i)\}_{i=1}^n$
     source feature extractor $\theta_s$

1: Initialize $\hat{\theta}_s = \theta_s$
2: **while** sampling mini-batches within one epoch **do**
3:   Generate mini-batch $B_s = \{(\theta_s(x_s^i), y_s^i)\}_i^M$
4:   Generate mini-batch $B_t = \{(\hat{\theta}_s(x_t^i), y_t^i)\}_i^N$
5:   Update $\hat{\theta}_s$ via maximizing F-OTCE$(B_s, B_t)$
6: **end while**
7: Initialize $\theta_t = \hat{\theta}_s$
8: Randomly initialize $h_t$
9: **while** $\theta_t, h_t$ not converge **do**
10:   Update $\theta_t, h_t$ using equation (14)
11: **end while**

---

where the $\pi^*$ is the optimal coupling matrix computed from Equation (4). Joint label distribution $\hat{P}(y_s, y_t)$ and marginal $\hat{P}(y_s)$ are computed from Equation (7), (8). The computation of solving the OT problem with entropic regularizer [28] (Equation (4)) is differentiable [30] since the iterations form a sequence of linear operations, so it can be implemented on the PyTorch framework as a specialized layer[1] of the neural network. After that, we initialize the target feature extractor $\theta_t$ from the optimized source weights $\hat{\theta}_s^*$, and then retrain the target model $(\theta_t, h_t)$ on the target training data using the cross-entropy loss function,

$$\theta_t^*, h_t^* = \arg\max_{\theta_t, h_t} \sum_{i=1}^m \sum_{l=1}^k \mathbf{1}\{y_t^i = l\} \log \frac{\exp(h_t^l(\theta_t(x_t^i)))}{\sum_{j=1}^k \exp(h_t^j(\theta_t(x_t^i)))},$$
(14)

where $m$ represents the number of target training samples, and $k$ is the number of the categories of the target task.

Note that we do not make it a one-step framework, i.e., simultaneously maximize the transferability and minimize the classification loss. Because optimizing two objectives simultaneously may cause gradient conflicts in mini-batch training, which will deteriorate the final classification performance.

*B. OTCE-based Domain Generalization*
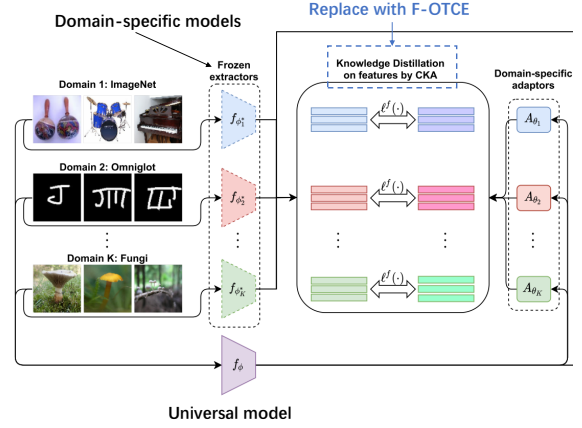
---

[1]https://github.com/dfdazac/wassdistance



Fig. 6. Partial illustration of the URL framework [27]. We replace the CKA similarity with our F-OTCE metric.

In contrast to the model finetuning task, the domain generalization (DG) task aims to learn the generalizable feature representation exhibiting domain-irrelevant and task-irrelevant characteristics from multiple training domains. Therefore, the learned model can also achieve high classification accuracy when transferred to the unseen tasks from unseen domains. We integrate our F-OTCE metric into a state-of-the-art domain generalization method URL [27], [31] as a loss function to illustrate its effectiveness in boosting the DG algorithm.

More specifically, URL learns a universal model via distilling the common knowledge from multiple pretrained domain-specific models corresponding to each training domain. The universal model is required to achieve high classification accuracy in all training domains as well. Once the universal model is obtained, we can use it to extract feature representations for unseen few-shot classification tasks and make predictions via the nearest neighbor classifier (NCC).

In our opinion, the process of distilling knowledge from domain-specific models can be interpreted as maximizing the transferability between the domain-specific models and the universal model. Therefore, we propose to replace the knowledge distillation objective Centered Kernel Alignment (CKA) [32] similarity used in URL with our F-OTCE metric, as illustrated in Fig. 6. Unlike CKA which solely focuses on minimizing feature differences, F-OTCE considers a wider range of task-specific information to minimize the label uncertainty of the universal model. We follow the default configuration of the URL algorithm. Please refer to [27], [31] and the official codebase[2] for more details about the URL algorithm.

*C. Few-shot Classification Task Definition*

We evaluate the effectiveness of our algorithms based on their transfer accuracy on few-shot classification tasks across domains. A few-shot classification task known as *C-way-K-shot* means that the support (training) set $S = \{(x^i, y^i)\}_{i=1}^{k \times C}$ contains $k$ labeled instances from each of the $C$ categories. The query set $Q = \{(x^i, y^i)\}_{i=1}^{q \times C}$ contains $q$ samples per category and serves as the testing set to evaluate the classification accuracy of the model finetuned on the support set.

---
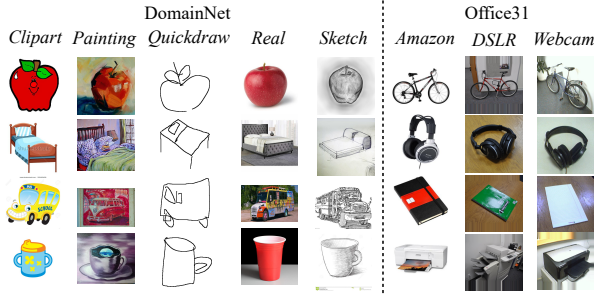
[2]https://github.com/VICO-UoE/URL

Fig. 7. Examples from the cross-domain datasets DomainNet and Office31, where images from different domains exhibit different image styles or are captured by different devices.

## VI. EXPERIMENTS

In this section, we begin by conducting quantitative evaluations of our proposed transferability metrics under various cross-domain cross-task transfer settings. We also explore their applications in source model selection and multi-source feature fusion, as well as provide further analyses on computation efficiency, memory consumption and hyperparameters. Additionally, we conduct extensive evaluations of our proposed transferability-guided transfer learning methods including the OTCE-based finetune algorithm and the OTCE-based URL algorithm.

### A. Evaluation on Transferability Estimation

**Datasets.** Our experiments are conducted on the data from the largest-to-date cross-domain dataset DomainNet [33] and popular Office31 [34] dataset. The DomainNet dataset contains 345-category images in five domains (image styles), i.e., Clipart (C), Painting (P), Quickdraw (Q), Real (R), and Sketch (S), and the Office31 contains 31-category images in three domains including Amazon (A), DSLR (D) and Webcam (W). Data examples are shown in Fig. 7.

**Evaluation criteria.** To quantitatively evaluate the effectiveness of transferability metrics, we adopt the commonly-used Spearman's rank correlation coefficient (Spearman's $\rho$ coefficient) and the Kendall rank correlation coefficient (Kendall's $\tau$ coefficient) [35] to assess the correlation between the transfer accuracy and predicted transferability scores. Specifically, the Spearman's $\rho$ coefficient is defined as:

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}, \tag{15}$$

where $d_i = R(\text{Acc}_i) - R(\text{Trf}_i)$ is the difference between the rankings of transfer accuracy $\text{Acc}_i$ and transferability score $\text{Trf}_i$ for the $i$th source-target task pair, and $n$ represents the total number of task pairs.

The Kendall's $\tau$ coefficient in our experiments is defined as:

$$\tau = \frac{2}{n(n-1)} \sum_{i<j} \text{sgn}(\text{Acc}_i - \text{Acc}_j)\text{sgn}(\text{Trf}_i - \text{Trf}_j). \tag{16}$$

The Kendall's $\tau$ coefficient computes the number of concordant pairs minus the number of discordant pairs divided by

the number of total pairs. A higher rank correlation indicates the more accurate transferability estimation result.

**Transfer settings.** In DomainNet dataset, we successively take each domain as the source domain and use the rest as target domains. For each target domain, we generate 100 target tasks by randomly sampling images in different categories. Then we transfer the source models (ResNet-18 [36]) pretrained on all source domain data to each target task to obtain the ground-truth transfer accuracy. To investigate the performance of transferability metrics under various transfer configurations, three different transfer settings are considered, i.e., the *standard setting*, the *few-shot setting*, and the *fixed category size* setting.

- *Standard setting.* We keep all the training samples of the target task for retraining the source model. Meanwhile, the number of categories of target tasks ranges from 10 to 100. Thus we totally conduct $5 \times 4 \times 100 = 2000$ cross-domain cross-task transfer tests.
- *Few-shot setting.* As transfer learning is commonly used in scenario where only a few labeled data are provided, it is worth evaluating the accuracy of transferability metrics on few-shot cases. The only difference with the *standard setting* is that we limit the target tasks to have only 10 training samples per category.
- *Fixed category size setting.* As studied in [22], the intrinsic complexity of the target task, e.g. category size (number of categories), also affects the transfer accuracy. Usually, a larger category size makes the target task more difficult to learn from limited data. As a result, in the previous two settings, the intrinsic complexity of target tasks with different category sizes may overshadow the more subtle variations in the relatedness with the source task. To investigate whether the transferability metrics are capable of capturing those subtle variations, we propose a more challenging *fixed category size* setting where all target tasks have the same $category\_size = 50$. Other configurations are the same as the *standard setting*.

Moreover, in the Office31 dataset, the DSLR and Webcam domains contain very few samples ($\sim$15 samples per category) and suffer from severe category imbalance. Consequently, we construct two different configurations: *data-imbalanced* and *data-balanced* settings. Both of these two settings are few-shot, but the data-balanced setting permits a maximum of 10 samples per category. Here we only use Amazon as the source domain since the other two domains lack sufficient data to train generalizable source models. It is worth noting that we use the *average per-class accuracy* instead of the *overall accuracy* for representing the transfer performance under the data-imbalanced setting.

For all the settings above, we adopt an SGD optimizer with a learning rate of 0.01 to optimize the cross-entropy loss for 100 epochs during the transfer training phase.

**Results.** Quantitative comparisons with state-of-the-art auxiliary-free transferability metrics including LEEP [21], NCE [20], H-score [7], LogME [23] are shown in Table II, and visual comparisons are illustrated in Fig. 8. Firstly, we can see that both our JC-OTCE and F-OTCE metrics consistently outperform recent LEEP, NCE, H-score, and

TABLE II
QUANTITATIVE COMPARISONS EVALUATED BY SPEARMAN'S $\rho$ COEFFICIENT AND KENDALL'S $\tau$ COEFFICIENT BETWEEN TRANSFERABILITY METRICS AND TRANSFER ACCURACY UNDER DIFFERENT CROSS-DOMAIN CROSS-TASK TRANSFER SETTINGS FOR IMAGE CLASSIFICATION TASKS. OUR PROPOSED JC-OTCE AND F-OTCE METRICS CONSISTENTLY OUTPERFORM STATE-OF-THE-ART AUXILIARY-FREE METRICS. MEANWHILE, THE JC-OTCE ACHIEVES COMPARABLE PERFORMANCE TO THE AUXILIARY-BASED OTCE.

| Setting | Source domain | Target domain | Spearman / Kendall correlation coefficient | | | | | | |
| | | | Auxiliary-based | Auxiliary-free | | | | | |
| | | | OTCE [22] | JC-OTCE | F-OTCE | LEEP [21] | NCE [20] | H-score [7] | LogME [23] |
|---|---|---|---|---|---|---|---|---|---|
| Standard (Retrain head) | C | P,Q,R,S | 0.976 / 0.861 | <u>0.965</u> / <u>0.836</u> | **0.966 / 0.839** | 0.932 / 0.779 | 0.825 / 0.670 | 0.920 / 0.748 | 0.867 / 0.667 |
| | P | C,Q,R,S | 0.977 / 0.868 | **0.966 / 0.837** | <u>0.960</u> / <u>0.822</u> | 0.906 / 0.743 | 0.849 / 0.686 | 0.937 / 0.777 | 0.929 / 0.761 |
| | Q | C,P,R,S | 0.961 / 0.826 | <u>0.962</u> / **0.833** | **0.963** / <u>0.832</u> | 0.953 / 0.810 | 0.943 / 0.793 | 0.942 / 0.784 | 0.912 / 0.744 |
| | R | C,P,Q,S | 0.975 / 0.863 | **0.965 / 0.836** | <u>0.951</u> / <u>0.808</u> | 0.910 / 0.747 | 0.872 / 0.707 | 0.942 / 0.786 | 0.855 / 0.670 |
| | S | C,P,Q,R | 0.969 / 0.842 | <u>0.965</u> / <u>0.834</u> | **0.967 / 0.839** | <u>0.965</u> / <u>0.834</u> | 0.962 / 0.830 | 0.950 / 0.802 | 0.908 / 0.733 |
| Standard (Finetune) | C | P,Q,R,S | 0.932 / 0.766 | **0.900 / 0.713** | 0.884 / 0.689 | 0.814 / 0.618 | 0.664 / 0.517 | 0.889 / <u>0.704</u> | <u>0.890</u> / 0.695 |
| | P | C,Q,R,S | 0.803 / 0.612 | 0.874 / <u>0.698</u> | **0.880** / <u>0.698</u> | 0.850 / 0.655 | 0.797 / 0.613 | <u>0.876</u> / **0.716** | 0.848 / 0.664 |
| | Q | C,P,R,S | 0.896 / 0.719 | **0.906 / 0.732** | 0.895 / <u>0.719</u> | 0.880 / 0.696 | 0.874 / 0.684 | 0.873 / 0.686 | 0.891 / 0.699 |
| | R | C,P,Q,S | 0.912 / 0.732 | **0.905** / <u>0.725</u> | 0.882 / 0.689 | 0.821 / 0.616 | 0.770 / 0.571 | <u>0.902</u> / **0.727** | 0.876 / 0.681 |
| | S | C,P,Q,R | 0.923 / 0.752 | **0.932 / 0.767** | <u>0.929</u> / <u>0.763</u> | 0.927 / <u>0.766</u> | 0.925 / 0.757 | 0.915 / 0.747 | 0.894 / 0.706 |
| | | Average | 0.932 / 0.784 | **0.934** / 0.782 | <u>0.928</u> / <u>0.770</u> | 0.896 / 0.727 | 0.849 / 0.682 | 0.915 / 0.748 | 0.887 / 0.702 |
| Few-shot (Retrain head) | C | P,Q,R,S | 0.926 / 0.756 | **0.926 / 0.757** | <u>0.909</u> / <u>0.729</u> | 0.836 / 0.640 | 0.745 / 0.576 | 0.762 / 0.567 | 0.731 / 0.524 |
| | P | C,Q,R,S | 0.931 / 0.772 | **0.928 / 0.769** | <u>0.886</u> / <u>0.701</u> | 0.803 / 0.618 | 0.746 / 0.575 | 0.811 / 0.608 | 0.849 / 0.649 |
| | Q | C,P,R,S | 0.821 / 0.631 | <u>0.856</u> / <u>0.673</u> | 0.829 / 0.636 | 0.798 / 0.602 | 0.782 / 0.584 | 0.813 / 0.614 | **0.866 / 0.682** |
| | R | C,P,Q,S | 0.929 / 0.769 | **0.897 / 0.724** | <u>0.853</u> / <u>0.666</u> | 0.770 / 0.589 | 0.728 / 0.559 | 0.845 / 0.652 | 0.774 / 0.574 |
| | S | C,P,Q,R | 0.914 / 0.742 | **0.902 / 0.725** | <u>0.895</u> / <u>0.710</u> | 0.872 / 0.680 | 0.872 / 0.679 | 0.838 / 0.645 | 0.867 / 0.684 |
| | | Average | 0.905 / 0.734 | **0.902 / 0.729** | <u>0.875</u> / <u>0.689</u> | 0.815 / 0.625 | 0.775 / 0.595 | 0.814 / 0.618 | 0.818 / 0.623 |
| Fixed category size (Retrain head) | C | P,Q,R,S | 0.701 / 0.500 | **0.695 / 0.498** | <u>0.687</u> / <u>0.487</u> | 0.685 / 0.486 | 0.666 / 0.472 | -0.438 / -0.290 | -0.222 / -0.151 |
| | P | C,Q,R,S | 0.670 / 0.485 | **0.665 / 0.479** | <u>0.631</u> / <u>0.448</u> | 0.630 / 0.446 | 0.612 / 0.430 | -0.529 / -0.371 | -0.043 / -0.039 |
| | Q | C,P,R,S | 0.341 / 0.225 | **0.381 / 0.261** | <u>0.316</u> / <u>0.211</u> | 0.210 / 0.136 | 0.291 / 0.191 | -0.256 / -0.172 | 0.066 / 0.037 |
| | R | C,P,Q,S | 0.637 / 0.455 | **0.695 / 0.498** | <u>0.598</u> / <u>0.415</u> | 0.587 / 0.407 | 0.586 / 0.406 | -0.094 / -0.063 | -0.382 / -0.252 |
| | S | C,P,Q,R | 0.428 / 0.292 | **0.497 / 0.343** | <u>0.436</u> / <u>0.299</u> | 0.404 / 0.277 | 0.432 / 0.298 | -0.247 / -0.164 | 0.027 / 0.006 |
| | | Average | 0.555 / 0.391 | **0.587 / 0.416** | <u>0.534</u> / <u>0.372</u> | 0.503 / 0.350 | 0.517 / 0.359 | -0.313 / -0.212 | -0.111 / -0.080 |
| Imbalanced (Retrain head) | A | D | - / - | **0.844 / 0.646** | <u>0.829</u> / <u>0.627</u> | 0.822 / 0.616 | 0.801 / 0.589 | 0.674 / 0.476 | 0.785 / 0.593 |
| | A | W | - / - | 0.847 / 0.651 | 0.850 / 0.653 | **0.862 / 0.665** | <u>0.859</u> / <u>0.663</u> | 0.657 / 0.489 | 0.787 / 0.590 |
| Balanced (Retrain head) | A | D | - / - | <u>0.822</u> / **0.627** | **0.824** / <u>0.625</u> | 0.796 / 0.592 | 0.783 / 0.572 | 0.574 / 0.393 | 0.747 / 0.536 |
| | A | W | - / - | **0.879 / 0.686** | 0.871 / 0.673 | <u>0.872</u> / <u>0.674</u> | 0.856 / 0.656 | 0.669 / 0.477 | 0.797 / 0.604 |
| | | Average | - / - | **0.848 / 0.653** | <u>0.844</u> / <u>0.645</u> | 0.838 / 0.637 | 0.825 / 0.620 | 0.644 / 0.459 | 0.779 / 0.581 |

**Bold** denotes the best result, and <u>underline</u> denotes the 2nd best result.

LogME metrics on all three transfer settings. In particular, our JC-OTCE metric achieves $(7.3\%, 14.7\%, 4.5\%, 11.4\%)$ and $(16.6\%, 22.5\%, 18.0\%, 17.0\%)$ average gains on Kendall correlation compared to LEEP, NCE, H-score, and LogME respectively under the *standard* setting and the *few-shot* setting. Moreover, the H-score metric and the LogME metric failed under the more challenging *fixed category size* setting, where they showed a negative correlation with the transfer accuracy.

Secondly, the JC-OTCE metric outperforms the F-OTCE metric with an average $5.4\%$ gain on Kendall correlation, which shows that involving the label distance in computing the data correspondences makes the transferability estimation more accurate. Meanwhile, the JC-OTCE metric performs comparably to the original OTCE metric in accuracy, while the former one is evidently more efficient and has fewer restrictions.

Basically, we can conclude that OTCE $\approx$ JC-OTCE $>$ F-OTCE in accuracy and OTCE $<$ JC-OTCE $<$ F-OTCE in efficiency. These three metrics can be applied flexibly according to the needs of different practical situations.

### B. Efficiency Analysis

Given $d$-dimensional extracted features of $m$ source samples and $n$ target samples, assuming that $|\mathcal{Y}_s|, |\mathcal{Y}_t| < \min(m, n)$, the computational complexity of F-OTCE is $O(mn \max\{d, k\})$, where $k$ is the number of Sinkhorn iterations in the OT computation. Specifically, the worst-case complexity of computing the cost matrix between source and target samples is $O(mnd)$. Solving the OT problem by Sinkhorn algorithm with $\epsilon$-accuracy has complexity $O(mnk) = O(2mn\|c\|_\infty^2/(\lambda\epsilon))$ [37] , where $\|c\|_\infty = \sup_{(z_s,z_t)\in\mathcal{Z}^2} c(z_s, z_t)$ is the maximum cost between source and target sample features and $\lambda$ is the weighting coefficient of the entropic regularizer. In practice, we usually set a constant parameter $\lambda$ and a stopping criteria with a maximum iteration. Finally, the conditional entropy computation takes $O(mn)$ time.

Compared to F-OTCE, the additional computation of JC-OTCE lies in computing pair-wise label distances, which needs to solve $|\mathcal{Y}_s| \times |\mathcal{Y}_t|$ OT problems between the samples with given labels and finally produce Wasserstein distances.

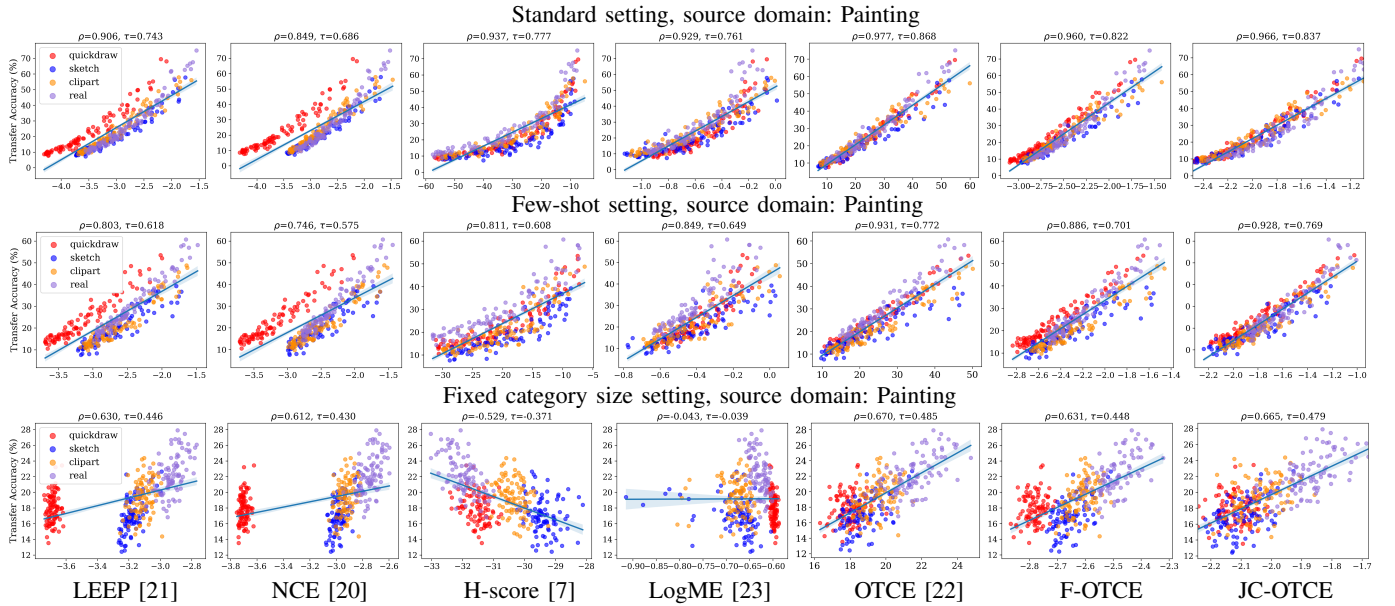The experimental computation time statistics of the empir-

Fig. 8. Visualization of the correlation between the transfer accuracy and transferability metrics, where the vertical axis denotes the transfer accuracy and the horizontal axis represents the transferability scores. Points in the figure represent different target tasks. Our JC-OTCE and F-OTCE metrics show significantly better correlation (more compact) with the transfer accuracy compared to state-of-the-art auxiliary-free metrics, especially under the challenging *fixed category size setting*. Meanwhile, the JC-OTCE metric performs comparably well to the auxiliary-based OTCE metric.

ical transferability and analytical transferability metrics are presented in Table III. Specifically, the empirical transferability calculation is performed on a GPU (NVIDIA GTX1080Ti) by retraining the source model on the target training data and subsequently evaluating its transfer accuracy on the testing set. In order to ensure a fair comparison among analytical transferability metrics, we establish a standardized evaluation setting. For each target classification task, we randomly select 1,000 samples for computation on CPU, which follows the same configuration as in our transferability experiments. To accurately reflect the real wall-clock time required for trans-ferability estimation, we have excluded the time spent on I/O operations. We conducted these computations over 10 random tasks and calculated the average time as the final results. The overall computation time consists of two components: feature extraction time and transferability prediction time. Importantly, the computation time remains consistent across different types of datasets once the source model architecture and the number of samples have been determined. The memory requirements of F-OTCE and JC-OTCE are 395MB and 407MB respectively, which can be easily met by a typical personal computer.

Results show that analytical transferability metrics are $\sim 90\times$ faster than the empirical transferability, without the requirement of GPU. Meanwhile, auxiliary-free metrics perform comparably on efficiency costing $\sim$10s for a pair of tasks. Although the original OTCE metric contains additional auxiliary time since it requires at least three auxiliary tasks with known empirical transferability to determine weighting coefficients, it is still worth using the OTCE metric for more accurate transferability estimation when there are many target tasks under the same cross-domain configuration needing evaluations.

TABLE III
COMPUTATION TIME STATISTICS UNDER THE STANDARD EVALUATION SETTING. AUXILIARY-FREE METRICS ACHIEVE COMPARABLE EFFICIENCY ($\sim$ 10S), WHICH IS EVIDENTLY MORE EFFICIENT THAN THE AUXILIARY-BASED OTCE AND THE EMPIRICAL TRANSFERABILITY.

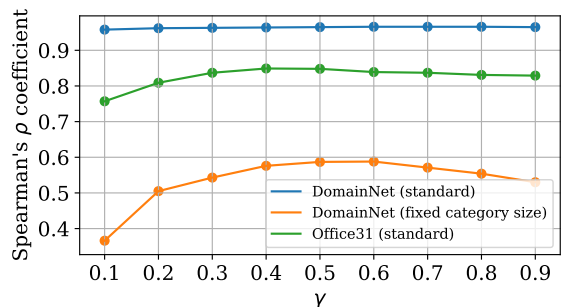| Metric | Auxiliary time | Wall-clock time | Correlation (Spearman) |
|---|---|---|---|
| Empirical transferability | - | 858s (14.3min) | 1.000 |
| LEEP [21] | - | 8.97s | 0.896 |
| NCE [20] | - | 8.92s | 0.849 |
| H-score [7] | - | 9.02s | 0.915 |
| LogME [23] | - | 9.23s | 0.887 |
| OTCE [22] | 858s $\times$ 3 | 9.32s | 0.932 |
| F-OTCE | - | 9.32s | 0.928 |
| JC-OTCE | - | 10.78s | 0.934 |



Fig. 9. We study how the hyper parameter $\gamma$ affects the performance of JC-OTCE and find that let $\gamma = 0.5$ achieving the highest performance.

### C. Effect of Parameter $\gamma$

Here we investigate the effect of the hyperparameter $\gamma$ in JC-OTCE (see Equation (10)), which is a coefficient to balance the impacts of the sample distance and the label distance

TABLE IV
THE TOP-k ACCURACY OF SELECTING THE BEST SOURCE MODEL FROM
15 CANDIDATE MODELS FOR 7 TARGET TASKS RESPECTIVELY ACCORDING
TO TRANSFERABILITY SCORES.

| Method | Top-1 | Top-2 | Top-3 |
|--------|-------|-------|-------|
| JC-NCE | **3 / 7** | **5 / 7** | **6 / 7** |
| F-OTCE | **3 / 7** | 4 / 7 | 5 / 7 |
| H-score [7] | 2 / 7 | 4 / 7 | **6 / 7** |
| LogME [23] | 0 / 7 | 0 / 7 | 0 / 7 |

in computing the ground cost. As shown in Fig. 9, the JC-OTCE metric consistently achieves the highest performance on different transfer settings when $\gamma = 0.5$.

### D. Application in Source Model Selection

One of the most straightforward applications of transferability metrics is choosing the optimal pretrained source model from a set of candidate models for a target task. In this experiment, we evaluate the effectiveness of transferability metrics in source model selection using the Visual Task Adaptation Benchmark (VTAB) [38]. More specifically, the model zoo contains 15 models trained on ImageNet [60] by various algorithms, e.g., supervised learning (Sup-100%) , semi-supervised learning (Semi-rotation-10% and Semi-exemplar-10% [39]), self-supervised learning (Rotation [40] and Jigsaw [41]), generative method (Cond-biggan [42]) and VAEs [43], etc. Meanwhile, VTAB provides the transfer accuracy of these models on 7 target datasets including Caltech101 [44], CIFAR-100 [45], DTD [46], Flowers102 [47], Pets [48], SVHN [49] and Camelyon [50]. In this setting, transferability metrics must identify the best source model with the highest transfer accuracy from 15 candidates. To evaluate the performance of transferability metrics in this task, we compute the Top-k (k=1,2,3) model selection accuracy, as shown in Table IV. We observe that JC-OTCE, F-OTCE and H-score perform well in the model selection task, with JC-OTCE achieving the best accuracy. Meanwhile, in most cases, we notice that the best source model can be chosen from the predicted Top-3 highest transferable models. Note that LEEP and NCE are infeasible here since VTAB only releases the feature extractors of models.

### E. Application in Multi-source Feature Fusion

When multiple source models are accessible, one can transfer them to a target task by merging their inferred features to obtain a fused feature representation [51]. However, different source models may yield different transfer performances on the target task. Therefore, simple average fusion may not effectively leverage the most useful information provided by source models. As a result, we utilize transferability scores to weigh feature fusion and improve its transfer accuracy.

We adopt the same experimental setting as proposed in [22], and employ the JC-OTCE and F-OTCE scores (normalized by a softmax function) as weighting coefficients to fuse four source models trained on different domains. Specifically, we randomly sample 50 few-shot classification tasks from the
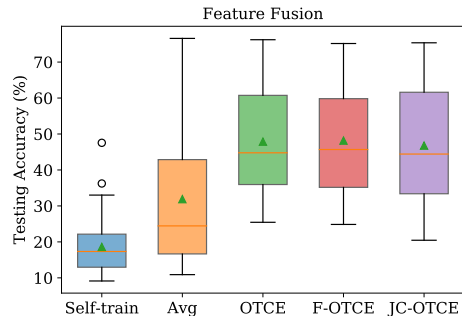


Fig. 10. Testing accuracy comparisons among "Self-train" (directly training on target data), "Avg" (average fusion) and the feature fusion weighted by "OTCE", "F-OTCE", "JC-OTCE", respectively.

TABLE V
DATASETS USED FOR EVALUATING OTCE-BASED FINETUNE ALGORITHM.

| Type | Dataset | Categories | Training samples | Content |
|------|---------|-----------|------------------|---------|
| Source | MNIST | 10 | 60,000 | handwritten digits |
| | Caltech101 | 101 | 9,146 | natural image |
| Target | Omniglot | 1,623 | 32,460 | handwritten character |
| | MiniImageNet | 100 | 60,000 | natural image |

*Real* domain of the DomainNet dataset as target tasks. We multiply the features generated by multiple source models with the corresponding transferability-weighted coefficients and concatenate them. Then a new head classifier is trained on the fused feature representation to produce the final predictions. The results presented in Fig. 10 indicate that JC-OTCE and F-OTCE perform comparably to OTCE, while significantly outperforming the average fusion as expected, in the multi-source feature fusion task.

### F. Evaluation on OTCE-based Finetune Algorithm

One significant use of transfer learning is to address the few-shot classification problem, in which the target task has limited labeled training data, such as 1-shot or 5-shot scenarios. Earlier few-shot learning approaches [52]–[55] driven by meta learning only show their effectiveness in the intra-domain generalization, where the tasks used for meta-training and meta-testing are drawn from the same data distribution. Therefore, we concentrate on the more challenging few-shot classification problem across domains and tasks.

**Task generation**. Specifically, we generate few-shot target tasks using the character recognition dataset Omniglot [56] and the natural image classification dataset MiniImageNet [53], which are commonly-used benchmarks in few-shot learning. And we generate their respective source tasks using MNIST [57] and Caltech101 [44] datasets. The details of datasets are introduced in Table V, and data examples are visualized in Fig. 11. We randomly generate 100 few-shot image classification tasks (5-way-5-shot) from each target dataset respectively.

**Implementation details.** We train the source model (for transfer learning approaches) or apply meta-training (for meta learning approaches) on the source dataset, and then finetune
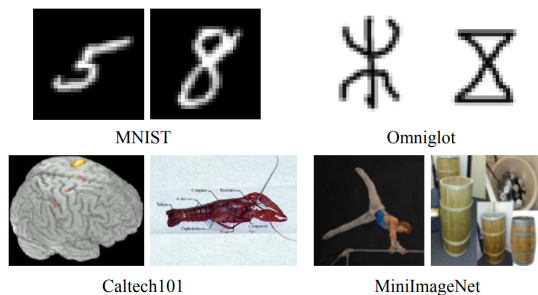
Fig. 11. Examples from the datasets used for model finetuning evaluations.

TABLE VI
MODEL ARCHITECTURE OF THE CONV4 NEURAL NETWORK.

| Layer name | Parameter |
|---|---|
| conv1 | $3 \times 3$ conv, 64 filters, batch norm, ReLU, $2 \times 2$ maxpooling. |
| conv2 | $3 \times 3$ conv, 64 filters, batch norm, ReLU, $2 \times 2$ maxpooling. |
| conv3 | $3 \times 3$ conv, 64 filters, batch norm, ReLU, $2 \times 2$ maxpooling. |
| conv4 | $3 \times 3$ conv, 64 filters, batch norm, ReLU, $2 \times 2$ maxpooling. |
| fc1 | fully connected layer, feature dim $\times$ k. |

the pretrained model or apply meta-test on the target task. To make a fair comparison with previous few-shot learning methods, we first evaluate performances using the widely-used Conv4 [52]–[55] architecture, which comprises four convolutional layers and one fully connected layer, as described in Table VI. Besides, we further examine the effectiveness of our OTCE-based finetune algorithm with different model architectures including the famous LeNet [58] for character recognition (MNIST→Omniglot) and the ResNet-18 [36] for natural image classification (Caltech101→MiniImageNet).

During the training phase of the OTCE-based finetune algorithm, we first optimize the source feature extractor over the source and target datasets for one epoch with the source batch size 256 and the target batch size 25. We use an Adam optimizer with learning rate of 0.0001. Then we initialize the target model with the optimized source weights and continue the finetuning on the target training set for 300 epochs, using the same optimizer supervised by the classification loss function. We adopt the same finetuning strategy for the vanilla finetune algorithm.

**Results.** Table VII demonstrates that our proposed OTCE-based finetune method consistently improves the transfer accuracy of the vanilla finetuning under all transfer settings, with up to $4.41\%$ classification accuracy gain. On the other hand, the vanilla finetune method even outperforms existing representative few-shot learning approaches including ProtoNet [54], MatchingNet [53], MAML [52] and Relation-Net [55] under the cross-domain setting, suggesting that current meta learning methods need further enhancements to improve their cross-domain generalization performance. Moreover, we analyze the running time of our OTCE-based finetune algorithm, and find that optimizing the F-OTCE score only accounts for 16% of the total training time (628s) for a task pair under the Caltech101→MiniImageNet setting using ResNet-18 model.

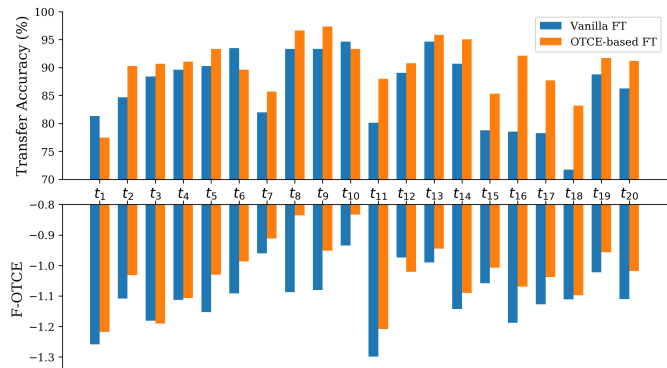Fig. 12 also visually demonstrates that our OTCE-based



Fig. 12. We randomly take 20 target tasks from the *MNIST to Omniglot (LeNet)* setting to visually compare the original source model (in blue) with the optimized source model (in orange) by our F-OTCE targeting to their transfer accuracy and transferability scores. It can be seen that the optimized model shows both higher transfer accuracy and higher F-OTCE score for most cases.

finetune method indeed improves both the transferability score and the transfer accuracy of the source model for most target tasks as expected. In summary, for a given target task, our OTCE-based finetune algorithm provides a straightforward yet effective approach to enhance the transferability of the source model, ultimately leading to more accurate classification results.

### G. Evaluation on OTCE-based URL Algorithm

**Dataset.** We conduct the evaluations of the OTCE-based URL algorithm on the popular Meta-Dataset [59] benchmark, which contains 8 training domains (datasets) such as ImageNet [60], Omniglot [56], Aircraft [61], etc, and 5 testing domains including TrafficSign [62], MS-COCO [63], MNIST [57], CIFAR-10 [45] and CIFAR-100 [45]. The universal model learned on training domains will be evaluated on the few-shot classification tasks randomly sampled from testing domains.

**Implementation details.** We follow the default configuration of the URL algorithm, which employs ResNet-18 [36] as the backbone for the universal model and domain-specific models. The universal model shares its feature extractor but not the head classifiers across domains. We optimize the F-OTCE objective using an SGD optimizer with a learning rate of 0.0001 and momentum of 0.9.

**Results.** Quantitative comparisons presented in Table VIII demonstrate that the generalization performances of URL incorporating with our F-OTCE metric are evidently improved, with up to 2.34% classification accuracy gain.

### VII. CONCLUSION

Transferability estimation under the cross-domain cross-task transfer setting is a practical and challenging problem in transfer learning. Our proposed transferability-guided transfer learning framework not only provides two accurate and efficient auxiliary-free transferability metrics, F-OTCE and JC-OTCE, without the need of retraining the source model, but also offers an useful objective function (F-OTCE) to enhance

TABLE VII
TESTING ACCURACY (%) OF THE CROSS-DOMAIN CROSS-TASK FEW-SHOT CLASSIFICATION EXPERIMENTS, AVERAGED OVER 100 TARGET TASKS AND WITH 95% CONFIDENCE INTERVALS.

| Model | Method | MNIST → Omniglot | Caltech101 → MiniImageNet |
|---|---|---|---|
| Conv4 | MAML [52] | $88.60 \pm 1.14$ | $28.23 \pm 0.44$ |
| | MatchingNet [53] | $87.92 \pm 1.10$ | $44.75 \pm 1.30$ |
| | ProtoNet [54] | $83.11 \pm 1.34$ | $50.40 \pm 1.35$ |
| | RelationNet [55] | $69.35 \pm 1.62$ | $29.55 \pm 0.61$ |
| | Vanilla finetune | $91.30 \pm 0.95$ | $49.49 \pm 1.27$ |
| | OTCE-based finetune | $\mathbf{92.32 \pm 0.87}$ | $\mathbf{51.36 \pm 1.33}$ |
| LeNet | Vanilla finetune | $86.11 \pm 1.10$ | - |
| | OTCE-based finetune | $\mathbf{90.52 \pm 0.94}$ | - |
| ResNet-18 | Vanilla finetune | - | $48.48 \pm 1.39$ |
| | OTCE-based finetune | - | $\mathbf{50.02 \pm 1.34}$ |

TABLE VIII
CLASSIFICATION ACCURACY (%) OF PRETRAINED URL MODEL GENERALIZING TO 100 FEW-SHOT TASKS (5-WAY-5-SHOT) FROM EACH UNSEEN DOMAIN RESPECTIVELY.

| Method | Traffic Sign | MS-COCO | MNIST | CIFAR-10 | CIFAR-100 |
|---|---|---|---|---|---|
| URL | $77.24 \pm 1.88$ | $\mathbf{72.36 \pm 1.90}$ | $90.84 \pm 1.03$ | $67.72 \pm 1.48$ | $78.30 \pm 1.76$ |
| OTCE-based URL | $\mathbf{79.58 \pm 1.68}$ | $71.66 \pm 1.56$ | $\mathbf{91.54 \pm 0.96}$ | $\mathbf{68.40 \pm 1.80}$ | $\mathbf{79.54 \pm 1.72}$ |

the generalizability of the source model, ultimately leading to higher transfer accuracy in downstream model finetuning and domain generalization tasks. F-OTCE is computed as the negative conditional entropy between the source and target labels when given the optimal coupling produced by Optimal Transport. The conditional entropy measures the predicted label uncertainty of the target task under the given pretrained source model and source data, which is negatively correlated with the ground-truth transfer accuracy such that it can serve as an accurate indicator of transferability. Furthermore, JC-OTCE includes the additional label distance in building more accurate data correspondences, which trade-offs a minor efficiency drop for more accurate transferability estimation under diverse transfer configurations.

Our proposed F-OTCE and JC-OTCE metrics drastically reduce the computation time of the auxiliary-based OTCE from 43 minutes to 9.32s and 10.78s respectively, while consistently showing higher accuracy in predicting ground-truth transfer performance than state-of-the-art auxiliary-free metrics, achieving average correlation gains of 21.1% and 25.8% respectively. In particular, JC-OTCE performs comparably to the original OTCE in transferability accuracy, with greater flexibility and efficiency. Additionally, our transferability-guided transfer learning algorithms improves the transfer accuracy of the source model with up to 4.41% and 2.34% gains in few-shot classification tasks. We believe that our OTCE framework can inspire various downstream tasks in transfer learning, multi-task learning and other related applications.

## VIII. ACKNOWLEDGEMENT

## REFERENCES

[1] L. Y. Pratt, "Discriminability-based transfer between neural networks," in *Advances in neural information processing systems*, 1993, pp. 204–211.

[2] Q. Sun, Y. Liu, T.-S. Chua, and B. Schiele, "Meta-transfer learning for few-shot learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 403–412.

[3] L. Shao, F. Zhu, and X. Li, "Transfer learning for visual categorization: A survey," *IEEE transactions on neural networks and learning systems*, vol. 26, no. 5, pp. 1019–1034, 2014.

[4] L. Zhang and X. Gao, "Transfer adaptation learning: A decade survey," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[5] W. Zhang, L. Deng, and D. Wu, "Overcoming negative transfer: A survey," *arXiv preprint arXiv:2009.00909*, 2020.

[6] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on knowledge and data engineering*, vol. 22, no. 10, pp. 1345–1359, 2009.

[7] Y. Bao, Y. Li, S.-L. Huang, L. Zhang, L. Zheng, A. Zamir, and L. Guibas, "An information-theoretic approach to transferability in task transfer learning," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 2309–2313.

[8] A. Maurer, "Transfer bounds for linear feature learning," *Machine learning*, vol. 75, no. 3, pp. 327–350, 2009.

[9] S. Ben-David and R. Schuller, "Exploiting task relatedness for multiple task learning," in *Learning Theory and Kernel Machines*. Springer, 2003, pp. 567–580.

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

[11] S. Ben-David, J. Blitzer, K. Crammer, and F. Pereira, "Analysis of representations for domain adaptation," *Advances in neural information processing systems*, vol. 19, pp. 137–144, 2006.

[12] J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. Wortman, "Learning bounds for domain adaptation," in *Advances in neural information processing systems*, 2008, pp. 129–136.

[13] Y. Mansour, M. Mohri, and A. Rostamizadeh, "Domain adaptation: Learning bounds and algorithms," *arXiv preprint arXiv:0902.3430*, 2009.

[14] W. Wang, H. Li, Z. Ding, F. Nie, J. Chen, X. Dong, and Z. Wang, "Rethinking maximum mean discrepancy for visual domain adaptation," *IEEE Transactions on Neural Networks and Learning Systems*, 2021.

[15] Z. Wang, B. Du, and Y. Guo, "Domain adaptation with neural embedding matching," *IEEE transactions on neural networks and learning systems*, vol. 31, no. 7, pp. 2387–2397, 2019.

[16] A. R. Zamir, A. Sax, W. Shen, L. J. Guibas, J. Malik, and S. Savarese, "Taskonomy: Disentangling task transfer learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 3712–3722.

[17] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?" in *Advances in neural information processing systems*, 2014, pp. 3320–3328.

[18] A. Achille, M. Lam, R. Tewari, A. Ravichandran, S. Maji, C. C. Fowlkes, S. Soatto, and P. Perona, "Task2vec: Task embedding for meta-learning," in *Proceedings of the IEEE international conference on computer vision*, 2019, pp. 6430–6439.

[19] W. Ying, Y. Zhang, J. Huang, and Q. Yang, "Transfer learning via learning to transfer," in *International Conference on Machine Learning*, 2018, pp. 5085–5094.

[20] A. T. Tran, C. V. Nguyen, and T. Hassner, "Transferability and hardness of supervised classification tasks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1395–1405.

[21] C. V. Nguyen, T. Hassner, C. Archambeau, and M. Seeger, "Leep: A new measure to evaluate transferability of learned representations," in *International Conference on Machine Learning*, 2020.

[22] Y. Tan, Y. Li, and S.-L. Huang, "Otce: A transferability metric for cross-domain cross-task representations," in *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*, June 2021, pp. 15 779–15 788.

[23] K. You, Y. Liu, J. Wang, and M. Long, "Logme: Practical assessment of pre-trained models for transfer learning," in *International Conference on Machine Learning*. PMLR, 2021, pp. 12 133–12 143.

[24] J. Huang, N. Xiao, and L. Zhang, "Balancing transferability and discriminability for unsupervised domain adaptation," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.

[25] L. Kantorovich, "On the translocation of masses, cr (dokl.) acad," *Sci. URSS (NS)*, vol. 37, p. 199, 1942.

[26] G. Peyré, M. Cuturi *et al.*, "Computational optimal transport: With applications to data science," *Foundations and Trends® in Machine Learning*, vol. 11, no. 5-6, pp. 355–607, 2019.

[27] W.-H. Li, X. Liu, and H. Bilen, "Universal representation learning from multiple domains for few-shot classification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9526–9535.

[28] M. Cuturi, "Sinkhorn distances: Lightspeed computation of optimal transport," in *Advances in neural information processing systems*, 2013, pp. 2292–2300.

[29] D. Alvarez-Melis and N. Fusi, "Geometric dataset distances via optimal transport," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 21 428–21 439.

[30] A. Genevay, G. Peyré, and M. Cuturi, "Learning generative models with sinkhorn divergences," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2018, pp. 1608–1617.

[31] W.-H. Li, X. Liu, and H. Bilen, "Universal representations: A unified look at multiple task and domain learning," *arXiv preprint arXiv:2204.02744*, 2022.

[32] S. Kornblith, M. Norouzi, H. Lee, and G. Hinton, "Similarity of neural network representations revisited," in *International Conference on Machine Learning*. PMLR, 2019, pp. 3519–3529.

[33] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang, "Moment matching for multi-source domain adaptation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1406–1415.

[34] K. Saenko, B. Kulis, M. Fritz, and T. Darrell, "Adapting visual category models to new domains," in *European conference on computer vision*. Springer, 2010, pp. 213–226.

[35] M. G. Kendall, "A new measure of rank correlation," *Biometrika*, vol. 30, no. 1/2, pp. 81–93, 1938.

[36] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[37] L. Chizat, P. Roussillon, F. Léger, F.-X. Vialard, and G. Peyré, "Faster wasserstein distance estimation with the sinkhorn divergence," in *Neural Information Processing Systems*, 2020.

[38] X. Zhai, J. Puigcerver, A. Kolesnikov, P. Ruyssen, C. Riquelme, M. Lucic, J. Djolonga, A. S. Pinto, M. Neumann, A. Dosovitskiy *et al.*, "A large-scale study of representation learning with the visual task adaptation benchmark," *arXiv preprint arXiv:1910.04867*, 2019.

[39] X. Zhai, A. Oliver, A. Kolesnikov, and L. Beyer, "S4l: Self-supervised semi-supervised learning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1476–1485.

[40] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in *International Conference on Learning Representations*, 2018.

[41] M. Noroozi and P. Favaro, "Unsupervised learning of visual representations by solving jigsaw puzzles," in *European Conference on Computer Vision*, 2016.

[42] A. Brock, J. Donahue, and K. Simonyan, "Large scale gan training for high fidelity natural image synthesis," *International Conference on Learning Representations*, 2019.

[43] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.

[44] L. Fei-Fei, R. Fergus, and P. Perona, "One-shot learning of object categories," *IEEE transactions on pattern analysis and machine intelligence*, vol. 28, no. 4, pp. 594–611, 2006.

[45] A. Krizhevsky, "Learning multiple layers of features from tiny images," *Master's thesis, University of Tront*, 2009.

[46] M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi, "Describing textures in the wild," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2014.

[47] M.-E. Nilsback and A. Zisserman, "Automated flower classification over a large number of classes," in *Indian Conference on Computer Vision, Graphics and Image Processing*, Dec 2008.

[48] O. M. Parkhi, A. Vedaldi, A. Zisserman, and C. V. Jawahar, "Cats and dogs," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2012.

[49] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," in *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.

[50] B. S. Veeling, J. Linmans, J. Winkens, T. Cohen, and M. Welling, "Rotation equivariant cnns for digital pathology," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2018.

[51] S. Hou, X. Liu, and Z. Wang, "Dualnet: Learn complementary features for image recognition," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 502–510.

[52] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *International Conference on Machine Learning*. PMLR, 2017, pp. 1126–1135.

[53] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra *et al.*, "Matching networks for one shot learning," *Advances in neural information processing systems*, vol. 29, pp. 3630–3638, 2016.

[54] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," *Advances in Neural Information Processing Systems*, vol. 30, pp. 4077–4087, 2017.

[55] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. Torr, and T. M. Hospedales, "Learning to compare: Relation network for few-shot learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1199–1208.

[56] B. Lake, R. Salakhutdinov, J. Gross, and J. Tenenbaum, "One shot learning of simple visual concepts," in *Proceedings of the annual meeting of the cognitive science society*, vol. 33, no. 33, 2011.

[57] Y. LeCun, C. Cortes, and C. Burges, "Mnist handwritten digit database," *ATT Labs [Online]. Available: http://yann.lecun.com/exdb/mnist*, vol. 2, 2010.

[58] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[59] E. Triantafillou, T. Zhu, V. Dumoulin, P. Lamblin, U. Evci, K. Xu, R. Goroshin, C. Gelada, K. Swersky, P.-A. Manzagol *et al.*, "Meta-dataset: A dataset of datasets for learning to learn from few examples," *arXiv preprint arXiv:1903.03096*, 2019.

[60] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.

[61] S. Maji, E. Rahtu, J. Kannala, M. Blaschko, and A. Vedaldi, "Fine-grained visual classification of aircraft," *arXiv preprint arXiv:1306.5151*, 2013.

[62] S. Houben, J. Stallkamp, J. Salmen, M. Schlipsing, and C. Igel, "Detection of traffic signs in real-world images: The german traffic sign detection benchmark," in *The 2013 international joint conference on neural networks (IJCNN)*. Ieee, 2013, pp. 1–8.

[63] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Computer Vision–ECCV 2014: 13th European Conference,*

*Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*. Springer, 2014, pp. 740–755.