

---

# MOST: Multi-Source Domain Adaptation via Optimal Transport for Student-Teacher Learning (Supplementary Material)

---

Tuan Nguyen<sup>1</sup>    Trung Le<sup>1</sup>    He Zhao<sup>1</sup>    Quan Hung Tran<sup>2</sup>    Truyen Nguyen<sup>3</sup>    Dinh Phung<sup>1,4</sup>

<sup>1</sup>Department of Data Science and AI, Monash University, Australia

<sup>2</sup>Adobe Research, San Jose, CA, USA

<sup>3</sup>University of Akron, USA

<sup>4</sup>VinAI Research, Vietnam

This is the supplementary material for “MOST: Multi-Source Domain Adaptation via Optimal Transport for Student-Teacher Learning”. Our supplementary material is organized into four main sections. In the first and second sections, we provide complete proofs for our theory. The third section is dedicated to the clustering view of optimal transport with a link to our imitation learning view, whilst the detail of network architecture, implementation specification, and further ablation studies are presented in the fourth section.

## 1 THEORETICAL DEVELOPMENTS

In what follows, we present all proofs of our theory. Given a pair of data distribution  $\mathbb{P}^A$  and labeling function  $h^A$ , we define a distribution  $\mathbb{P}_{A,h^A}$  over  $\mathcal{X}^A \times \mathcal{Y}_\Delta$  including sample pair  $(\mathbf{x}, h^A(\mathbf{x}))$  by first sampling  $\mathbf{x} \sim \mathbb{P}^A$  and then computing  $h^A(\mathbf{x})$ . Similarly, we can define another distribution  $\mathbb{P}_{B,h^B}$  over  $\mathcal{X}^B \times \mathcal{Y}_\Delta$  using the data distribution  $\mathbb{P}^B$  and the labeling function  $h^B$ . We now investigate the Wasserstein distance between  $\mathbb{P}_{A,h^A}$  and  $\mathbb{P}_{B,h^B}$  w.r.t the cost (metric) function  $d = \lambda d_{\mathcal{X}} + d_{\mathcal{Y}}$ . Proposition 1 is crucial for us to develop an imitation learning view based on optimal transport.

**Proposition 1.** *The WS distance of interest can be expressed as:*

$$\begin{aligned} \mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{B,h^B}) &= \min_{L: L_{\#}\mathbb{P}^A = \mathbb{P}^B} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}^A} [\lambda d_{\mathcal{X}}(\mathbf{x}, L(\mathbf{x})) + d_{\mathcal{Y}}(h^A(\mathbf{x}), h^B(L(\mathbf{x})))] \\ &= \min_{H: H_{\#}\mathbb{P}^B = \mathbb{P}^A} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}^B} [\lambda d_{\mathcal{X}}(\mathbf{x}, H(\mathbf{x})) + d_{\mathcal{Y}}(h^B(\mathbf{x}), h^A(H(\mathbf{x})))]. \end{aligned}$$

*Proof.* Observe first that for any  $U_A \subset \mathcal{X}^A \times \mathcal{Y}_\Delta$ , we have  $\mathbb{P}_{A,h^A}(U_A) = \mathbb{P}^A(V_A)$  where  $V_A := \{\mathbf{x} \in \mathcal{X}^A \mid (\mathbf{x}, h^A(\mathbf{x})) \in U_A\}$ . Similarly, we have for any  $U_B \subset \mathcal{X}^B \times \mathcal{Y}_\Delta$  that  $\mathbb{P}_{B,h^B}(U_B) = \mathbb{P}^B(V_B)$  where  $V_B := \{\mathbf{x} \in \mathcal{X}^B \mid (\mathbf{x}, h^B(\mathbf{x})) \in U_B\}$ .

Let  $K: \text{supp}(\mathbb{P}_{A,h^A}) \rightarrow \text{supp}(\mathbb{P}_{B,h^B})$  (i.e.,  $\text{supp}$  indicates the support of a distribution) be such that  $K_{\#}\mathbb{P}_{A,h^A} = \mathbb{P}_{B,h^B}$ . We can express  $K$  as

$$K(\mathbf{x}, h^A(\mathbf{x})) := (K_1(\mathbf{x}, h^A(\mathbf{x})), K_2(\mathbf{x}, h^A(\mathbf{x}))),$$

with  $K_1(\mathbf{x}, h^A(\mathbf{x})) \in \mathcal{X}^B$  and  $K_2(\mathbf{x}, h^A(\mathbf{x})) \in \mathcal{Y}_\Delta$ . Define  $L(\mathbf{x}) := K_1(\mathbf{x}, h^A(\mathbf{x}))$ . We claim that  $L_{\#}\mathbb{P}^A = \mathbb{P}^B$ . Indeed, let  $V_B \subset \mathcal{X}^B$  be any measurable set and take  $U_B := V_B \times \mathcal{Y}_\Delta$ . Then by using the observation above and the fact  $K_{\#}\mathbb{P}_{A,h^A} = \mathbb{P}_{B,h^B}$ , we obtain

$$\mathbb{P}^B(V_B) = \mathbb{P}_{B,h^B}(U_B) = \mathbb{P}_{A,h^A}(K^{-1}(U_B)) = \mathbb{P}_{A,h^A}(L^{-1}(V_B) \times \mathcal{Y}_\Delta) = \mathbb{P}^A(L^{-1}(V_B)).$$

Thus the claim is proved.

It also follows from  $K_{\#}\mathbb{P}_{A,h^A} = \mathbb{P}_{B,h^B}$  and the claim that  $K_2(\mathbf{x}, h^A(\mathbf{x})) = h_B(L(\mathbf{x}))$ , which gives

$$d((\mathbf{x}, h^A(\mathbf{x})), K(\mathbf{x}, h^A(\mathbf{x}))) = \lambda d_{\mathcal{X}}(\mathbf{x}, L(\mathbf{x})) + d_{\mathcal{Y}}(h^A(\mathbf{x}), h^B(L(\mathbf{x}))). \quad (1)$$

Therefore, we deduce that

$$\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{B,h^B}) \geq \min_{L: L_{\#}\mathbb{P}^A = \mathbb{P}^B} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}^A} [\lambda d_{\mathcal{X}}(\mathbf{x}, L(\mathbf{x})) + d_{\mathcal{Y}}(h^A(\mathbf{x}), h^B(L(\mathbf{x})))].$$

In order to prove the reverse inequality, let us consider any map  $L$  satisfying  $L_{\#}\mathbb{P}^A = \mathbb{P}^B$ . Define  $K(\mathbf{x}, h^A(\mathbf{x})) := (L(\mathbf{x}), h^B(L(\mathbf{x})))$ . Then (1) holds and  $K_{\#}\mathbb{P}_{A,h^A} = \mathbb{P}_{B,h^B}$ . To verify the latter, let  $U_B \subset \mathcal{X}^B \times \mathcal{Y}^B$  be any measurable set and take  $V_B := \{\mathbf{x} \in \mathcal{X}^B \mid (\mathbf{x}, h^B(\mathbf{x})) \in U_B\}$ . Then as

$$K^{-1}(U_B) = \{(\mathbf{x}, h^A(\mathbf{x})) \mid L(\mathbf{x}) \in V_B\} = \{(\mathbf{x}, h^A(\mathbf{x})) \mid \mathbf{x} \in L^{-1}(V_B)\},$$

we have

$$\mathbb{P}_{A,h^A}(K^{-1}(U_B)) = \mathbb{P}^A(L^{-1}(V_B)) = \mathbb{P}^B(V_B) = \mathbb{P}_{B,h^B}(U_B).$$

Thus it follows that

$$\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{B,h^B}) \leq \min_{L: L_{\#}\mathbb{P}^A = \mathbb{P}^B} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}^A} [\lambda d_{\mathcal{X}}(\mathbf{x}, L(\mathbf{x})) + d_{\mathcal{Y}}(h^A(\mathbf{x}), h^B(L(\mathbf{x})))].$$

By combining the above two inequalities, we obtain the equality

$$\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{B,h^B}) = \min_{L: L_{\#}\mathbb{P}^A = \mathbb{P}^B} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}^A} [\lambda d_{\mathcal{X}}(\mathbf{x}, L(\mathbf{x})) + d_{\mathcal{Y}}(h^A(\mathbf{x}), h^B(L(\mathbf{x})))].$$

Symmetrically, we achieve the other equality. □

**Theorem 2.** *The following statements hold*

i)  $\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{B,h^B}) \geq \lambda \mathcal{W}_{d_{\mathcal{X}}}(\mathbb{P}^A, \mathbb{P}^B)$ .

ii) Given  $\mathcal{X}^A = \mathcal{X}^B = \mathcal{X}$ ,  $\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{B,h^B}) = 0$  if and only if  $\mathbb{P}^A = \mathbb{P}^B$  and  $h^A = h^B$ .

iii) Consider the problem:  $\min_{h^A} \mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{A,f^A})$ , the optimal solution is  $h^A_* = f^A$  obtained with the optimal mover  $L^* : L^*_{\#}\mathbb{P}^A = \mathbb{P}^A$  to be the identity map.

iv)  $\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{A,f^A}) \leq \mathcal{L}(h^A, f^A, \mathbb{P}^A)$ .

*Proof.* i) Using Proposition 1, we have

$$\begin{aligned} \mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{B,h^B}) &= \min_{L: L_{\#}\mathbb{P}^A = \mathbb{P}^B} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}^A} [\lambda d_{\mathcal{X}}(\mathbf{x}, L(\mathbf{x})) + d_{\mathcal{Y}}(h^A(\mathbf{x}), h^B(L(\mathbf{x})))] \\ &\geq \min_{L: L_{\#}\mathbb{P}^A = \mathbb{P}^B} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}^A} [\lambda d_{\mathcal{X}}(\mathbf{x}, L(\mathbf{x}))] = \lambda \mathcal{W}_{d_{\mathcal{X}}}(\mathbb{P}^A, \mathbb{P}^B). \end{aligned}$$

ii) It is obvious that if  $\mathbb{P}^A = \mathbb{P}^B$  and  $h^A = h^B$ , then  $\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{B,h^B}) = 0$ . For the converse, note that

$$\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{B,h^B}) \geq \lambda \mathcal{W}_{d_{\mathcal{X}}}(\mathbb{P}^A, \mathbb{P}^B)$$

by statement (i). Thus,  $\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{B,h^B}) = 0$  implies that  $\mathcal{W}_{d_{\mathcal{X}}}(\mathbb{P}^A, \mathbb{P}^B) = 0$  and hence  $\mathbb{P}^A = \mathbb{P}^B$ . Next, let  $L^*$  be the optimal transport:

$$L^* = \operatorname{argmin}_{L: L_{\#}\mathbb{P}^A = \mathbb{P}^B} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}^A} [\lambda d_{\mathcal{X}}(\mathbf{x}, L(\mathbf{x})) + d_{\mathcal{Y}}(h^A(\mathbf{x}), h^B(L(\mathbf{x})))].$$

Then thanks to Proposition 1, we have

$$\mathbb{E}_{\mathbf{x} \sim \mathbb{P}^A} [\lambda d_{\mathcal{X}}(\mathbf{x}, L^*(\mathbf{x})) + d_{\mathcal{Y}}(h^A(\mathbf{x}), h^B(L^*(\mathbf{x})))] = 0$$

which yield  $L^*(\mathbf{x}) = \mathbf{x}$  and hence  $h^A(\mathbf{x}) = h^B(\mathbf{x})$  almost everywhere w.r.t the measure  $\mathbb{P}^A$ .

iii) We have

$$\min_{h^A} \mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{A,f^A}) = 0,$$

and the minimum value is attained at  $h_*^A = f^A$  and  $L^* = id$ . From ii), this optimal solution is unique.

iv) Using Proposition 1, we have:

$$\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{A,f^A}) = \min_{L: L_{\#} \mathbb{P}^A = \mathbb{P}^A} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}^A} [\lambda d_{\mathcal{X}}(\mathbf{x}, L(\mathbf{x})) + d_{\mathcal{Y}}(h^A(\mathbf{x}), f^A(L(\mathbf{x})))].$$

Then by choosing  $L$  as the identity map (i.e.,  $L(\mathbf{x}) = \mathbf{x}$  for all  $\mathbf{x}$ ), we obtain:

$$\mathcal{W}_d(\mathbb{P}_{A,h^A}, \mathbb{P}_{A,f^A}) \leq \mathbb{E}_{\mathbf{x} \sim \mathbb{P}^A} [d_{\mathcal{Y}}(h^A(\mathbf{x}), f^A(\mathbf{x}))] = \mathcal{L}(h^A, f^A, \mathbb{P}^A).$$

□

## 2 OUR PROPOSED METHOD

### 2.1 PROBLEM FORMULATION

In multi-source domain adaptation, we have  $K$  multiple source domains with collected data and labels, and single target domain with only collected data. We wish to transfer a model learned on labeled source domains to an unlabeled target domain.

Let us denote the collected data and labels for the source domains by  $\mathcal{D}_k^S = \{(\mathbf{s}\mathbf{x}_i^k, y_i^k)\}_{i=1}^{N_k^S}$  with label  $y_i^k \in \{1, 2, \dots, M\}$  and collected data without labels for the target domain  $\mathcal{D}^T = \{\mathbf{t}\mathbf{x}_i\}_{i=1}^{N^T}$ .

For the sake of simplification, we denote the common space for source domains by  $\mathcal{X}^S$ . Note that if source domains have different input spaces, we can resize either input images or use corresponding transformations to map them to a common space. We further equip source domains with data distribution  $\mathbb{P}_{1:K}^S$  whose density functions are  $p_{1:K}^S(\mathbf{x})$ . Let us denote the ground-truth labeling functions for source domains by  $f_{1:K}^S(\cdot) \in \mathcal{Y}_{\Delta}$ , implying that  $p_k^S(y | \mathbf{x}) = f_k^S(\mathbf{x}, y)$  (i.e.,  $f_k^S(\mathbf{x}, y)$  represents the  $y$ -th value of  $f_k^S(\mathbf{x})$ ). Therefore, the joint distribution to generate data instance  $\mathbf{x}$  and categorical label  $y \in \{1, \dots, M\}$  is  $p_k^S(\mathbf{x}, y) = p_k^S(\mathbf{x}) f_k^S(\mathbf{x}, y)$ .

Regarding the target domain, we define its data space as  $\mathcal{X}^T$ , data distribution and density function as  $\mathbb{P}^T$  and  $p^T(\mathbf{x})$ , respectively. We further define the ground-truth labeling function for the target domain by  $f^T$  which subsequently implies  $p^T(y | \mathbf{x}) = f^T(\mathbf{x}, y)$  for a categorical label  $y \in \{1, \dots, M\}$ .

Given a discrete distribution  $\pi$  over  $\{1, \dots, K\}$ , we define  $\mathbb{P}_{\pi}^S := \sum_{k=1}^K \pi_k \mathbb{P}_k^S$  which is a mixture of  $\mathbb{P}_{1:K}^S$ . For a data instance  $\mathbf{x} \sim \mathbb{P}_{\pi}^S$  (i.e., we sample a hidden index  $t \sim \text{Cat}(\pi)$  (i.e., the category distribution) and then sample  $\mathbf{x} \sim \mathbb{P}_t^S$ ), we further define  $f^S$  as a labeling function such that  $f^S(\mathbf{x})$  is identical to  $f_t^S(\mathbf{x})$ . By this definition,  $f^S$  can be viewed as the ground-truth labeling function over the mixture distribution  $\mathbb{P}_{\pi}^S$ . Finally, the mixing proportion  $\pi$  can be the uniform distribution  $[\frac{1}{K}, \dots, \frac{1}{K}]$  or proportional to the number of training examples in the source domains (i.e.,  $N_{1:K}^S$ ).

### 2.2 MULTI-SOURCE EXPERT TEACHER

Using the labeled source training sets  $\mathcal{D}_{1:K}^S$ , we can train qualified domain expert classifier  $h_{1:K}^S$  (i.e.,  $h_k^S(\mathbf{x}) \in \mathcal{Y}_{\Delta}$  represents the prediction probability of  $h_k^S$  for a data instance  $\mathbf{x}$ ) with good generalization capacity (e.g.,  $\mathcal{L}(h_k^S, f_k^S, \mathbb{P}_k^S) \leq \epsilon$  for some small  $\epsilon > 0$ ). The next arising question is how to combine those domain experts to achieve a multi-source expert teacher  $h^S$  that can work well on  $\mathbb{P}_{\pi}^S$  (i.e.,  $\mathcal{L}(h^S, f^S, \mathbb{P}_{\pi}^S) \leq \epsilon$ ). We leverage the domain experts to achieve a more powerful multi-source expert teacher by a weighted ensembling as follows:

$$h^S(\mathbf{x}, y) = \sum_{k=1}^K \frac{\pi_k p_k^S(\mathbf{x}, y)}{\sum_{j=1}^K \pi_j p_j^S(\mathbf{x}, y)} h_k^S(\mathbf{x}, y), \quad (2)$$

where  $y \in \{1, 2, \dots, M\}$ , and  $h_k^S(\mathbf{x}, y)$  and  $h^S(\mathbf{x}, y)$  specify the  $y$ -th values of  $h_k^S(\mathbf{x})$  and  $h^S(\mathbf{x})$  respectively.

The following theorem shows that the multi-source domain expert teacher  $h^S$  can work well on the mixture joint distribution  $\mathbb{P}_{\pi}^S$ . More specifically, it works better than the worst domain expert on its source domain, hence if each domain expert is an  $\epsilon$ -qualified classifier (i.e.,  $\mathcal{L}(h_k^S, f_k^S, \mathbb{P}_k^S) \leq \epsilon$ ), the multi-source expert teacher  $h^S$  is also an  $\epsilon$ -qualified classifier (i.e.,  $\mathcal{L}(h^S, f^S, \mathbb{P}_{\pi}^S) \leq \epsilon$ ).

**Theorem 3.** If  $d_Y$  can be decomposed as  $d_Y(\boldsymbol{\alpha}, \boldsymbol{\beta}) := \sum_{i=1}^M \beta_i \ell(\alpha_i)$  where  $\boldsymbol{\alpha}, \boldsymbol{\beta} \in \mathcal{Y}_\Delta$  and  $\ell$  is a convex function, the following statements hold true:

i)  $\mathcal{L}(h^S, f^S, \mathbb{P}_\pi^S) \leq \max_{1 \leq k \leq K} \mathcal{L}(h_k^S, f_k^S, \mathbb{P}_k^S)$ .

ii) If each domain expert is an  $\epsilon$ -qualified classifier (i.e.,  $\mathcal{L}(h_k^S, f_k^S, \mathbb{P}_k^S) \leq \epsilon$ ), the multi-source expert teacher  $h^S$  is also an  $\epsilon$ -qualified classifier (i.e.,  $\mathcal{L}(h^S, f^S, \mathbb{P}_\pi^S) \leq \epsilon$ ).

*Proof.* i) We have

$$\begin{aligned}
\mathcal{L}(h^S, f^S, \mathbb{P}_\pi^S) &= \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_\pi^S} [d_Y(h^S(\mathbf{x}), f^S(\mathbf{x}))] = \sum_{k=1}^K \pi_k \int d_Y(h^S(\mathbf{x}), f_k^S(\mathbf{x})) d\mathbb{P}_k^S \\
&= \sum_{k=1}^K \pi_k \int d_Y(h^S(\mathbf{x}), f_k^S(\mathbf{x})) d\mathbb{P}_k^S = \sum_{k=1}^K \pi_k \int \sum_{i=1}^M f_k^S(\mathbf{x}, i) \ell(h^S(\mathbf{x}, i)) d\mathbb{P}_k^S \\
&= \sum_{i=1}^M \int \sum_{k=1}^K \pi_k p_k^S(y=i | \mathbf{x}) \ell(h^S(\mathbf{x}, i)) d\mathbb{P}_k^S = \sum_{i=1}^M \int \ell(h^S(\mathbf{x}, i)) \sum_{k=1}^K \pi_k p_k^S(y=i | \mathbf{x}) p_\pi^S(\mathbf{x}) d\mathbf{x} \\
&= \sum_{i=1}^M \int \ell(h^S(\mathbf{x}, i)) p_\pi^S(\mathbf{x}, y=i) d\mathbf{x} = \sum_{i=1}^M \int \ell \left( \sum_{k=1}^K \frac{\pi_k p_k^S(\mathbf{x}, y=i)}{\sum_{j=1}^K \pi_j p_j^S(\mathbf{x}, y=i)} h_k^S(\mathbf{x}, i) \right) p_\pi^S(\mathbf{x}, y=i) d\mathbf{x} \\
&\leq \sum_{i=1}^M \int \sum_{k=1}^K \frac{\pi_k p_k^S(\mathbf{x}, y=i)}{\sum_{j=1}^K \pi_j p_j^S(\mathbf{x}, y=i)} \ell(h_k^S(\mathbf{x}, i)) p_\pi^S(\mathbf{x}, y=i) d\mathbf{x} \quad (\text{thanks to the convexity of } \ell) \\
&= \sum_{i=1}^M \int \sum_{k=1}^K \frac{\pi_k p_k^S(\mathbf{x}, y=i)}{p_\pi^S(\mathbf{x}, y=i)} \ell(h_k^S(\mathbf{x}, i)) p_\pi^S(\mathbf{x}, y=i) d\mathbf{x} \\
&= \sum_{i=1}^M \int \sum_{k=1}^K \pi_k p_k^S(\mathbf{x}, y=i) \ell(h_k^S(\mathbf{x}, i)) d\mathbf{x} = \sum_{k=1}^K \pi_k \sum_{i=1}^M \int \ell(h_k^S(\mathbf{x}, i)) p_k^S(\mathbf{x}, y=i) d\mathbf{x} \\
&= \sum_{k=1}^K \pi_k \sum_{i=1}^M \int \ell(h_k^S(\mathbf{x}, i)) p_k^S(y=i | \mathbf{x}) p_k^S(\mathbf{x}) d\mathbf{x} \\
&= \sum_{k=1}^K \pi_k \int \sum_{i=1}^M \ell(h_k^S(\mathbf{x}, i)) f_k^S(\mathbf{x}, i) d\mathbb{P}_k^S = \sum_{k=1}^K \pi_k \int d_Y(h_k^S(\mathbf{x}), f_k^S(\mathbf{x})) d\mathbb{P}_k^S \\
&= \sum_{k=1}^K \pi_k \mathcal{L}(h_k^S, f_k^S, \mathbb{P}_k^S) \leq \max_{1 \leq k \leq K} \mathcal{L}(h_k^S, f_k^S, \mathbb{P}_k^S).
\end{aligned}$$

Note that we define  $p_\pi^S(\mathbf{x}, y=i) := \sum_{k=1}^K \pi_k p_k^S(\mathbf{x}, y=i) = \sum_{k=1}^K \pi_k p_k^S(y=i | \mathbf{x}) p_k^S(\mathbf{x})$ .

ii) It is trivial from (i). □

In what follows, we present how to train the multi-source expert teacher  $h^S$ . Our workaround to train  $h^S$  comes from the following theoretical observation. Assume that we have  $K$  distributions  $\mathbb{R}_{1:K}$  with density functions  $r_{1:K}(\mathbf{z})$ . We form a joint distribution  $\mathcal{D}$  of a data instance  $\mathbf{z}$  and label  $t \in \{1, \dots, K\}$  by sampling an index  $t \sim \text{Cat}(\boldsymbol{\pi})$ , sampling  $\mathbf{x} \sim \mathbb{R}_t$ , and collecting  $(\mathbf{z}, t)$  as a sample from  $\mathcal{D}$ . With this setting, we have the following corollary.

**Corollary 4.** If we train a source domain discriminator  $\mathcal{C}$  to classify samples from the joint distribution  $\mathcal{D}$  using the cross-entropy loss (i.e.,  $CE(\cdot, \cdot)$ ), the optimal source domain discriminator  $\mathcal{C}^*$  defined as

$$\mathcal{C}^* = \underset{\mathcal{C}}{\text{argmin}} \mathbb{E}_{(\mathbf{z}, t) \sim \mathcal{D}} [CE(\mathcal{C}(\mathbf{z}), t)]$$

satisfies  $\mathcal{C}^*(\mathbf{z}) = \left[ \frac{\pi_i r_i(\mathbf{z})}{\sum_j \pi_j r_j(\mathbf{z})} \right]_{i=1}^K$ .

*Proof.* We have

$$\begin{aligned}\mathbb{E}_{(\mathbf{z},t)\sim\mathcal{D}}[CE(\mathcal{C}(\mathbf{z}),t)] &= \sum_{t=1}^K \pi_t \int CE(\mathcal{C}(\mathbf{z}),t) r_t(\mathbf{z}) d\mathbf{z} \\ &= - \int \sum_{t=1}^K \log \mathcal{C}(\mathbf{z},t) \pi_t r_t(\mathbf{z}) d\mathbf{z}.\end{aligned}$$

Given  $\mathbf{z}$ , we now find  $\mathcal{C}^* = [\mathcal{C}_t^*]_{t=1}^K$  subjected to  $\|\mathcal{C}^*\|_1 = 1$  and  $\mathcal{C}^* \geq \mathbf{0}$  to maximize

$$\max_{\mathcal{C}:\|\mathcal{C}\|_1=1} \sum_{t=1}^K \log \mathcal{C}_t \pi_t r_t(\mathbf{z}).$$

The Lagrange function is as follows:

$$\mathcal{L} = \sum_{t=1}^K \log \mathcal{C}_t \pi_t r_t(\mathbf{z}) - \lambda \left( \sum_{t=1}^K \mathcal{C}_t - 1 \right).$$

Setting the derivatives to 0, we obtain

$$\frac{\partial \mathcal{L}}{\partial \mathcal{C}_t} = \frac{\pi_t r_t(\mathbf{z})}{\mathcal{C}_t} - \lambda = 0, \quad t = 1, \dots, K.$$

Note that  $\sum_{t=1}^K \mathcal{C}_t = 1$ , we arrive at

$$\mathcal{C}_t^* = \frac{\pi_t r_t(\mathbf{z})}{\sum_j \pi_j r_j(\mathbf{z})}, \quad t = 1, \dots, K.$$

Finally, we reach the conclusion. □

Inspired by the statement (ii) in Theorem 2, recall that  $f^T$  is the ground-truth labeling function on the target domain, we propose to learn a classifier  $h^T$  on this domain to further minimize with the aim to obtain  $h^T = f^T$ :

$$\min_{h^T} \mathcal{W}_d(\mathbb{P}_{T,h^T}, \mathbb{P}_{T,f^T}).$$

To proceed our theory, we assume that  $d_{\mathcal{Y}}$  is a metric over  $\mathcal{Y}_{\Delta}$ , which together with the metric  $d_{\mathcal{X}}$  forms the metric  $d = \lambda d_{\mathcal{X}} + d_{\mathcal{Y}}$ , implying that  $\mathcal{W}_d(\mathbb{P}_{\cdot,\cdot}, \mathbb{P}_{\cdot,\cdot})$  is a proper metric. We can thus bound the quantity of interest  $\mathcal{W}_d(\mathbb{P}_{T,h^T}, \mathbb{P}_{T,f^T})$  as

$$\mathcal{W}_d(\mathbb{P}_{T,h^T}, \mathbb{P}_{T,f^T}) \leq \mathcal{W}_d(\mathbb{P}_{T,h^T}, \mathbb{P}_{S,h^S}^{\pi}) + \mathcal{W}_d(\mathbb{P}_{S,h^S}^{\pi}, \mathbb{P}_{S,f^S}^{\pi}) + \mathcal{W}_d(\mathbb{P}_{S,f^S}^{\pi}, \mathbb{P}_{T,f^T}), \quad (3)$$

where  $\mathbb{P}_{S,f^S}^{\pi}$ , a joint distribution over  $\mathcal{X}^S \times \mathcal{Y}_{\Delta}$ , consists of pairs  $(\mathbf{x}, y_{\Delta})$  in which  $\mathbf{x} \sim \mathbb{P}_{\pi}^S$  and  $y_{\Delta} = f^S(\mathbf{x})$ .

In the upper-bound in (3),  $\mathcal{W}_d(\mathbb{P}_{S,f^S}^{\pi}, \mathbb{P}_{T,f^T})$  is constant. We employ the multi-source expert teacher  $h^S$  as in Section 5.2, which can operate well on  $\mathbb{P}_{\pi}^S$  as long as we can train good domain experts  $h_{1:K}^S$ . Noting that  $\mathcal{W}_d(\mathbb{P}_{S,h^S}^{\pi}, \mathbb{P}_{S,f^S}^{\pi})$  is upper-bounded by  $\mathcal{L}(h^S, f^S, \mathbb{P}_{\pi}^S)$  (thanks to the statement (iii) in Theorem 2), we arrive at the following optimization problem:

$$\min_{h^T} \left\{ \mathcal{W}_d(\mathbb{P}_{T,h^T}, \mathbb{P}_{S,h^S}^{\pi}) + \mathcal{L}(h^S, f^S, \mathbb{P}_{\pi}^S) \right\}. \quad (4)$$

The optimization problem in (4) is in line with the context of imitation learning for which the teacher classifier  $h^S$  has been trained effectively on the mixture source domain (i.e.,  $\mathbb{P}_{\pi}^S$ ) and the student classifier  $h^T$  tries to imitate the teacher on the target domain. Specifically, Proposition 1 implies finding the optimal transport map  $H^*: H_{\#}^* \mathbb{P}^T = \mathbb{P}_{\pi}^S$  so that for

any  $\mathbf{x} \sim \mathbb{P}^T$ ,  $h^T(\mathbf{x})$  should mimic the prediction of the expert teacher  $h^S$  over  $H^*(\mathbf{x}) \sim \mathbb{P}_\pi^S$ . This observation forms the foundation of our proposed Multi-Source Domain Adaptation via Optimal Transport for Student-Teacher Learning (MOST).

Proposition 1 further means that among the transport maps  $H$  transporting  $\mathbb{P}^T$  to  $\mathbb{P}_\pi^S$ , which incurs a minimal label shift and enables the student  $h^T$  easiest to imitate its teacher  $h^S$ . Inspired by the statement (iv) in Theorem 2 (i.e.,  $\mathcal{W}_d(\mathbb{P}_{S,h^S}^\pi, \mathbb{P}_{T,h^T})$  is lower-bounded by  $\lambda \mathcal{W}_{d,\mathcal{X}}(\mathbb{P}_\pi^S, \mathbb{P}^T)$ —the discrepancy gap between the mixture of source distributions and the target one), to reduce the data shift, we propose to map both  $(\mathcal{X}^S, \mathbb{P}_\pi^S)$  and  $(\mathcal{X}^T, \mathbb{P}^T)$  to a common joint space via two generators  $G^S$  and  $G^T$  and solve the following optimization problem:

$$\min_{h^T, G^T} \left\{ \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S) + \mathcal{W}_d(\mathbb{Q}_{T,h^T}, \mathbb{Q}_{S,h^S}^\pi) \right\}, \quad (5)$$

where  $\mathbb{Q}_{T,h^T}$  is similar to  $\mathbb{P}_{T,h^T}$  but on the joint space and consists of the pairs  $(G^T(\mathbf{x}), h^T(G^T(\mathbf{x})))$  for  $\mathbf{x} \sim \mathbb{P}^T$  and  $\mathbb{Q}_{S,h^S}^\pi$  is similar to  $\mathbb{P}_{S,h^S}^\pi$  but on the joint space and consists of the pairs  $(G^S(\mathbf{x}), h^S(G^S(\mathbf{x})))$  for  $\mathbf{x} \sim \mathbb{P}^S$ . Note that both  $h^S$  and  $h_{1:K}^S$  now act on  $G^S(\cdot)$ .

**Theorem 5.** Let  $h_*^S \circ G_*^S$  be the optimal teacher and  $h_*^T, G_*^T$  be the optimal solutions of the optimization problem in (5). Assume that  $G^T, h^T$  are in the families having infinite capacity (i.e., those can approximate any continuous function up to any level of precision, e.g., neural nets), we have<sup>1</sup>

$$\min_{h^T, G^T} \mathcal{W}_d(\mathbb{P}_{T,h^T}^{G^T}, \mathbb{P}_{T,f_T^T}^{G^T}) \leq \mathcal{L}(h_*^S \circ G_*^S, f^S, \mathbb{P}_\pi^S) + \mathcal{W}_d(\mathbb{P}_{S,f_*^S}^{G^S}, \mathbb{P}_{T,f_*^T}^{G^T}), \quad (6)$$

where  $f_*^S := f_{S_*}^{G^S}$  and  $f_*^T := f_{T_*}^{G^T}$ . where  $f_*^S := f_{S_*}^{G^S}$  and  $f_*^T := f_{T_*}^{G^T}$ .

*Proof.* Let  $f_T^{G^T}$  be the induced ground-truth labeling function over the joint space for which  $f_T^{G^T}$  predicts  $G^T(\mathbf{x})$  as same as  $f^T$  predicts  $\mathbf{x}$  for  $\mathbf{x} \sim \mathbb{P}^T$ . Let  $f_S^{G^S}$  be the induced ground-truth labeling function over the joint space for which  $f_S^{G^S}$  predicts  $G^S(\mathbf{x})$  as same as  $f^S$  predicts  $\mathbf{x}$  for  $\mathbf{x} \sim \mathbb{P}_\pi^S$ . We have the following inequality

$$\begin{aligned} \mathcal{W}_d(\mathbb{P}_{T,h^T}^{G^T}, \mathbb{P}_{T,f_T^{G^T}}^{G^T}) &\leq \mathcal{W}_d(\mathbb{P}_{T,h^T}^{G^T}, \mathbb{P}_{S,h^S}^{G^S}) + \mathcal{W}_d(\mathbb{P}_{S,h^S}^{G^S}, \mathbb{P}_{S,f_S^{G^S}}^{G^S}) + \mathcal{W}_d(\mathbb{P}_{S,f_S^{G^S}}^{G^S}, \mathbb{P}_{T,f_T^{G^T}}^{G^T}) \\ &\stackrel{(1)}{\leq} \mathcal{W}_d(\mathbb{P}_{T,h^T}^{G^T}, \mathbb{P}_{S,h^S}^{G^S}) + \mathcal{L}(h^S, f_S^{G^S}, \mathbb{P}_\pi^S) + \mathcal{W}_d(\mathbb{P}_{S,f_S^{G^S}}^{G^S}, \mathbb{P}_{T,f_T^{G^T}}^{G^T}) \\ &\stackrel{(2)}{=} \mathcal{W}_d(\mathbb{P}_{T,h^T}^{G^T}, \mathbb{P}_{S,h^S}^{G^S}) + \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S) + \mathcal{W}_d(\mathbb{P}_{S,f_S^{G^S}}^{G^S}, \mathbb{P}_{T,f_T^{G^T}}^{G^T}), \end{aligned} \quad (7)$$

where  $\mathbb{P}_{S,h^S}^{G^S}$  consists of the pairs  $(G^S(\mathbf{x}), h^S(G^S(\mathbf{x})))$  for which  $\mathbf{x} \sim \mathbb{P}_\pi^S$ ,  $\mathbb{P}_{S,f_S^{G^S}}^{G^S}$  consists of the pairs  $(G^S(\mathbf{x}), f_S^{G^S}(G^S(\mathbf{x})))$  for which  $\mathbf{x} \sim \mathbb{P}_\pi^S$ ,  $\mathbb{P}_{T,f_T^{G^T}}^{G^T}$  consists of the pairs  $(G^T(\mathbf{x}), f_T^{G^T}(G^T(\mathbf{x})))$  for which  $\mathbf{x} \sim \mathbb{P}^T$ ,  $\mathbb{P}_{T,h^T}^{G^T}$  consists of the pairs  $(G^T(\mathbf{x}), h^T(G^T(\mathbf{x})))$  for which  $\mathbf{x} \sim \mathbb{P}^T$ , and  $\mathbb{P}_\pi^{G^S} = G^S \# \mathbb{P}_\pi^S$ .

We note that the derivation in (1) is from the statement (iv) in Theorem 2 and to reach (2), we derive as follows:

$$\begin{aligned} \mathcal{L}(h^S, f_S^{G^S}, \mathbb{P}_\pi^{G^S}) &= \int d\mathbf{y} (h^S(\mathbf{y}), f_S^{G^S}(\mathbf{y})) d\mathbb{P}_\pi^{G^S} = \int d\mathbf{y} (h^S(G^S(\mathbf{x})), f_S^{G^S}(G^S(\mathbf{x}))) d\mathbb{P}_\pi^S \\ &= \int d\mathbf{y} (h^S(G^S(\mathbf{x})), f^S(\mathbf{x})) d\mathbb{P}_\pi^S = \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S), \end{aligned}$$

where the second identity is due to  $G^S \# \mathbb{P}_\pi^S = \mathbb{P}_\pi^{G^S}$ .

To proceed, we first observe that

$$\min_{h^S, h^T, G^S, G^T} \left\{ \mathcal{W}_d(\mathbb{P}_{T,h^T}^{G^T}, \mathbb{P}_{S,h^S}^{G^S}) + \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S) \right\} = \min_{h^S, G^S} \left\{ \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S) \right\}, \quad (8)$$

<sup>1</sup>We define  $f^G$  as the induced labeling function over the joint space such that  $f^G$  predicts  $G(\mathbf{x})$  as same as  $f$  predicts  $\mathbf{x}$ .

where we consider  $h^S$  in the family of expert teachers (i.e., those which are a combination of the domain experts).

Indeed, let  $(\bar{h}^S, \bar{G}^S) := \operatorname{argmin}_{h^S, G^S} \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S)$ . Thanks to the infinite capacity of the families, we can choose the

generator  $\bar{G}^T$  and the target classifier  $\bar{h}^T$  such that  $\bar{G}^T \# \mathbb{P}^T = \bar{G}^S \# \mathbb{P}^S$  and  $\bar{h}^T = \bar{h}^S$  (i.e.,  $\mathcal{W}_d(\mathbb{P}_{T, \bar{h}^T}^{\bar{G}^T}, \mathbb{P}_{S, \bar{h}^S}^{\bar{G}^S}) = 0$ ).

Then we have

$$\begin{aligned} \min_{h^S, G^S} \{ \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S) \} &= \mathcal{W}_d(\mathbb{P}_{T, \bar{h}^T}^{\bar{G}^T}, \mathbb{P}_{S, \bar{h}^S}^{\bar{G}^S}) + \mathcal{L}(\bar{h}^S \circ \bar{G}^S, f^S, \mathbb{P}_\pi^S) \\ &\geq \min_{h^S, h^T, G^S, G^T} \left\{ \mathcal{W}_d(\mathbb{P}_{T, h^T}^{G^T}, \mathbb{P}_{S, h^S}^{G^S}) + \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S) \right\}. \end{aligned}$$

Therefore, (8) holds as the reverse inequality is obvious.

Now let  $h_*^S, h_*^T, G_*^S, G_*^T$  be the optimal solutions of the left side of (8). Then by (8), we must have

$$\mathcal{W}_d(\mathbb{P}_{T, h_*^T}^{G_*^T}, \mathbb{P}_{S, h_*^S}^{G_*^S}) + \mathcal{L}(h_*^S \circ G_*^S, f^S, \mathbb{P}_\pi^S) = \min_{h^S, G^S} \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S).$$

This implies that  $\mathcal{W}_d(\mathbb{P}_{T, h_*^T}^{G_*^T}, \mathbb{P}_{S, h_*^S}^{G_*^S}) = 0$  and  $(h_*^S, G_*^S) = \operatorname{argmin}_{h^S, G^S} \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S)$ . From Theorem (2), we see that the first identity means that  $G_*^S \# \mathbb{P}_\pi^S = G_*^T \# \mathbb{P}^T$  and  $h_*^T = h_*^S$ .

To further proceed, we refer to the triangle inequality in (7) as

$$\mathcal{W}_d(\mathbb{P}_{T, h^T}^{G^T}, \mathbb{P}_{T, f_T^{G^T}}^{G^T}) \leq \mathcal{W}_d(\mathbb{P}_{T, h^T}^{G^T}, \mathbb{P}_{S, h^S}^{G^S}) + \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S) + \mathcal{W}_d(\mathbb{P}_{S, f_S^{G^S}}^{G^S}, \mathbb{P}_{T, f_T^{G^T}}^{G^T}).$$

Taking a minimization and using  $\mathcal{W}_d(\mathbb{P}_{T, h_*^T}^{G_*^T}, \mathbb{P}_{S, h_*^S}^{G_*^S}) = 0$ , we obtain

$$\begin{aligned} \min_{h^T, G^T} \mathcal{W}_d(\mathbb{P}_{T, h^T}^{G^T}, \mathbb{P}_{T, f_T^{G^T}}^{G^T}) &\leq \min_{h^T, h^S, G^T, G^S} \left\{ \mathcal{W}_d(\mathbb{P}_{T, h^T}^{G^T}, \mathbb{P}_{S, h^S}^{G^S}) + \mathcal{L}(h^S \circ G^S, f^S, \mathbb{P}_\pi^S) + \mathcal{W}_d(\mathbb{P}_{S, f_S^{G^S}}^{G^S}, \mathbb{P}_{T, f_T^{G^T}}^{G^T}) \right\} \\ &\leq \mathcal{W}_d(\mathbb{P}_{T, h_*^T}^{G_*^T}, \mathbb{P}_{S, h_*^S}^{G_*^S}) + \mathcal{L}(h_*^S \circ G_*^S, f^S, \mathbb{P}_\pi^S) + \mathcal{W}_d(\mathbb{P}_{S, f_S^{G_*^S}}^{G_*^S}, \mathbb{P}_{T, f_T^{G_*^T}}^{G_*^T}) \\ &= \mathcal{L}(h_*^S \circ G_*^S, f^S, \mathbb{P}_\pi^S) + \mathcal{W}_d(\mathbb{P}_{S, f_S^{G_*^S}}^{G_*^S}, \mathbb{P}_{T, f_T^{G_*^T}}^{G_*^T}). \end{aligned}$$

□

### 3 CLUSTERING VIEW OF OPTIMAL TRANSPORT

We now present the clustering view of optimal transport which assists us to intuitively explain the training of teacher and student. Let us start with the clustering view of optimal transport.

Let  $\mathbb{P}$  and  $\mathbb{Q}$  be two discrete distributions defined as

$$\mathbb{P} := \frac{1}{m} \sum_{i=1}^m \delta_{\mathbf{u}_i} \text{ and } \mathbb{Q} := \frac{1}{n} \sum_{j=1}^n \delta_{\mathbf{v}_j},$$

where  $\delta_{\mathbf{x}}$  indicates a Dirac measure centered at  $\mathbf{x}$ .

Without loss of generality, we can assume that  $n \leq m$  and consider the Wasserstein distance  $\mathcal{W}_d(\mathbb{P}, \mathbb{Q})$  w.r.t a metric  $d$ . The following theorem characterizes the clustering view of optimal transport.

Table 1: Small, medium and large network architecture of MOST. The Leaky ReLU (lReLU) parameter  $a$  is set to 0.1.  $K$  and  $M$  denote the number of source domains and the number of classes respectively.

Architecture	Small	Large
Input size	2048	$32 \times 32 \times 3$
Generator $G_{1:K}^S, G^T$ (Share weights)		instance normalization
	256 dense, ReLU	$3 \times 3$ conv. 64 lReLU
	dropout, $p = 0.5$	$3 \times 3$ conv. 64 lReLU
	Gaussian noise, $\sigma = 1$	$3 \times 3$ conv. 64 lReLU
		$2 \times 2$ max-pool, stride 2
		dropout, $p = 0.5$
		Gaussian noise, $\sigma = 1$
		$3 \times 3$ conv. 64 lReLU
		$3 \times 3$ conv. 64 lReLU
		$3 \times 3$ conv. 64 lReLU
Classifier $h_{1:K}^S, h^T$		$2 \times 2$ max-pool, stride 2
		dropout, $p = 0.5$
		Gaussian noise, $\sigma = 1$
	$M$ dense, softmax	$3 \times 3$ conv. 64 lReLU
Source domain discriminator $\mathcal{C}$	$K$ dense, ReLU	$3 \times 3$ conv. 64 lReLU
		$3 \times 3$ conv. 64 lReLU
$\phi$	1 dense, linear	$3 \times 3$ conv. 64 lReLU
		global average pool
		$M$ dense, softmax
		100 dense, ReLU
		$K$ dense, ReLU
		100 dense, ReLU
		1 dense, linear

**Theorem 6.** Consider the following optimization problem:

$$\min_{\mathbf{v}_{1:n}} \mathcal{W}_d(\mathbb{P}, \mathbb{Q}).$$

Let  $\mathbf{v}_{1:n}^*$  and  $\mathbb{Q}^* := \frac{1}{n} \sum_{j=1}^n \delta_{\mathbf{v}_j^*}$  be its optimal solution and  $T^*$  be the optimal transport map as

$$T^* = \operatorname{argmin}_{T: T_{\#}\mathbb{P}=\mathbb{Q}^*} \sum_{i=1}^m d(\mathbf{u}_i, T(\mathbf{u}_i)).$$

Furthermore, let  $\mathbf{c}_{1:n}^*$  and  $\sigma^*$  denote the optimal solution of the following clustering problem

$$\min_{\mathbf{c}_{1:n}, \sigma \in \Pi(m, n)} \sum_{i=1}^m d(\mathbf{u}_i, \mathbf{v}_{\sigma(i)}),$$

where  $\Pi(m, n)$  is the set of surjective maps from  $\{1, \dots, m\}$  to  $\{1, \dots, n\}$ . We then have  $\mathbf{c}_{1:n}^* = \mathbf{v}_{1:n}^*$  and  $T^*(\mathbf{u}_i) = \mathbf{v}_{\sigma^*(i)}^*$ .

*Proof.* The proof is quite obvious from the definition of optimal transport and the fact that  $T_{\#}\mathbb{P} = \mathbb{Q}$  means  $T(\mathbf{u}_i) = \mathbf{v}_{\sigma(i)}, \forall i = 1, \dots, n$  for some  $\sigma \in \Pi(m, n)$ .  $\square$

## 4 IMPLEMENTATION SPECIFICATION AND ADDITIONAL EXPERIMENTAL RESULTS

### 4.1 DATA PREPARATION AND PREPROCESSING

**Digits-five** consists of five-digit datasets: MNIST (**mt**), MNIST-M (**mm**), USPS (**up**), SVHN (**sv**), Synthetic Digits (**sy**). There are 10 classes corresponding to digits ranging from 0 to 9 in each domain. We resize the resolution of digit images to  $32 \times 32$ , and normalize the value of each pixel to the range of  $[-1, 1]$ .



Table 2: The experimental settings in all transfer tasks. The *batch size* column denotes the total batch sizes of all source domains in each iteration.  $\alpha, \beta, \gamma$  are trade-off parameters of  $\mathcal{L}^{WS}, \mathcal{L}^{pl}, \mathcal{L}^{clus}$  respectively in our final objective function (11).

dataset	domains	classes	pretrained model	input size	batch size	iterations	learning rate	$\lambda$	update $\phi$ (times)	$\epsilon$	$\alpha$	$\beta$	$\gamma$
Digits-five	5	10	None	$32 \times 32$	200	80000	$2 \times 10^{-4}$	10.0	5	0.1	0.1	{0.1, 1.0}	{0.0, 0.1}
Office-Caltech10	4	10	ResNet-101	$224 \times 224$	60	20000	$1 \times 10^{-4}$	10.0	5	0.1	0.1	0.1	0.1
Office-31	3	31	AlexNet	$227 \times 227$	62	20000	$1 \times 10^{-4}$	10.0	5	0.1	0.1	0.1	0.1

**Office-Caltech10** is categorized in four different domains: Amazon (**A**), Caltech (**C**), DSLR (**D**), and Webcam (**W**) with 10 common classes and 2533 images in total. The resolutions of images are resized to  $224 \times 224$  for finetuning pre-trained ResNet-101.

**Office-31** contains 4,110 images with 31 classes, and is categorized into three domains: Amazon (**A**), DSLR (**D**), and Webcam (**W**). To perform finetuning on AlexNet, we resized the resolutions of images in this dataset to  $227 \times 227$ .

## 4.2 ARCHITECTURE/HYPERPARAMETERS

### 4.2.1 Network architecture

We use a small network to finetune pretrained AlexNet and ResNet-101, while a larger network is employed to train on digit images. These network types are described in Table 1. We empirically find that using a shared-weight generator improves the imitation learning capability of the student to the teacher since target and source samples are closer in the joint space. To reduce the overfitting, excluding dense layers in  $\phi$  network, we add the batch normalization layers on top of convolutional and dense layers. Finally, we implement our MOST in Python using TensorFlow (version 1.9.0), an open-source software library for Machine Intelligence developed by the Google Brain Team. All experiments are run on a computer with an NVIDIA Tesla V100 SXM2 with 16 GB memory.

### 4.2.2 Hyperparameters

In the loss  $\mathcal{L}^{WS}$ , a number of parameters are crucial to make the cross-domain imitation possible such as  $\lambda$  and  $\epsilon$ . Specifically,  $\lambda$  is set to 10.0 on all transfer tasks, while the coefficient  $\epsilon$  specifies the shape of joint density for the optimal transport plan and is set to 0.1 in all settings. Finally, we apply Adam optimizer ( $\beta_1 = 0.5, \beta_2 = 0.999$ ) with Polyak averaging. The learning rate along with the number of training iterations are depicted in Table 2.

## 4.3 ADDITIONAL ABLATION STUDY

### 4.3.1 Parameter Sensitivity

We evaluate the effects of the trade-off parameters  $\alpha, \beta, \gamma$  in Figure 1. We search all  $\alpha, \beta$  and  $\gamma$  in the grid of  $\{0.005, 0.01, 0.05, 0.1, 0.5, 1.0, 5.0\}$  and report the test accuracy on two transfer tasks “ $\rightarrow \mathbf{A}$ ” and “ $\rightarrow \mathbf{C}$ ” on *Office-Caltech10*. The results show that the model yields the stable performance when  $\alpha$  from 0.005 to 0.5,  $\beta$  from 0.05 to 5.0, and  $\gamma$  from 0.005 to 0.1. We find that our MOST can achieve high performance when  $\alpha = \beta = \gamma = 0.1$ , hence we suggest to pick this value on most of our experiments.

### 4.3.2 Feature visualization

We visualize the features of ResNet-101 and our methods on “ $\rightarrow \mathbf{C}$ ” tasks by *t*-SNE in Figure 2. Figure 2a shows that ResNet-101 classifies quite well on the mixture of source domains (**A**, **D**, **W**) but poorly on the target domain (**C**), while the representation in Figure 2b is generated by our method with better alignment. MOST represents the ability of reducing the data shift and label shift between two domains to exactly achieve 10 clusters corresponding to 10 classes of *Office-Caltech10*.

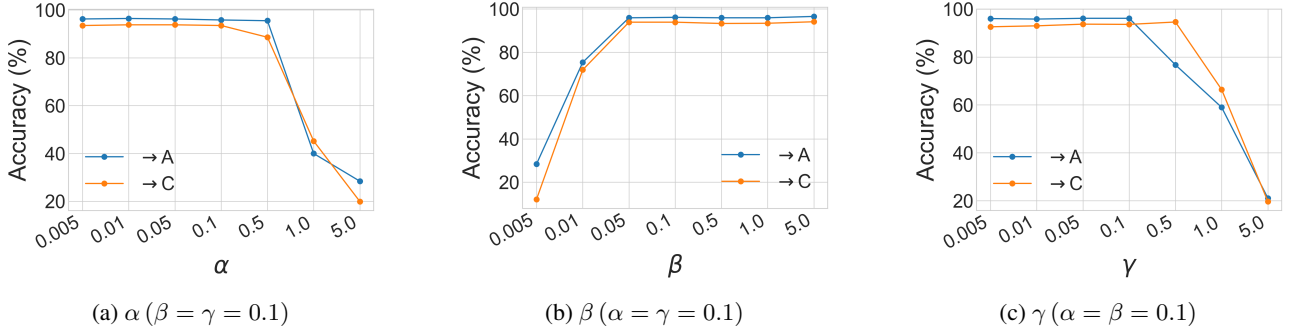


Figure 1: Test accuracy (%) when tweaking  $\alpha$ ,  $\beta$  and  $\gamma$  on transfer tasks “ $\rightarrow A$ ” and “ $\rightarrow C$ ”.

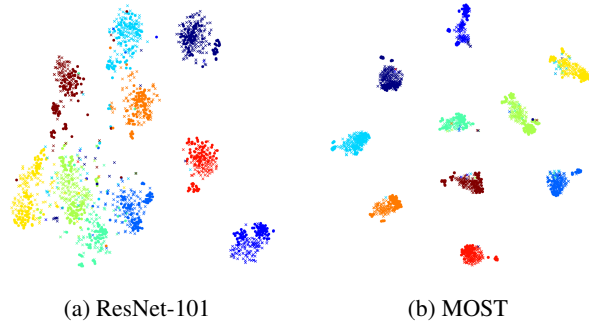


Figure 2: The  $t$ -SNE visualization of the transfer task “ $\rightarrow C$ ” with label and domain information. Each color denotes a class while the circle and cross markers represent the mixture of source and target data respectively.

### 4.3.3 Effect of parameters on imitation learning

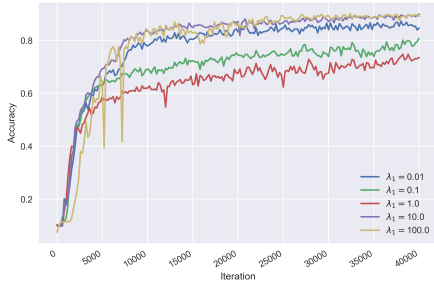
The loss  $\mathcal{L}^{WS}$  is introduced in Section 5.4.2 of the main paper. In this section, we further investigate sensitivity of two parameters,  $\lambda_1$  and  $\lambda_2$ , in  $d(G^S(\mathbf{x}^S), G^T(\mathbf{x}^T))$  to testify the imitation ability of MOST. They are trade-off parameters of  $d_X(\cdot, \cdot)$  and  $d_Y(\cdot, \cdot)$ , respectively as follows:

$$\mathcal{L}^{WS} = \max_{\phi} \left\{ \mathbb{E}_{\mathbb{P}^T} \left[ -\epsilon \log \left( \mathbb{E}_{\mathbb{P}_{\pi}^S} \left[ \exp \left\{ \frac{\phi(G^S(\mathbf{x}^S)) - d(G^S(\mathbf{x}^S), G^T(\mathbf{x}^T))}{\epsilon} \right\} \right] \right) \right] + \mathbb{E}_{\mathbb{P}_{\pi}^S} [\phi(G^S(\mathbf{x}^S))] \right\},$$

where  $d(G^S(\mathbf{x}^S), G^T(\mathbf{x}^T)) = \lambda_1 \|G^T(\mathbf{x}^T) - G^S(\mathbf{x}^S)\| + \lambda_2 d_Y(h^T(G^T(\mathbf{x}^T)), h^S(G^S(\mathbf{x}^S)))$ , while  $\mathbf{x}^T \sim \mathbb{P}^T$ ,  $\mathbf{x}^S \sim \mathbb{P}_{\pi}^S$ . Adjusting  $\lambda_1$  and  $\lambda_2$  helps to mitigate the data shift and label shift effectively. Figure 3a and 3b portray the model performance with a diverse range of values  $\lambda_1$  and  $\lambda_2$ , respectively. After searching them in the grid of  $\{0.01, 0.1, \dots, 100.0\}$ , we observe from Figure 3 that our model is more sensitive to  $\lambda_2$  since it directly controls the imitation capability of the student. In general, the training process of MOST is stable with  $\lambda_1 = 10.0$  and  $\lambda_2 = 1.0$ .

### 4.3.4 Choosing the mixing proportion $\pi$

This ablation study aims to testify the effect of choosing the mixture weights  $\pi$ . We conduct experiments on two scenarios: i)  $\pi$  is a uniform distribution  $[\frac{1}{K}, \dots, \frac{1}{K}]$  (**U**), and ii)  $\pi$  is proportional to the number of training examples in the source domains (i.e.,  $N_{1:K}^S$ ) (**P**). The experimental results in Table 3 show the competitive results between two approaches, which means that the imitation capability is not affected by changing the mixture distribution  $\mathbb{P}_{S, h^S}^{\pi}$  in (8).



(a)  $\lambda_1$  ( $\lambda_2 = 1.0$ )



(b)  $\lambda_2$  ( $\lambda_1 = 10.0$ )

Figure 3: Test accuracy (%) when tweaking  $\lambda_1$  and  $\lambda_2$  on “ $\rightarrow\mathbf{mm}$ ” task.

Table 3: Results (%) of different choices of  $\pi$  on Office-Caltech10.

Settings	$\rightarrow\mathbf{W}$	$\rightarrow\mathbf{D}$	$\rightarrow\mathbf{C}$	$\rightarrow\mathbf{A}$	Avg
<b>P</b>	99.7	100.0	95.1	96.5	97.8
<b>U</b>	100.0	100.0	95.7	96.3	98.0