

Keypoint-Guided Optimal Transport with Applications in Heterogeneous Domain Adaptation

A Mathematical Deductions

A.1 Proof of Proposition 1

Proposition 1 in the paper is for the case that $p_i = q_j$, for all $(i, j) \in \mathcal{K}$. As stated in the paper, the mask-based modeling of the transport plan is applicable even for the case that there exist some $(i, j) \in \mathcal{K}$ such that $p_i \neq q_j$. To see this, we first mathematically give the definition of preserving the matching of keypoint pairs and then prove Proposition A-1, a generalization of Proposition 1.

Definition A-1. Given the marginal distributions \mathbf{p} and \mathbf{q} , we say that the transport plan $\pi \in \Pi(\mathbf{p}, \mathbf{q})$ preserves the matching of a keypoint pair with index $(i, j) \in \mathcal{K}$, if π satisfies one of the following conditions:

1. If $p_i = q_j$, π satisfies that $\pi_{i,j'} = 0, \forall j' \neq j; \pi_{i',j} = 0, \forall i' \neq i; \pi_{i,j} = p_i = q_j$.
2. If $p_i > q_j$, π satisfies that $\pi_{i',j} = 0, \forall i' \neq i; \pi_{i,j} = q_j$.
3. If $p_i < q_j$, π satisfies that $\pi_{i,j'} = 0, \forall j' \neq j; \pi_{i,j} = p_i$.

□

The left part of Fig. A-1 illustrates these conditions. Specifically, the first condition implies that if $p_i = q_j$ (e.g., (i, j) is taken as (4, 4) in Fig. A-1), the all mass p_i of x_i will be transported to y_j and y_j can only receive mass from x_i . The second condition implies that if $p_i > q_j$ (e.g., (i, j) is taken as (3, 2) in Fig. A-1), y_j can only receive mass from x_i and consequently the partial mass $p_i - q_j$ of x_i is allowed to be transported to the target points apart from y_j . The third condition indicates that if $p_i < q_j$ (e.g., (i, j) is taken as (6, 5) in Fig. A-1), the all mass p_i of x_i will be transported to y_j and y_j is enabled to receive partial mass $q_j - p_i$ from the source points apart from x_i . For the convenience of description, for each pair $(i, j) \in \mathcal{K}$, we denote $j = \kappa(i)$ and $i = \kappa'(j)$.

Proposition A-1. Suppose that the mask matrix M satisfies that

$$M_{i,j} = \begin{cases} 1, & \text{if } (i, j) \in \mathcal{K}, \\ 0, & \text{if } i \in \mathcal{I}, p_i \leq q_{\kappa(i)}, \text{ and } (i, j) \notin \mathcal{K}, \\ 0, & \text{if } j \in \mathcal{J}, p_{\kappa'(j)} \geq q_j, \text{ and } (i, j) \notin \mathcal{K}, \\ 1, & \text{if } i \in \mathcal{I}, p_i > q_{\kappa(i)}, \text{ and } (i, j) \notin \mathcal{K}, \\ 1, & \text{if } j \in \mathcal{J}, p_{\kappa'(j)} < q_j, \text{ and } (i, j) \notin \mathcal{K}, \\ 1, & \text{otherwise (i.e., } i \notin \mathcal{I}, j \notin \mathcal{J}). \end{cases} \quad (\text{A-1})$$

Then, the transport plan $\tilde{\pi} = M \odot \pi$ with $\pi \in \Pi(\mathbf{p}, \mathbf{q}; M)$ preserves the matching of keypoint pairs with index in \mathcal{K} .

According to the definition of M , $M_{i,j} = 1$ for the keypoint pair $(i, j) \in \mathcal{K}$, implying that $\tilde{\pi}_{i,j}$ could take non-zero value. For $i \in \mathcal{I}$, $(i, j) \notin \mathcal{K}$ and $p_i \leq q_{\kappa(i)}$, $M_{i,j}$ is set to 0, enforcing that the i -th row of $\tilde{\pi}$ are zeros except for the location $\kappa(i)$ of the target keypoint paired with i (e.g., the 4-th and 6-th rows of $\tilde{\pi}$ in Fig. A-1). Similarly, for $j \in \mathcal{J}$, $(i, j) \notin \mathcal{K}$, and $p_{\kappa'(j)} \leq q_j$, we set $M_{i,j} = 0$, enforcing that the j -th column of $\tilde{\pi}$ are zeros except for the location $\kappa'(j)$ of the source keypoint paired with j (e.g., the 2-th and 4-th columns of $\tilde{\pi}$ in Fig. A-1). For the other points (corresponding to the last three cases in Eq. (A-1)), we set $M_{i,j} = 1$, indicating that there is no additional constraint on $\tilde{\pi}_{i,j}$. If $p_i = q_j$, for all $(i, j) \in \mathcal{K}$ (i.e., $p_i = q_{\kappa(i)}, \forall i \in \mathcal{I}$), Proposition A-1 degenerates to Proposition 1 in the paper.

Proof:

For any $(i, j) \in \mathcal{K}$, we next prove that $\tilde{\pi}$ preserves the matching of keypoint pair (i, j) .

- If $p_i = q_j$, from the definition of M , we have $M_{i,j'} = 0$ for all $j' \neq j$ and $M_{i,j} = 1$. Then, we have $\tilde{\pi}_{i,j'} = M_{i,j'} \pi_{i,j'} = 0$ for all $j' \neq j$. Since $\sum_{j'=1}^n \tilde{\pi}_{i,j'} = p_i$, we have $\tilde{\pi}_{i,j} = p_i$. Similarly, we have $M_{i',j} = 0$ for all $i' \neq i$. Then, we have $\tilde{\pi}_{i',j} = M_{i',j} \pi_{i',j} = 0$ for all $i' \neq i$, and $\tilde{\pi}_{i,j} = \sum_{i'=1}^m \tilde{\pi}_{i',j} = q_j$.

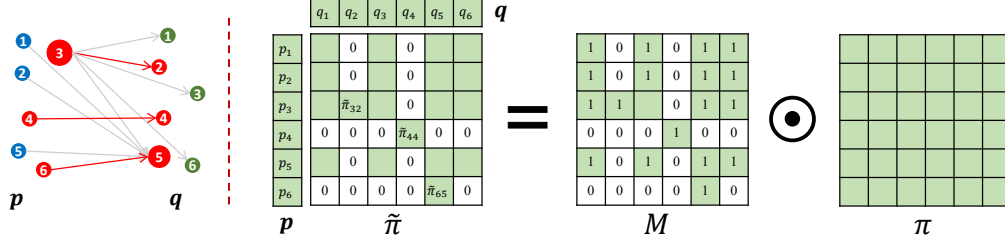


Figure A-1: Example of modeling the matching of keypoints (red) using mask, where $\mathcal{K} = \{(3, 2), (4, 4), (6, 5)\}$ with $p_3 > q_2, p_4 = q_4, p_6 < q_5$.

- If $p_i > q_j$, from the definition of M , we have $M_{i',j} = 0$ for all $i' \neq i$ and $M_{i,j} = 1$. Then, we have $\tilde{\pi}_{i',j} = M_{i',j}\pi_{i',j} = 0$ for all $i' \neq i$, and $\tilde{\pi}_{i,j} = \sum_{i'=1}^m \tilde{\pi}_{i',j} = q_j$.
- $p_i < q_j$, from the definition of M , we have $M_{i,j'} = 0$ for all $j' \neq j$ and $M_{i,j} = 1$. Then $\tilde{\pi}_{i,j'} = M_{i,j'}\pi_{i,j'} = 0$ for all $j' \neq j$, and $\tilde{\pi}_{i,j} = \sum_{j'=1}^n \tilde{\pi}_{i,j'} = q_i$.

Thus, for any keypoint pair with index $(i, j) \in \mathcal{K}$, $\tilde{\pi}$ satisfies the conditions in Definition A-1. This means that $\tilde{\pi}$ preserves the matching of keypoint pairs with index in \mathcal{K} .

A.2 Linear Programming for Solving KPG-RL

We cast the matrix G (resp. M, π) as the vector \mathbf{c} (resp. \mathbf{m}, \mathbf{x}) $\in \mathbb{R}^{mn}$, such that the $(i + m(j - 1))$ -th element of \mathbf{c} is G_{ij} . By denoting

$$A = \begin{bmatrix} \mathbb{1}_n^\top \otimes I_m \\ I_n \otimes \mathbb{1}_m \end{bmatrix} * \text{diag}(\mathbf{m}), \mathbf{h} = \begin{bmatrix} \mathbf{p} \\ \mathbf{q} \end{bmatrix}, \text{ and } \tilde{\mathbf{c}} = \mathbf{c} \odot \mathbf{m}, \quad (\text{A-2})$$

where I_m is the identity matrix of size n and \otimes is the Kronecker product, the KPG-RL model in Eq. (9) in the paper reads

$$\begin{aligned} \min_{\mathbf{x}} \tilde{\mathbf{c}}^\top \mathbf{x} \\ \text{s.t. } \mathbf{x} \geq 0, \\ A\mathbf{x} = \mathbf{h}. \end{aligned} \quad (\text{A-3})$$

With the standard form of linear programming in Eq. (A-3), the Simplex algorithm can be directly used to solve the KPG-RL model.

A.3 Sinkhorn's Algorithm for Solving KPG-RL

The entropy-regularized model for KPG-RL is

$$\begin{aligned} \min_{\pi} \langle M \odot \pi, G \rangle_F - \epsilon H(M \odot \pi) \\ \text{s.t. } \pi \geq 0, (M \odot \pi) \mathbb{1}_n = \mathbf{p}, (M \odot \pi)^\top \mathbb{1}_m = \mathbf{q}, \end{aligned} \quad (\text{A-4})$$

where $H(M \odot \pi) = -(\langle M \odot \pi, \log(M \odot \pi) \rangle_F - \mathbb{1}_m^\top (M \odot \pi) \mathbb{1}_n)$ is the entropy of the transport plan $M \odot \pi$. The Lagrangian function is

$$\begin{aligned} L(\pi, \mathbf{f}, \mathbf{g}) = \langle M \odot \pi, G \rangle_F + \epsilon (\langle M \odot \pi, \log(M \odot \pi) \rangle_F - \mathbb{1}_m^\top (M \odot \pi) \mathbb{1}_n) \\ - \langle \mathbf{f}, (M \odot \pi) \mathbb{1}_n - \mathbf{p} \rangle_F - \langle \mathbf{g}, (M \odot \pi)^\top \mathbb{1}_m - \mathbf{q} \rangle_F, \end{aligned} \quad (\text{A-5})$$

where $\mathbf{f} \in \mathbb{R}^m$ and $\mathbf{g} \in \mathbb{R}^n$. The first-order conditions then yield

$$\frac{\partial L}{\partial \pi_{i,j}} = M_{i,j} G_{i,j} + \epsilon M_{i,j} \log(M_{i,j} \pi_{i,j}) - M_{i,j} f_i - M_{i,j} g_j = 0. \quad (\text{A-6})$$

If $M_{i,j} = 0$, $\pi_{i,j}$ could be arbitrary non-negative value, and if $M_{i,j} = 1$, we have $\pi_{i,j} = e^{f_i/\epsilon} e^{-G_{i,j}/\epsilon} e^{g_j/\epsilon}$. Therefore, we can unify the expression as $\pi_{i,j} = M_{i,j} e^{f_i/\epsilon} e^{-C_{ij}/\epsilon} e^{g_j/\epsilon}$, in

which we enforce $\pi_{i,j} = 0$ if $M_{i,j} = 0$. The matrix form is $\pi = \text{diag}(\mathbf{u})K\text{diag}(\mathbf{v})$ where $\mathbf{u} = e^{\mathbf{f}/\epsilon}$, $K = M \odot e^{-G/\epsilon}$, and $\mathbf{v} = e^{\mathbf{g}/\epsilon}$. The constraints are

$$\text{diag}(\mathbf{u})K\text{diag}(\mathbf{v})\mathbb{1}_n = \mathbf{p}, \quad (\text{diag}(\mathbf{u})K\text{diag}(\mathbf{v}))^\top \mathbb{1}_m = \mathbf{q}. \quad (\text{A-7})$$

Since the entries of K are non-negative, the Sinkhorn's algorithm can be applied [15]. The iteration formulas are

$$\mathbf{u}^{(l+1)} = \frac{\mathbf{p}}{K\mathbf{v}^l}, \quad \mathbf{v}^{(l+1)} = \frac{\mathbf{q}}{K^\top \mathbf{u}^{(l+1)}}. \quad (\text{A-8})$$

The division operator used above is entry-wise.

Log-domain Sinkhorn iteration. For KP, the Sinkhorn iteration in the log-domain is more stable [60]. We next deduce the log-domain Sinkhorn iteration for our KPG-RL. In the log-domain, the left equation in Eq. (A-8) is

$$\frac{1}{\epsilon} f_i^{(l+1)} = \log(p_i) - \log \left(\sum_{j=1}^n M_{i,j} e^{-\frac{G_{i,j} + g_j^{(l)}}{\epsilon}} \right). \quad (\text{A-9})$$

Let $H(\mathbf{f}, \mathbf{g}) = \frac{1}{\epsilon}(-G + \mathbf{f}\mathbb{1}_n^\top + \mathbf{g}\mathbb{1}_m^\top)$, then

$$\begin{aligned} f_i^{(l+1)} &= \epsilon \log(p_i) - \epsilon \log \left(\sum_{j=1}^n M_{i,j} e^{H(\mathbf{f}^{(l)}, \mathbf{g}^{(l)})_{i,j}} e^{-f_i^{(l)}/\epsilon} \right) \\ &= \epsilon \log(p_i) - \epsilon \log \left(e^{-f_i^{(l)}/\epsilon} \sum_{j=1}^n M_{i,j} e^{H(\mathbf{f}^{(l)}, \mathbf{g}^{(l)})_{i,j}} \right) \\ &= \epsilon \log(p_i) - \epsilon \log \left(\sum_{j=1}^n M_{i,j} e^{H(\mathbf{f}^{(l)}, \mathbf{g}^{(l)})_{i,j}} \right) + f_i^{(l)} \\ &= \epsilon \log(p_i) - \epsilon \log \left(\sum_{j=1}^n e^{\log(M_{i,j})H(\mathbf{f}^{(l)}, \mathbf{g}^{(l)})_{i,j}} \right) + f_i^{(l)} \end{aligned} \quad (\text{A-10})$$

If $M_{i,j} = 0$, $\log(M_{i,j})H(\mathbf{f}^{(l)}, \mathbf{g}^{(l)})_{i,j} = -\infty$. We define $\bar{H}(\mathbf{f}, \mathbf{g})$ as

$$\bar{H}(\mathbf{f}, \mathbf{g})_{i,j} = \begin{cases} H(\mathbf{f}, \mathbf{g})_{i,j} & \text{if } M_{i,j} = 1, \\ -\infty & \text{if } M_{i,j} = 0, \end{cases} \quad (\text{A-11})$$

and define the log-sum-exp function $\text{LogSumExp} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^m$ as

$$\text{LogSumExp}(A) = \left(\log \left(\sum_j e^{A_{1,j}} \right), \log \left(\sum_j e^{A_{2,j}} \right), \dots, \log \left(\sum_j e^{A_{m,j}} \right) \right)^\top. \quad (\text{A-12})$$

Then, the matrix form of Eq. (A-10) becomes

$$\mathbf{f}^{(l+1)} = \epsilon \log(\mathbf{p}) - \epsilon \text{LogSumExp} \left(\bar{H}(\mathbf{f}^{(l)}, \mathbf{g}^{(l)}) \right) + \mathbf{f}^{(l)}. \quad (\text{A-13})$$

Similarly, the corresponding iteration formula in the log-domain of the right equation in (A-8) is

$$\mathbf{g}^{(l+1)} = \epsilon \log(\mathbf{q}) - \epsilon \text{LogSumExp} \left(\bar{H}(\mathbf{f}^{(l+1)}, \mathbf{g}^{(l)})^\top \right) + \mathbf{g}^{(l)}. \quad (\text{A-14})$$

Equations (A-13) and (A-14) consists in the formulas of the Sinkhorn iteration in the log-domain.

A.4 Frank-Walfe Algorithm for Solving KPG-RL-GW

We define the 4-order tensor $L = (L_{i,j,k,l}) \in \mathbb{R}^{m \times n \times m \times n}$ by $L_{i,j,k,l} = (C_{i,k}^s - C_{j,l}^t)^2$, and define the tensor-matrix product $L \circ \pi = ((L \circ \pi)_{i,j}) \in \mathbb{R}^{m \times n}$ by $(L \circ \pi)_{i,j} = \sum_{k,l} L_{i,j,k,l} \pi_{k,l}$. Then, The KPG-RL-GW model in Eq. (12) in the paper reads

$$\min_{\pi \in \Pi(\mathbf{p}, \mathbf{q}; M)} \langle (M \odot \pi), \alpha L \circ (M \odot \pi) + (1 - \alpha)G \rangle_F \triangleq \mathcal{L}(\pi) \quad (\text{A-15})$$

The gradient of the objective function is

$$\nabla \mathcal{L}(\pi) = M \odot (2\alpha L \circ (M \odot \pi) + (1 - \alpha)G). \quad (\text{A-16})$$

In k -th iteration of the Frank-Walfe algorithm, it runs the following three steps:

Step 1. Compute a linear minimization oracle over the set $\Pi(\mathbf{p}, \mathbf{q}; M)$, i.e.,

$$\hat{\pi} \leftarrow \underset{\pi \in \Pi(\mathbf{p}, \mathbf{q}; M)}{\operatorname{argmin}} \langle \nabla \mathcal{L}(\pi^{(k)}), \pi \rangle_F. \quad (\text{A-17})$$

Equation (A-17) can be rewritten as

$$\hat{\pi} \leftarrow \underset{\pi \in \Pi(\mathbf{p}, \mathbf{q}; M)}{\operatorname{argmin}} \langle M \odot \pi, 2\alpha L \circ (M \odot \pi^{(k)}) + (1 - \alpha)G \rangle_F, \quad (\text{A-18})$$

Equation (A-18) is a KPG-RL-like problem and can be solved using linear programming or Sinkhorn's algorithm.

Step 2. Determine optimal step-size $\beta^{(k)}$ subject to

$$\beta^{(k)} \leftarrow \underset{\beta \in [0,1]}{\operatorname{argmin}} \mathcal{L}((1 - \beta)\pi^{(k)} + \beta\hat{\pi}). \quad (\text{A-19})$$

$\beta^{(k)}$ can be obtained by the line-search method in [61].

Step 3. Update

$$\pi^{(k+1)} = (1 - \beta^{(k)})\pi^{(k)} + \beta^{(k)}\hat{\pi}. \quad (\text{A-20})$$

A.5 Theoretical Properties of KPG-RL-KP and KPG-RL-GW

In this section, we show that given prior ‘‘correct’’ paired keypoints, the KPG-RL-KP model provides a proper metric for distributions supported in the same space, and the the KPG-RL-GW model provides a divergence for distributions in distinct spaces, under mild conditions. Since the discrete distributions $\mathbf{p} = \frac{1}{m} \sum_i^m \delta_{x_i}$ (resp. $\mathbf{q} = \frac{1}{n} \sum_j^n \delta_{y_j}$) are invariant to the permutation of $\{x_i\}_{i=1}^m$ (resp. $\{y_j\}_{j=1}^n$), we assume that any two paired keypoints across domains share the same index. Therefore, the index set of paired keypoints is $\mathcal{K} = \{(i_u, i_u)\}_{u=1}^U$. We assume $p_{i_u} = q_{i_u}, \forall i_u$, in this section. For the convenience of description, in this section, we denote $M^{\mathbf{p}\mathbf{q}}$ as the mask matrix for transporting \mathbf{p} to \mathbf{q} , and $\mathcal{P}_{\mathcal{I}}^{\mathcal{X}}$ as the set of discrete probability distributions supported on m points in ground space \mathcal{X} such that all distributions in $\mathcal{P}_{\mathcal{I}}^{\mathcal{X}}$ share the keypoint index set $\mathcal{I} = \{i_u\}_{u=1}^U$.

A.5.1 KPG-RL-KP Providing a Proper Metric

For distributions supported in the same space, the ‘‘correct’’ paired keypoints indicates that if $\mathbf{p} = \mathbf{q}$, each source keypoint is equal to its paired target keypoint, i.e., $x_{i_u} = y_{i_u}$, for any $i_u \in \mathcal{I}$. We denote

$$\mathcal{S}_{krk}(\mathbf{p}, \mathbf{q}) = \min_{\pi \in \Pi(\mathbf{p}, \mathbf{q}; M^{\mathbf{p}\mathbf{q}})} \sum_{i,j} M_{i,j}^{\mathbf{p}\mathbf{q}} \pi_{i,j} (\alpha C_{i,j} + (1 - \alpha)G_{i,j}), \quad (\text{A-21})$$

where $\alpha \in (0, 1)$.

Theorem A-1. *Suppose c is a proper distance in space \mathcal{X} and d is a proper distance in probability simplex Σ_m . Then, for any \mathbf{p} and \mathbf{q} in $\mathcal{P}_{\mathcal{I}}^{\mathcal{X}}$, given the ‘‘correct’’ paired keypoints stated above, $\mathcal{S}_{krk}(\mathbf{p}, \mathbf{q})$ is a proper distance between \mathbf{p} and \mathbf{q} .*

Proof:

(1) **Show that $\mathcal{S}_{krk}(\mathbf{p}, \mathbf{q}) = 0$ if and only if $\mathbf{p} = \mathbf{q}$.** (a) If $\mathbf{p} = \mathbf{q}$, we have $x_i = y_i$ and $p_i = q_i$ for any $i \in [m]$ (since the permutation of support points does not change the distribution). Hence, $C_{i,i} = c(x_i, y_i) = 0$, and $C_{i,i_u}^s = c(x_i, x_{i_u}) = c(y_i, y_{i_u}) = C_{i,i_u}^t, \forall i \in [m]$ and $\forall i_u \in \mathcal{I}$, which implies that $R_i^s = R_i^t$. Then, we have $G_{i,i} = d(R_i^s, R_i^t) = 0$. We define π by $\pi_{i,j} = p_i$ if $i = j$, and 0 otherwise. Obviously, $M^{\mathbf{p}\mathbf{q}} \odot \pi$ is in $\Pi(\mathbf{p}, \mathbf{q}; M^{\mathbf{p}\mathbf{q}})$ and $\sum_{i,j} M_{i,j}^{\mathbf{p}\mathbf{q}} \pi_{i,j} (\alpha C_{i,j} + (1 - \alpha)G_{i,j}) = 0$. Therefore, $\mathcal{S}_{krk}(\mathbf{p}, \mathbf{q}) = 0$. (b) We denote π^* as the optimal solution of problem (A-21). If $\mathcal{S}_{krk}(\mathbf{p}, \mathbf{q}) = 0$, we have $\langle M^{\mathbf{p}\mathbf{q}} \odot \pi^*, C \rangle_F = 0$. This means that the KP problem $\min_{\pi \in \Pi(\mathbf{p}, \mathbf{q})} \langle \pi, C \rangle_F = 0$. Using the Proposition 2.2 in [62], we have $\mathbf{p} = \mathbf{q}$.

(2) **Show that** $\mathcal{S}_{krk}(\mathbf{p}, \mathbf{q}) = \mathcal{S}_{krk}(\mathbf{q}, \mathbf{p})$. From the definition of mask matrix in Proposition 1 in the paper, we have $M_{i,j}^{\mathbf{p}\mathbf{q}} = M_{j,i}^{\mathbf{q}\mathbf{p}}$. C and G are symmetric because c and d are distances. For any $\pi \in \Pi(\mathbf{p}, \mathbf{q}; M^{\mathbf{p}\mathbf{q}})$, we define π' as $\pi'_{i,j} = \pi_{j,i}$, and then $\pi' \in \Pi(\mathbf{q}, \mathbf{p}; M^{\mathbf{q}\mathbf{p}})$. Then, we have

$$\begin{aligned} \mathcal{S}_{krk}(\mathbf{p}, \mathbf{q}) &= \min_{\pi \in \Pi(\mathbf{p}, \mathbf{q}; M^{\mathbf{p}\mathbf{q}})} \sum_{i,j} M_{i,j}^{\mathbf{p}\mathbf{q}} \pi_{i,j} (\alpha C_{i,j} + (1-\alpha)G_{i,j}) \\ &= \min_{\pi \in \Pi(\mathbf{p}, \mathbf{q}; M^{\mathbf{p}\mathbf{q}})} \sum_{i,j} M_{j,i}^{\mathbf{q}\mathbf{p}} \pi_{i,j} (\alpha C_{j,i} + (1-\alpha)G_{j,i}) \\ &= \min_{\pi' \in \Pi(\mathbf{q}, \mathbf{p}; M^{\mathbf{q}\mathbf{p}})} \sum_{j,i} M_{j,i}^{\mathbf{q}\mathbf{p}} \pi'_{j,i} (\alpha C_{j,i} + (1-\alpha)G_{j,i}) \\ &= \mathcal{S}_{krk}(\mathbf{q}, \mathbf{p}). \end{aligned} \tag{A-22}$$

(3) **Show that** $\mathcal{S}_{krk}(\mathbf{p}, \mathbf{q}) \leq \mathcal{S}_{krk}(\mathbf{p}, \mathbf{r}) + \mathcal{S}_{krk}(\mathbf{r}, \mathbf{q})$. Let $M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}}$ and $M^{\mathbf{r}\mathbf{q}} \odot \pi^{\mathbf{r}\mathbf{q}}$ be the optimal transport plans corresponding to $\mathcal{S}_{krk}(\mathbf{p}, \mathbf{r})$ and $\mathcal{S}_{krk}(\mathbf{r}, \mathbf{q})$, respectively. We define

$$\tilde{\gamma} = (M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}}) \text{diag} \left(\frac{1}{\tilde{\mathbf{r}}} \right) (M^{\mathbf{r}\mathbf{q}} \odot \pi^{\mathbf{r}\mathbf{q}}), \tag{A-23}$$

where the element \tilde{r}_j of $\tilde{\mathbf{r}}$ is r_j if $r_j > 0$, and 1 otherwise. We notice that

$$\tilde{\gamma} \mathbb{1}_m = (M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}}) \text{diag} \left(\frac{1}{\tilde{\mathbf{r}}} \right) \mathbf{r} = (M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}}) \begin{pmatrix} \mathbf{r} \\ \tilde{\mathbf{r}} \end{pmatrix} = (M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}}) \tilde{\mathbb{1}}_m, \tag{A-24}$$

where the j -th location of $\tilde{\mathbb{1}}_m$ is 1 if $r_j > 0$, and 0 otherwise. Note that for j such that $r_j = 0$, we have $\sum_{i,j} (M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}})_{i,j} = r_j = 0$, which implies $(M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}})_{i,j} = 0$ for any i . Hence,

$$(M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}}) \tilde{\mathbb{1}}_m = (M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}}) \mathbb{1}_m = \mathbf{p}. \tag{A-25}$$

Similarity, $\tilde{\gamma}^\top \mathbb{1}_m = \mathbf{q}$. Since the indexes of paired keypoints across any two distribution in $\mathcal{P}_{\mathcal{I}}^{\mathcal{X}}$ are the same, for any $i_u \in \mathcal{I}$, the i_u -th row and column of $M^{\mathbf{p}\mathbf{r}}$ and $M^{\mathbf{r}\mathbf{q}}$ are zeros except for that $M_{i_u, i_u}^{\mathbf{p}\mathbf{r}} = M_{i_u, i_u}^{\mathbf{r}\mathbf{q}} = 1$. So the i_u -th row and column of $\tilde{\gamma}$ are zeros except for $\tilde{\gamma}_{i_u, i_u}$. Then, we can write $\tilde{\gamma}_{i_u, i_u} = M^{\mathbf{p}\mathbf{q}} \odot \gamma$ with $\gamma \in \mathbb{R}_+^{m \times m}$. Further, we have $\gamma \in \Pi(\mathbf{p}, \mathbf{q}; M^{\mathbf{p}\mathbf{q}})$. The triangle inequality follows then from

$$\begin{aligned} \mathcal{S}_{krk}(\mathbf{p}, \mathbf{q}) &= \min_{\pi \in \Pi(\mathbf{p}, \mathbf{q}; M^{\mathbf{p}\mathbf{q}})} \sum_{i,j} M_{i,j}^{\mathbf{p}\mathbf{q}} \pi_{i,j} (\alpha c(x_i, y_j) + (1-\alpha)d(R_i^s, R_j^t)) \\ &\leq \sum_{i,j} \tilde{\gamma}_{i,j} (\alpha c(x_i, y_j) + (1-\alpha)d(R_i^s, R_j^t)) \\ &= \sum_{i,j} (\alpha c(x_i, y_j) + (1-\alpha)d(R_i^s, R_j^t)) \sum_k \frac{(M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}})_{i,k} (M^{\mathbf{r}\mathbf{q}} \odot \pi^{\mathbf{r}\mathbf{q}})_{k,j}}{\tilde{r}_k} \\ &\leq \sum_{i,k,j} (\alpha(c(x_i, z_k) + c(z_k, y_j)) + (1-\alpha)(d(R_i^s, R_k^r) + d(R_k^r, R_j^t))) \frac{(M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}})_{i,k} (M^{\mathbf{r}\mathbf{q}} \odot \pi^{\mathbf{r}\mathbf{q}})_{k,j}}{\tilde{r}_k} \\ &= \sum_{i,k,j} (\alpha c(x_i, z_k) + (1-\alpha)d(R_i^s, R_k^r)) \frac{(M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}})_{i,k} (M^{\mathbf{r}\mathbf{q}} \odot \pi^{\mathbf{r}\mathbf{q}})_{k,j}}{\tilde{r}_k} \\ &\quad + \sum_{i,k,j} (\alpha c(z_k, y_j) + (1-\alpha)d(R_k^r, R_j^t)) \frac{(M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}})_{i,k} (M^{\mathbf{r}\mathbf{q}} \odot \pi^{\mathbf{r}\mathbf{q}})_{k,j}}{\tilde{r}_k} \\ &= \sum_{i,k} (\alpha c(x_i, z_k) + (1-\alpha)d(R_i^s, R_k^r)) (M^{\mathbf{p}\mathbf{r}} \odot \pi^{\mathbf{p}\mathbf{r}})_{i,k} \\ &\quad + \sum_{k,j} (\alpha c(z_k, y_j) + (1-\alpha)d(R_k^r, R_j^t)) (M^{\mathbf{r}\mathbf{q}} \odot \pi^{\mathbf{r}\mathbf{q}})_{k,j} \\ &= \mathcal{S}_{krk}(\mathbf{p}, \mathbf{r}) + \mathcal{S}_{krk}(\mathbf{r}, \mathbf{q}), \end{aligned} \tag{A-26}$$

where z_k is the support point of \mathbf{r} and R_k^r is the relation of z_k to the keypoints of \mathbf{r} . \square

A.5.2 KPG-RL-GW Providing a Divergence

For any distribution $\mathbf{p} \in \mathcal{P}_{\mathcal{I}}^{\mathcal{X}}$ and $\mathbf{q} \in \mathcal{P}_{\mathcal{I}}^{\mathcal{Y}}$, \mathbf{p} and \mathbf{q} are said to be isomorphic if there exists a bijection $\sigma : [m] \mapsto [m]$ such that $c(x_i, x_k) = c'(y_{\sigma(i)}, y_{\sigma(k)})$, and $p_i = q_{\sigma(i)}$, where $[m] = \{1, 2, \dots, m\}$, and c and c' are respectively proper distances in \mathcal{X} and \mathcal{Y} . The keypoints are ‘‘correct’’ means that if $\mathbf{p} = \mathbf{q}$, σ maps each source keypoint to its paired target keypoint, *i.e.*, $\sigma(i_u) = i_u$, for any $i_u \in \mathcal{I}$. We denote

$$\begin{aligned} \mathcal{S}_{krq}(\mathbf{p}, \mathbf{q}) = \min_{\pi \in \Pi(\mathbf{p}, \mathbf{q}; M^{\mathbf{p}\mathbf{q}})} \sum_{i,j} \left[\alpha \left(\sum_{k,l} (M^{\mathbf{p}\mathbf{q}} \odot \pi)_{i,j} (M^{\mathbf{p}\mathbf{q}} \odot \pi)_{k,l} |C_{i,k}^s - C_{j,l}^t|^2 \right) \right. \\ \left. + (1 - \alpha)(M^{\mathbf{p}\mathbf{q}} \odot \pi)_{i,j} G_{i,j} \right]. \end{aligned} \quad (\text{A-27})$$

Theorem A-2. *Suppose c and c' are proper distances in spaces \mathcal{X} and \mathcal{Y} . Suppose d is a divergence in probability simplex Σ_m . Then, for any \mathbf{p} in $\mathcal{P}_{\mathcal{I}}^{\mathcal{X}}$ and any \mathbf{q} in $\mathcal{P}_{\mathcal{I}}^{\mathcal{Y}}$, given the ‘‘correct’’ paired keypoints stated above, $\mathcal{S}_{krq}(\mathbf{p}, \mathbf{q}) = 0$ if and only if \mathbf{p} and \mathbf{q} are isomorphic.*

Proof:

(a) If \mathbf{p} and \mathbf{q} are isomorphic, for any $i \in [m]$ and any $i_u \in \mathcal{I}$, we have $c(x_i, x_{i_u}) = c'(y_{\sigma(i)}, y_{\sigma(i_u)}) = c'(y_{\sigma(i)}, y_{i_u})$, implying that $R_i^s = R_{\sigma(i)}^t$. We define π as $\pi_{i,j} = p_i$ if $j = \sigma(i)$, and 0 otherwise. We then have

$$\begin{aligned} & \sum_{i,j} \alpha \left(\sum_{k,l} (M^{\mathbf{p}\mathbf{q}} \odot \pi)_{i,j} (M^{\mathbf{p}\mathbf{q}} \odot \pi)_{k,l} |C_{i,k}^s - C_{j,l}^t|^2 \right) + (1 - \alpha)(M^{\mathbf{p}\mathbf{q}} \odot \pi)_{i,j} G_{i,j} \\ &= \sum_i \left[\alpha \left(\sum_k (M^{\mathbf{p}\mathbf{q}} \odot \pi)_{i,\sigma(i)} (M^{\mathbf{p}\mathbf{q}} \odot \pi)_{k,\sigma(k)} |C_{i,k}^s - C_{\sigma(i),\sigma(k)}^t|^2 \right) \right. \\ & \quad \left. + (1 - \alpha)(M^{\mathbf{p}\mathbf{q}} \odot \pi)_{i,\sigma(i)} G_{i,\sigma(i)} \right] \\ &= \sum_i \left[\alpha \left(\sum_k (M^{\mathbf{p}\mathbf{q}} \odot \pi)_{i,\sigma(i)} (M^{\mathbf{p}\mathbf{q}} \odot \pi)_{k,\sigma(k)} |c(x_i, x_k) - c'(x_{\sigma(i)}, x_{\sigma(k)})|^2 \right) \right. \\ & \quad \left. + (1 - \alpha)(M^{\mathbf{p}\mathbf{q}} \odot \pi)_{i,\sigma(i)} d(R_i^s, R_{\sigma(i)}^t) \right] \\ &= 0. \end{aligned} \quad (\text{A-28})$$

This implies $\mathcal{S}_{krq}(\mathbf{p}, \mathbf{q}) = 0$.

(b) Let $(M^{\mathbf{p}\mathbf{q}} \odot \pi^*)$ be the optimal transport plan corresponding to $\mathcal{S}_{krq}(\mathbf{p}, \mathbf{q}) = 0$. If $\mathcal{S}_{krq}(\mathbf{p}, \mathbf{q}) = 0$, we have

$$\sum_{i,j,k,l} (M^{\mathbf{p}\mathbf{q}} \odot \pi^*)_{i,j} (M^{\mathbf{p}\mathbf{q}} \odot \pi^*)_{k,l} |C_{i,k}^s - C_{j,l}^t|^2 = 0. \quad (\text{A-29})$$

This indicates that the Gromov-Wasserstein distance

$$\min_{\pi \in \Pi(\mathbf{p}, \mathbf{q})} \sum_{i,j,k,l} \pi_{i,j} \pi_{k,l} |C_{i,k}^s - C_{j,l}^t|^2 = 0. \quad (\text{A-30})$$

By virtue to Gromov-Wasserstein properties in [27], there exists a bijection $\sigma : [m] \mapsto [m]$ such that $c(x_i, x_k) = c'(y_{\sigma(i)}, y_{\sigma(k)})$, and $p_i = q_{\sigma(i)}$. \square

From Theorem A-2, the KPG-RL-GW model provides a divergence in the sense of isomorphism.

A.6 Motivations of the Solving Algorithm for Partial-KPG-RL Model

In the partial-KPG-RL model in Eq. (13) in the paper, only s -unit mass of source and target distributions is matched. Inspired by [20], we add a dummy point with mass $\|\mathbf{q}\|_1 - s$ for source domain (the left black circle in Fig. A-2) and a dummy point with mass $\|\mathbf{p}\|_1 - s$ for target domain (the right black circle in Fig. A-2). We denote $\bar{\mathbf{p}} = (\mathbf{p}^\top, \|\mathbf{q}\|_1 - s)^\top$ and $\bar{\mathbf{q}} = (\mathbf{q}^\top, \|\mathbf{p}\|_1 - s)^\top$. As illustrated in Fig. A-2, we aim to design the extended guiding matrix \bar{G} and extended mask matrix \bar{M} such that performing KPG-RL between $\bar{\mathbf{p}}$ and $\bar{\mathbf{q}}$ will transport $\|\mathbf{p}\|_1 - s$ mass from source real data points to the target dummy point and transport $\|\mathbf{q}\|_1 - s$ mass from source dummy point to target real data

points. As a sequence, only s mass of source and target real data points are matched. Meanwhile, the keypoints should not be matched to the dummy points because they are annotated data to guide the matching. To do this, we extend G, M by

$$\bar{G} = \begin{bmatrix} G & \xi \mathbb{1}_n \\ \xi \mathbb{1}_m^\top & 2\xi + A \end{bmatrix}, \bar{M} = \begin{bmatrix} M & \mathbf{a} \\ \mathbf{b}^\top & 1 \end{bmatrix},$$

where $A > 0, \xi > 0, \mathbf{a} \in \mathbb{R}^m, \mathbf{b} \in \mathbb{R}^n$, and M is constructed as in Proposition A-1. The element a_i of \mathbf{a} is 0 if $i \in \mathcal{I}$, and 1 otherwise. The element b_j of \mathbf{b} is 0 if $j \in \mathcal{J}$, and 1 otherwise. By Theorem A-3, solving partial-KPG-RL model boils down to solving the KPG-RL-like problem $\min_{\bar{\pi} \in \Pi(\bar{\mathbf{p}}, \bar{\mathbf{q}}, \bar{M})} \langle \bar{M} \odot \bar{\pi}, \bar{G} \rangle_F$.

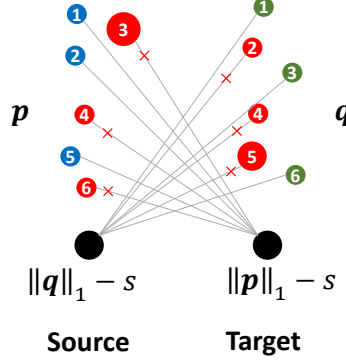


Figure A-2: Illustration of dummy points (black circles) for source and target domains.

A.7 Proof of Theorem 1

For the convenience of understanding, the Theorem 1 in the paper is for the case that $p_i = q_j$ for all $(i, j) \in \mathcal{K}$. In this appendix, we provide and prove the Theorem A-3, a generalization of Theorem 1 in the paper, for the general case that there could exist some $(i, j) \in \mathcal{K}$ such that $p_i \neq q_j$. Before that, we rewrite the partial-KPG-RL model first.

Partial-KPG-RL model:

$$\min_{\pi \in \Pi^s(\mathbf{p}, \mathbf{q}; M)} \{L_{kpg}(M \odot \pi) = \langle M \odot \pi, G \rangle_F\}, \quad (\text{A-31})$$

where $\Pi^s(\mathbf{p}, \mathbf{q}; M) = \{\pi \in \mathbb{R}_+^{m \times n} | (M \odot \pi) \mathbb{1}_n \leq \mathbf{p}, (M \odot \pi)^\top \mathbb{1}_m \leq \mathbf{q}, \mathbb{1}_m^\top (M \odot \pi) \mathbb{1}_n = s; (M \odot \pi)_{i,:} \mathbb{1}_n = p_i, \forall i \in \mathcal{I}; \mathbb{1}_m^\top (M \odot \pi)_{:,j} = q_j, \forall j \in \mathcal{J}\}$.

Theorem A-3. Suppose $A > 0, \xi > 0, \sum_{i \in \mathcal{I}} p_i + \max\{q_{\kappa(i)} - p_i, 0\} < s$, and $\sum_{j \in \mathcal{J}} q_j + \max\{p_{\kappa'(j)} - q_j, 0\} < s$, then the optimal transport plan $M \odot \pi^*$ of partial-KPG-RL model is the m -by- n block in the upper left corner of the optimal transport plan $M \odot \bar{\pi}^*$ of problem $\min_{\bar{\pi} \in \Pi(\bar{\mathbf{p}}, \bar{\mathbf{q}}, \bar{M})} \langle \bar{M} \odot \bar{\pi}, \bar{G} \rangle_F$.

The definitions of $\kappa(i)$ and $\kappa'(j)$ are given in Appendix A.1. The condition $\sum_{i \in \mathcal{I}} p_i + \max\{q_{\kappa(i)} - p_i, 0\} < s$ implies that the sum of the mass ($\sum_{i \in \mathcal{I}} p_i$) of source keypoints and the mass ($\sum_{i \in \mathcal{I}} \max\{q_{\kappa(i)} - p_i, 0\}$) of the other source points apart from keypoints that should be transported to target keypoints is less than s . The condition $\sum_{j \in \mathcal{J}} q_j + \max\{p_{\kappa'(j)} - q_j, 0\} < s$ implies that the sum of the mass ($\sum_{j \in \mathcal{J}} q_j$) of target keypoints and the mass ($\sum_{j \in \mathcal{J}} \max\{p_{\kappa'(j)} - q_j, 0\}$) of the other target points apart from keypoints received from source keypoints is less than s . The two conditions are reasonable to guarantee the admissible solutions of problem (A-31). If $p_i = q_j$ for all $(i, j) \in \mathcal{K}$ (i.e., $p_i = q_{\kappa(i)}, \forall i \in \mathcal{I}$ and $q_j = p_{\kappa'(j)}, \forall j \in \mathcal{J}$), Theorem A-3 degenerates to the Theorem 1 in the paper.

A.7.1 Proof

We denote $\bar{\pi} = \bar{\pi}_{1:m, 1:n}^*$ and $t = \bar{\pi}_{m+1, n+1}^*$. To prove Theorem A-3, we first give some preparations, and then conduct the following three steps. In step 1, we show that $t = 0$. In step 2, we show that

$\ddot{\pi} \in \Pi^s(\mathbf{p}, \mathbf{q}; M)$ which means that $\ddot{\pi}$ is a feasible solution of problem (A-31). In step 3, we show that $M \odot \ddot{\pi}$ is the optimal transport plan of problem (A-31). We next detail these steps.

Preparations:

Since $\bar{\pi}^* \in \Pi(\bar{\mathbf{p}}, \bar{\mathbf{q}}; \bar{M})$ and $t = \bar{\pi}_{m+1, n+1}^*$, we have

$$\begin{aligned} \mathbb{1}_{m+1}^\top (\bar{M} \odot \bar{\pi}^*) \mathbb{1}_{n+1} &= [\mathbb{1}_m^\top \quad 1] \left(\begin{bmatrix} M & \mathbf{a} \\ \mathbf{b} & 1 \end{bmatrix} \odot \begin{bmatrix} \ddot{\pi} & \bar{\pi}_{1:m, n+1}^* \\ \bar{\pi}_{m+1, 1:n}^* & \bar{\pi}_{m+1, n+1}^* \end{bmatrix} \right) \begin{bmatrix} \mathbb{1}_n \\ 1 \end{bmatrix} \\ &= \mathbb{1}_m^\top (M \odot \ddot{\pi}) \mathbb{1}_n + \sum_{i=1}^m a_i \bar{\pi}_{i, n+1}^* + \sum_{j=1}^n b_j \bar{\pi}_{m+1, j}^* + t. \end{aligned} \quad (\text{A-32})$$

Meanwhile, we have

$$\mathbb{1}_{m+1}^\top (\bar{M} \odot \bar{\pi}^*) \mathbb{1}_{n+1} = \mathbb{1}_{m+1}^\top \bar{\mathbf{p}} = \|\bar{\mathbf{p}}\|_1 = \|\mathbf{p}\|_1 + \|\mathbf{q}\|_1 - s, \quad (\text{A-33})$$

$$\sum_{i=1}^m a_i \bar{\pi}_{i, n+1}^* + t = \|\mathbf{p}\|_1 - s, \quad (\text{A-34})$$

and

$$\sum_{j=1}^n b_j \bar{\pi}_{m+1, j}^* + t = \|\mathbf{q}\|_1 - s. \quad (\text{A-35})$$

Combining the above four equations, we have

$$\mathbb{1}_m^\top (M \odot \ddot{\pi}) \mathbb{1}_n + \|\mathbf{p}\|_1 + \|\mathbf{q}\|_1 - 2s - t = \|\mathbf{p}\|_1 + \|\mathbf{q}\|_1 - s. \quad (\text{A-36})$$

Therefore,

$$\mathbb{1}_m^\top (M \odot \ddot{\pi}) \mathbb{1}_n = s + t. \quad (\text{A-37})$$

Step 1: show that $t = \bar{\pi}_{m+1, n+1}^* = 0$.

First, we have

$$\begin{aligned} \langle \bar{M} \odot \bar{\pi}^*, \bar{G} \rangle_F &= \sum_{i=1}^m \sum_{j=1}^n M_{i,j} \bar{\pi}_{i,j}^* G_{i,j} + \xi \sum_{i=1}^m a_i \bar{\pi}_{i, n+1}^* \\ &\quad + \xi \sum_{j=1}^n b_j \bar{\pi}_{m+1, j}^* + (2\xi + A) \bar{\pi}_{m+1, n+1}^* \\ &= \sum_{i=1}^m \sum_{j=1}^n M_{i,j} \bar{\pi}_{i,j}^* G_{i,j} + \xi (\|\mathbf{p}\|_1 + \|\mathbf{q}\|_1 - 2s - 2t) + (2\xi + A)t \\ &= \sum_{i=1}^m \sum_{j=1}^n M_{i,j} \bar{\pi}_{i,j}^* G_{i,j} + \xi (\|\mathbf{p}\|_1 + \|\mathbf{q}\|_1 - 2s) + At. \end{aligned} \quad (\text{A-38})$$

Suppose $\bar{\pi}_{m+1, n+1}^* > 0$, we next construct a solution γ such that $\gamma_{m+1, n+1} = 0$ and leads to conflict. We randomly select a set $S = \{(i, j) | \bar{\pi}_{i,j}^* > 0, i \leq m, j \leq n, i \notin \mathcal{I}, j \notin \mathcal{J}\}$ and a index pair (i_0, j_0) satisfying the constraints of elements in S , such that $\sum_{(i,j) \in S} \bar{\pi}_{i,j}^* \leq t$ and $\sum_{(i,j) \in S} \bar{\pi}_{i,j}^* + \bar{\pi}_{i_0, j_0}^* > t$. In the rest part of this section, the involved i, j satisfy $i \leq m$ and $j \leq n$. Such non-empty S and

(i_0, j_0) always exist, because

$$\begin{aligned}
\mathbb{1}_m^\top(M \odot \bar{\pi})\mathbb{1}_n &= \sum_{i=1}^m \sum_{j=1}^n M_{i,j} \bar{\pi}_{i,j}^* = \sum_{i \in \mathcal{I}, j} \bar{\pi}_{i,j}^* + \sum_{i \notin \mathcal{I}, j \in \mathcal{J}} M_{i,j} \bar{\pi}_{i,j}^* + \sum_{i \notin \mathcal{I}, j \notin \mathcal{J}} \bar{\pi}_{i,j}^* \\
&= \sum_{i \in \mathcal{I}, j} \bar{\pi}_{i,j}^* + \sum_{i \notin \mathcal{I}, j \in \mathcal{J}} M_{i,j} \bar{\pi}_{i,j}^* + \sum_{i \notin \mathcal{I}, j \notin \mathcal{J}} \bar{\pi}_{i,j}^* \\
&= \sum_{i \in \mathcal{I}, j} \bar{\pi}_{i,j}^* + \sum_{i \notin \mathcal{I}, i' \in \mathcal{I}} M_{i, \kappa(i')} \bar{\pi}_{i, \kappa(i')}^* + \sum_{i \notin \mathcal{I}, j \notin \mathcal{J}} \bar{\pi}_{i,j}^* \\
&\leq \sum_{i \in \mathcal{I}} p_i + \sum_{i' \in \mathcal{I}, i' \neq i} M_{i, \kappa(i')} \bar{\pi}_{i, \kappa(i')}^* + \sum_{i \notin \mathcal{I}, j \notin \mathcal{J}} \bar{\pi}_{i,j}^* \\
&= \sum_{i \in \mathcal{I}} p_i + \sum_{i' \in \mathcal{I}} \max\{q_{\kappa(i')} - p_{i'}, 0\} + \sum_{i \notin \mathcal{I}, j \notin \mathcal{J}} \bar{\pi}_{i,j}^*,
\end{aligned} \tag{A-39}$$

$\mathbb{1}_m^\top(M \odot \bar{\pi})\mathbb{1}_n = s + t$, and $\sum_{i \in \mathcal{I}} p_i + \max\{q_{\kappa(i)} - p_i, 0\} < s$, we have $\sum_{i \notin \mathcal{I}, j \notin \mathcal{J}} \bar{\pi}_{i,j}^* > t$. We now move the mass of index pairs in S and (i_0, j_0) to their marginal such that a total mass of t is moved. Specifically, for $(i, j) \in S$, we set $\gamma_{i,j} = 0, \gamma_{i, n+1} = \bar{\pi}_{i, n+1}^* + \bar{\pi}_{i,j}^*, \gamma_{m+1, j} = \bar{\pi}_{m+1, j}^* + \bar{\pi}_{i,j}^*$. For (i_0, j_0) , we set $\gamma_{i_0, j_0} = \bar{\pi}_{i_0, j_0}^* - (t - \sum_{i \notin \mathcal{I}, j \notin \mathcal{J}} \bar{\pi}_{i,j}^*), \gamma_{i_0, n+1} = \bar{\pi}_{i_0, n+1}^* + (t - \sum_{i \notin \mathcal{I}, j \notin \mathcal{J}} \bar{\pi}_{i,j}^*), \gamma_{m+1, j_0} = \bar{\pi}_{m+1, j_0}^* - (t - \sum_{i \notin \mathcal{I}, j \notin \mathcal{J}} \bar{\pi}_{i,j}^*)$. For $(i, j) \notin S$, we set $\gamma_{i,j} = \bar{\pi}_{i,j}^*, \gamma_{i, n+1} = \bar{\pi}_{i, n+1}^*, \gamma_{m+1, j} = \bar{\pi}_{m+1, j}^*$. It is easy to verify that $\gamma \in \Pi(\bar{\mathbf{p}}, \bar{\mathbf{q}}; \bar{M})$. Similar to Eq. (A-38), we have

$$\langle \bar{M} \odot \gamma, \bar{G} \rangle_F = \sum_{i=1}^m \sum_{j=1}^n M_{i,j} \gamma_{i,j} G_{i,j} + \xi(\|\mathbf{p}\|_1 + \|\mathbf{q}\|_1 - 2s). \tag{A-40}$$

Using the optimality of $\bar{M} \odot \bar{\pi}^*$, we have

$$\langle \bar{M} \odot \gamma, \bar{G} \rangle_F - \langle \bar{M} \odot \bar{\pi}^*, \bar{G} \rangle_F = \sum_{i=1}^m \sum_{j=1}^n M_{i,j} (\gamma_{i,j} - \bar{\pi}_{i,j}^*) G_{i,j} - At > 0. \tag{A-41}$$

From the definition of γ , we can see that $\gamma_{i,j} \leq \bar{\pi}_{i,j}^*$, and thus $\sum_{i=1}^m \sum_{j=1}^n M_{i,j} (\gamma_{i,j} - \bar{\pi}_{i,j}^*) G_{i,j} \leq 0$. Hence, from Eq. (A-41), we have $A < 0$, contradicting the assumption that $A > 0$. Therefore, $t = \bar{\pi}_{m+1, n+1}^* = 0$ holds.

Step 2: show that $\bar{\pi}$ is a feasible solution of problem in Eq. (A-31).

We verify the constraints as follows.

(1) Since $\bar{\pi}^* \geq 0$, we have $\bar{\pi} \geq 0$.

(2) $(\bar{M} \odot \bar{\pi}^*)\mathbb{1}_{n+1} = \begin{bmatrix} M \odot \bar{\pi}^* & \mathbf{a} \odot \bar{\pi}_{1:m, n+1}^* \\ \mathbf{b} \odot \bar{\pi}_{m+1, 1:n}^* & 0 \end{bmatrix} \begin{bmatrix} \mathbb{1}_n \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{p} \\ \|\mathbf{q}\|_1 - s \end{bmatrix}$, then $(M \odot \bar{\pi})\mathbb{1}_n + \mathbf{a} \odot \bar{\pi}_{1:m, n+1}^* = \mathbf{p}$, and $(M \odot \bar{\pi})\mathbb{1}_n \leq \mathbf{p}$.

(3) Similarly, from $\mathbb{1}_{m+1}^\top(\bar{M} \odot \bar{\pi}^*) = (\mathbf{q}, \|\mathbf{q}\|_1 - s)^\top$, we have $\mathbb{1}_m^\top(M \odot \bar{\pi}) \leq \mathbf{q}$.

(4) $\mathbb{1}_m^\top(M \odot \bar{\pi})\mathbb{1}_n = s$ holds, because $t = 0$ as in Step 1.

(5) $\forall i \in \mathcal{I}, (\bar{M} \odot \bar{\pi}^*)_{i, :} \mathbb{1}_{n+1} = (M \odot \bar{\pi})_{i, :} \mathbb{1}_n + a_i \bar{\pi}_{i, n+1}^* = p_i$. Since $a_i = 0$, we have $(M \odot \bar{\pi})_{i, :} \mathbb{1}_n = p_i$.

(6) $\forall j \in \mathcal{J}, \mathbb{1}_{m+1}^\top(\bar{M} \odot \bar{\pi}^*)_{:, j} = \mathbb{1}_m^\top(M \odot \bar{\pi})_{:, j} + b_j \bar{\pi}_{m+1, j}^* = q_j$. Since $b_j = 0$, we have $\mathbb{1}_m^\top(M \odot \bar{\pi})_{:, j} = q_j$.

Therefore, we have $\bar{\pi} \in \Pi^s(\mathbf{p}, \mathbf{q}; M)$, and $\bar{\pi}$ is a feasible solution of problem in Eq. (A-31).

Step 3: show that $M \odot \bar{\pi}$ is the optimal transport plan of problem in Eq. (A-31).

Suppose there exist a transport plan $M \odot \gamma$ with $\gamma \in \Pi^s(\mathbf{p}, \mathbf{q}; M)$ such that

$$\sum_{i=1}^m \sum_{j=1}^n M_{i,j} \gamma_{i,j} G_{i,j} < \sum_{i=1}^m \sum_{j=1}^n M_{i,j} \bar{\pi}_{i,j} G_{i,j}.$$

We construct $\bar{\gamma}$ as follows. For $i \leq m, j \leq n$, $\bar{\gamma}_{i,j} = \gamma_{i,j}$. $\bar{\gamma}_{i,n+1} = p_i - \sum_{j=1}^n \gamma_{i,j}, \forall i \leq m$. $\bar{\gamma}_{m+1,j} = q_j - \sum_{i=1}^m \gamma_{i,j}, \forall j \leq n$. $\bar{\gamma}_{m+1,n+1} = 0$. Easily, we can verify that $\bar{\gamma}$ is in $\Pi(\bar{\mathbf{p}}, \bar{\mathbf{q}}; \bar{M})$. Meanwhile,

$$\begin{aligned} \langle \bar{M} \odot \bar{\gamma}, \bar{G} \rangle_F &= \sum_{i=1}^m \sum_{j=1}^n M_{i,j} \gamma_{i,j} C_{i,j} + \xi(\|\mathbf{p}\|_1 + \|\mathbf{q}\|_1 - 2s) \\ &< \sum_{i=1}^m \sum_{j=1}^n M_{i,j} \bar{\pi}_{i,j} G_{i,j} + \xi(\|\mathbf{p}\|_1 + \|\mathbf{q}\|_1 - 2s) \\ &= \langle \bar{M} \odot \bar{\pi}^*, \bar{G} \rangle_F. \end{aligned} \tag{A-42}$$

This contradicts the fact that $\bar{M} \odot \bar{\pi}^*$ is the optimal transport plan of problem $\min_{\bar{\pi} \in \Pi(\bar{\mathbf{p}}, \bar{\mathbf{q}}; \bar{M})} \langle \bar{M} \odot \bar{\pi}, \bar{G} \rangle_F$. Therefore, $M \odot \bar{\pi}$ is the optimal transport plan of problem in Eq. (A-31). \square

B Additional Experimental Details and Results

B.1 Toy Experiment for Evaluating Partial-KPG-RL model

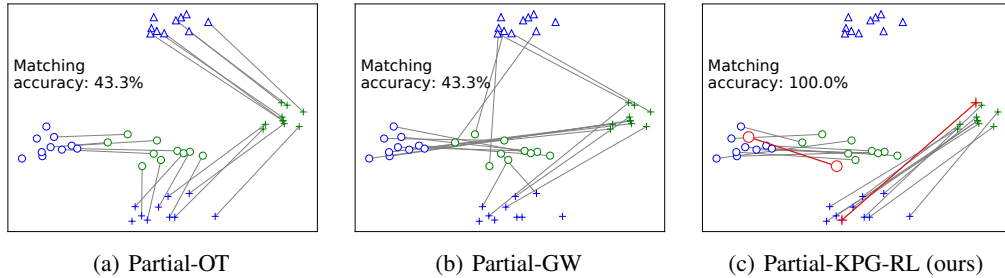


Figure A-3: Matching produced by (a) partial-OT model, (b) partial-GW model, and (c) our proposed partial-KPG-RL model.

Figure A-3 illustrates the toy data experiment for evaluating the partial-KPG-RL model. In Fig. A-3, the source (blue) and target (green) distributions are Gaussian mixtures. The source (resp. target) distribution is composed of three (resp. two) distinct Gaussian components indicated by different shapes where the same shapes indicate the points of the same class. When conducting OT, the source class data represented by “ Δ ” should not be transported. In Figs. A-3(a) and A-3(b), we can observe that both the partial-OT model (defined in Eq. (2) in the paper) and the partial-GW model [20] wrongly transport some source points of class “ Δ ” to target domain and lead to low matching accuracy. With the guidance of a few keypoints (red pairs), our proposed partial-KPG-RL model does not transport the source points of class “ Δ ” to target domain and apparently improves the matching accuracy as in Fig. A-3(c).

B.2 Additional Experimental Details and Results for Open-set HDA

More experimental details. In open-set HDA, we are given a large amount of labeled source domain data $\{(x_i, t_i)\}_{i=1}^m$, a few labeled target domain data $\{y_j, \bar{t}_j\}_{j=1}^{n_l}$, and a large number of unlabeled target domain data $\{y_j\}_{j=n_l+1}^n$. The fraction of unknown class data is η . To apply the partial-KPG-RL model defined in Eq. (13) in the paper to open-set HDA, for each labeled target domain data, we take its corresponding source class center to construct a keypoint pair. We then resample the source domain data such that the total number of resampled source domain data and the source keypoints is $m' = (1 - \eta)n$. We define the source distribution as $\mathbf{p} = \frac{1-\eta}{m'} (\sum_{j=1}^{n_l} \delta_{c_j} + \sum_{j=n_l+1}^{m'} \delta_{x'_j})$, where x'_j is a resampled source domain sample and c_j is the source class center corresponding to the target labeled sample y_j . The target distribution is defined as $\mathbf{q} = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$.

The partial-KPG-RL model is conducted to transport mass from p to q with $s = 1 - \eta$. After transport, the η -fraction unlabeled target data receiving smallest mass from source domain are detected as unknown class and the rest unlabeled target data are taken as common class ones. Finally, we train the kernel SVM on the transported source domain data and labeled target domain data to classify the unlabeled target domain common class data.

Table A-1: Results on Office-31 for open-set HDA with unknown η . $\hat{\eta}$ is the estimate of η (the true $\eta = 0.67$).

Method	A→A ($\hat{\eta} = 0.57$)			A→D ($\hat{\eta} = 0.48$)			A→W ($\hat{\eta} = 0.62$)			D→A ($\hat{\eta} = 0.57$)			D→D ($\hat{\eta} = 0.48$)		
	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS
Baseline	38.2	61.9	47.2	20.0	69.3	31.0	28.2	80.1	41.7	38.2	61.9	47.2	20.0	69.3	31.0
Partial-KPG-RL	49.1	70.1	57.8	61.8	59.3	60.5	54.5	73.2	62.5	59.1	73.6	65.5	83.6	66.7	74.2

Method	D→W ($\hat{\eta} = 0.62$)			W→A ($\hat{\eta} = 0.57$)			W→D ($\hat{\eta} = 0.48$)			W→W ($\hat{\eta} = 0.62$)			Avg		
	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS	OS*	UNK	HOS
Baseline	28.2	80.1	41.7	38.2	61.9	47.2	20.0	69.3	31.0	28.2	80.1	41.7	28.8	70.4	40.0
Partial-KPG-RL	78.2	87.4	82.6	60.9	74.5	67.0	81.8	66.7	73.5	78.2	87.0	82.4	67.5	73.2	69.5

Results for open-set HDA with unknown η . For the more practical open-set HDA setting that η is unknown, researchers can design methods to estimate η and then apply our method using the estimate of η , or take η as a hyper-parameter and design methods to tune it. We directly use the positive-unlabeled learning [63] method [64] to estimate the fraction of common class data among the target domain unlabeled data, by taking the labeled target data as positive samples. The results of different methods for open-set HDA using the estimate $\hat{\eta}$ of η are given in Table A-1. According to Table A-1, the positive transfer is achieved by our method. We can see that $\hat{\eta}$ in all tasks is lower than the true η , implying that less unknown class samples are detected. Correspondingly, the UNK value (73.2%) achieved by partial-KPG-RL using $\hat{\eta}$ in Table A-1 is smaller than that (83.5%) using η in Table 3 in the paper. Surprisingly, the OS* value (67.5%) of partial-KPG-RL in Table A-1 is higher than that (59.7%) in Table 3 in the paper. As a balance, the HOS value (69.5%) achieved by partial-KPG-RL using $\hat{\eta}$ is similar to the HOS value (69.1%) of partial-KPG-RL using the true η .

In the following Table A-2, we take η as a hyper-parameter and show the average HOS achieved by partial-KPG-RL using varying magnitude of η . It is observed that the average HOS is stable to η in a relatively large range of [0.50, 0.80].

Table A-2: Average HOS of partial-KPG-RL using varying magnitude of η (the unknown true value of η is 0.67).

η	0.50	0.55	0.60	0.65	0.70	0.75	0.80
Average HOS	67.2	69.2	69.9	68.7	69.2	68.5	65.5

B.3 Application in Deep Unsupervised Domain Adaptation

In this section, we apply our method to deep unsupervised domain adaptation where the mini-batch-based implementation is required. The main challenge is that some of the samples in the mini-batch may not be matched. For instance, the categories of some samples in the source mini-batch may not be present in the target mini-batch, and thus these source samples should not be transported/matched. Inspired by [65] that uses partial OT over the mini-batch data to implement deepJDOT [66], we use our partial KPG-RL-KP model to partially match the mini-batch data in the training of the deep network. The partial KPG-RL-KP model is modified from Eq. (13) by replacing G by $\alpha C + (1 - \alpha)G$. As an experimental example, we apply the partial KPG-RL-KP to the unsupervised domain adaptation experiment on the Office-Home dataset [67]. We take the source and target class centers of the same class as a keypoint pair. The centers are online updated by exponential moving average in training, same as in [68]. We use the pseudo labels of target data to update the target class centers, due to the

lack of target labels. The protocol is the same as that in [65]. The batch size is set to 65 and the total transport mass (s in Eq. (13)) is set to 0.6, which are the same as those in [65]. The results are reported in the following Table A-3.

Table A-3: Results for unsupervised domain adaptation. “A”, “C”, “P”, and “R” are the domains of “Art”, “Clipart”, “Product”, and “RealWorld” in Office-Home dataset.

Method	A→C	A→P	A→R	C→A	C→P	C→R	P→A	P→C	P→R	R→A	R→C	R→P	Avg
ROT [23]	47.20	71.80	76.40	58.60	68.10	70.20	56.50	45.00	75.80	69.40	52.10	80.60	64.30
m-OT [66]	51.75	70.01	75.79	59.60	66.46	70.07	57.60	47.88	75.29	66.82	55.71	78.11	64.59
m-UOT [69]	54.99	74.45	80.78	65.66	74.93	74.91	64.70	53.42	80.01	74.58	59.88	83.73	70.17
m-POT [65]	55.65	73.80	80.76	66.34	74.88	76.16	64.46	53.38	80.60	74.55	59.71	83.81	70.34
m-KPG-RL-KP	52.13	63.65	74.53	61.12	67.84	67.88	59.84	52.93	76.90	71.92	59.21	82.55	65.88
m-PKPG-RL-KP	57.96	74.45	78.75	66.30	75.22	74.39	66.87	58.47	80.47	75.15	61.15	84.23	71.12

In Table A-3, ROT [23] is a robust OT method. m-OT is the direct mini-batch implementation of deepJDOT [66]. m-UOT [69] and m-POT [65] are respectively unbalanced deepJDOT and partial deepJDOT on mini-batch data. m-KPG-RL-KP is the direct mini-batch implementation of our KPG-RL-KP model. m-PKPG-RL-KP is the mini-batch implementation of our partial KPG-RL-KP model. We can see that by partially matching the samples in the mini-batches, m-KPG-RL-KP outperforms m-KPG-RL-KP by a margin of 6.24%. Our partial KPG-RL-KP (m-PKPG-RL-KP) outperforms partial DeepJDOT (m-POT) by 0.68%, indicating that using partial matching, our approach is effective for unsupervised domain adaptation under mini-batch implementation.

B.4 Additional Details for HDA Experiments

Kernel SVM. In the kernel SVM, we use the radial basis function kernel $k(x, y) = \exp(-\gamma\|x - y\|^2)$, where γ is set to the reciprocal of the feature dimension. We use the scikit-learn package of python to implement it by simply running the following codes:

```
clf = SVC(gamma='auto')
clf.fit(feats_train, label_train)
```

Barycentric mapping. The barycentric mapping is defined as follows. Given the transport plan $\pi \in \Sigma_{m \times n}$ and source data point x_{i_0} , the barycentric mapping [49] is defined as $B_\pi(x_{i_0}) = \arg \min_y \sum_{j=1}^n \pi_{i_0, j} c(x_{i_0}, y_j)$. Since c is the squared L_2 -distance in our paper, $B_\pi(x_{i_0})$ has closed-form expression of

$$B_\pi(x_{i_0}) = \frac{1}{\sum_{j=1}^n \pi_{i_0, j}} \sum_{j=1}^n \pi_{i_0, j} y_j. \quad (\text{A-43})$$

B.5 Additional Ablation Studies

Matching accuracy on real data in HDA application. We compare the matching accuracy of different OT models on Office-31 dataset in Table A-4. To compute the matching accuracy, for each transported source data point, we find its nearest neighbor among the target data points to construct a matched pair. If the two points in a pair have the same class labels, the matching is correct, otherwise the matching is incorrect. The matching accuracy is the ratio of correctly matched pairs. In Table A-4, we can see that without the guidance of keypoints, the matching accuracy of GW model is less than 3%. SGW improves the matching accuracy of GW. Our proposed KPG-RL and KPG-RL-GW models achieve better matching accuracy than SGW.

Table A-4: Matching accuracy of different OT models on Office-31 for HDA tasks.

OT models	A→A	A→D	A→W	D→A	D→D	D→W	W→A	W→D	W→W	Avg
GW [27]	2.5	0.7	1.4	2.9	1.8	1.8	2.5	0.4	0.4	1.6
SGW [37]	43.4	61.7	64.3	43.8	73.1	68.9	43.4	72.5	72.7	60.4
KPG-RL	48.7	67.7	66.1	50.7	86.8	81.0	51.0	82.9	83.3	68.7
KPG-RL-GW	50.3	67.4	66.2	49.7	86.5	81.0	51.9	82.9	82.6	68.7

Comparison of different choices for d . Since R_k^s and R_l^t are in the probability simplex, it is reasonable to measure their difference by a distribution divergence/distance. The widely used distribution divergences/distances include the KL-divergence, JS-divergence, and Wasserstein distance. The KL-divergence is not symmetric, so we need to determine the order of inputs. For the Wasserstein distance, one should define the ground metric first. A possible strategy is to set the ground metric to 0 if the two keypoints are paired, otherwise 1. Such a ground metric makes the Wasserstein distance equal to the L_1 -distance. In this work, d is taken as the JS-divergence. We compare the performance of different choices of d in the experiment of HDA on Office-31, as in Table A-5.

Table A-5: Results of different choices of d in HDA experiment on Office-31.

Choices of d	A→A	A→D	A→W	D→A	D→D	D→W	W→A	W→D	W→W	Avg
KL-ST	59.0	89.7	83.6	56.8	95.2	89.0	57.7	93.6	88.1	79.2
KL-TS	58.1	89.0	82.3	54.2	93.9	88.1	54.2	93.2	89.4	78.0
L_1 -distance	57.4	85.8	79.0	58.0	85.8	82.9	58.4	92.6	83.6	75.9
L_2 -distance	52.3	85.8	81.3	53.2	91.3	82.3	52.6	90.3	82.9	74.7
GW	42.0	71.6	70.0	41.6	71.0	69.4	42.3	71.3	70.0	61.0
JS	60.0	91.6	83.6	57.4	95.8	87.7	59.1	95.2	88.4	79.9

In Table A-5, KL-ST and KL-TS denote the KL-divergence $KL(R_k^s, R_l^t)$ and $KL(R_l^t, R_k^s)$ respectively. GW is the Gromov-Wasserstein distance between R_k^s and R_l^t where the source/target cost is taken as the L_2 -distance of source/target keypoints. We find that the JS-divergence achieves the best performance, compared with KL-ST, KL-TS, L_1 -distance, L_2 -distance, and Gromov-Wasserstein.

Results of KPG-RL without using the guiding matrix G . The guiding matrix is the core to impose the relation preservation. We below show in Table A-6 the results for our KPG-RL without using the guiding matrix G , i.e., $L_{kpg}(\pi) = \langle M, \pi \rangle_F$.

Table A-6: Results for different definitions of $L_{kpg}(\pi)$ in HDA experiment.

Definition of $L_{kpg}(\pi)$	$\langle M, \pi \rangle_F$	$\langle M \odot \pi, G \rangle_F$
KPG-RL	60.7	79.9
KPG-RL-GW	60.6	79.6

We can see that without G , both the results of KPG-RL model and KPG-RL-GW model decrease. This may be because $L_{kpg}(\pi) = \langle M, \pi \rangle_F$ may not well impose the guidance of keypoints, since it does not model the “relation” of each point to keypoints.

Sensitivity to source keypoints. In the experiments of the paper, the source keypoints are taken as the source class centers. To study the sensitivity to the location of source keypoints, we randomly sample one data point from each class as a keypoint to construct the source keypoints. We run the experiments with five different samplings for constructing the source keypoints (these five runs are denoted as S1, S2, S3, S4, S5 respectively). The results are reported in the Table A-7. We can see that using the class center as the keypoints achieves the best results, compared with randomly sampling one data point per class as the keypoints. This may be because the class centers are estimated using all the data of each class, and these centers can better represent each class than a randomly sampled data point of each class.

Table A-7: Results for different locations of source keypoints.

S1	S2	S3	S4	S5	Centers
76.8	77.5	78.2	77.8	76.9	79.9

We next study the sensitivity to the number of source keypoints, of which the results are reported in Table A-8. In this experiment, we randomly sample 3/5/7/9 samples (keypoints) or use all the source samples (keypoints) for each class in the source domain to compute the source class centers, which are paired with labeled target samples for constructing the keypoint pairs. The results in Table A-8 show that as the number of source keypoints increases, the accuracy gradually increases. The best result is obtained when all source samples are used to compute the class centers.

Table A-8: Results for different numbers of source keypoints.

Number	3	5	7	9	All
Accuracy	78.4	79.2	79.6	79.8	79.9

On defining keypoints in other practical applications. According to the results in Table A-7, the class centers are better to be the keypoints than the randomly selected samples. For other practical applications, there may not be “class labels” available. We could first cluster the points and then annotate the points near to the center of the clusters as the keypoints.

Time and memory cost. We report the memory and time cost of KPG-RL with different sizes of the guiding matrix G in the bottom row of Tables A-9 and A-10 respectively. For comparisons, we also report the memory and time cost of the Kantorovich Problem (KP). KP needs to calculate the pair-wised cost matrix C , as in Eq. (1). KPG-RL calculates the relation score, and then computes the guiding matrix G . Since we have deduced Sinkhorn’s algorithm for solving KPG-RL, we solve both KP and KPG-RL using Sinkhorn’s algorithm with $\epsilon = 0.005$. Table A-9 shows that KPG-RL costs a slightly larger memory than KP. Table A-10 shows the computational time for solving KPG-RL and KP problems. In the experiment on Office-31, the maximum memory of G is 38M, and the peak memory of the running process is 780M.

Table A-9: Peak memory for computing KP and KPG-RL.

Size ($m \times n$) of C/G	500×500	1000×1000	2000×2000
KP	201M	218M	330M
KPG-RL	207M	232M	378M

Table A-10: Time cost for computing KP and KPG-RL.

Size ($m \times n$) of C/G	500×500	1000×1000	2000×2000
KP	6.9s	27.7s	60.1s
KPG-RL	10.8s	42.1s	76.5s

Sensitivity to hyper-parameters. We show the sensitivity of our method to hyper-parameters τ , τ' in Table A-11, ϵ in Table A-12, and α in Table A-13. ϵ is the the coefficient of entropy regularization. τ and τ' are used to define the relation in Eqs. (7) and (8) in the paper. We set $\tau = \rho \max_{i,j} \{C_{i,j}^s\}$ and $\tau' = \rho \max_{i,j} \{C_{i,j}^t\}$. We then show the results with varying values of ρ . It can be observed that the best value of α is 0.4 in this task, and the results are relatively stable when α ranges in [0.2, 0.5].

Table A-11: Sensitivity of KPG-RL to hyper-parameters τ and τ' in HDA task A \rightarrow W. We set $\tau = \rho \max_{i,j} \{C_{i,j}^s\}$ and $\tau' = \rho \max_{i,j} \{C_{i,j}^t\}$. We then show the results with varying values of ρ .

ρ	0.05	0.07	0.09	0.1	0.2	0.3	0.4	0.5
Accuracy	82.3	83.2	83.2	83.6	83.2	82.9	82.6	82.6

Table A-12: Sensitivity of KPG-RL to hyper-parameter ϵ in HDA task A \rightarrow W.

ϵ	0.0001	0.0005	0.001	0.005	0.01	0.05	0.1	1
Accuracy	83.2	83.2	83.3	83.6	82.3	76.5	74.8	71.0

Table A-13: Sensitivity of KPG-RL-GW to hyper-parameter α in HDA task A \rightarrow W.

α	0.9	0.8	0.7	0.6	0.5	0.4	0.3	0.2	0.1
Accuracy	74.3	78.1	81.5	82.9	84.2	84.5	84.0	84.0	83.7