

GRAPH OPTIMAL TRANSPORT FOR CROSS-DOMAIN ALIGNMENT

Liqun Chen¹ Zhe Gan² Yu Cheng² Linjie Li² Lawrence Carin¹ Jingjing Liu²

Xuehan Chen
06/28/2021



OUTLINE

- Cross-domain alignment
- Graph construction:
 - Image Object detection: Faster R-CNN
 - Word: Bi-GRU
- Graph matching(Optimal Transport Distances)
 - Wasserstein Distance
 - Gromov-Wasserstein Distance
 - Graph Matching via OT Distances
- Experimental Results
 - Image-text retrieval, VQA



GRAPH CONSTRUCTION

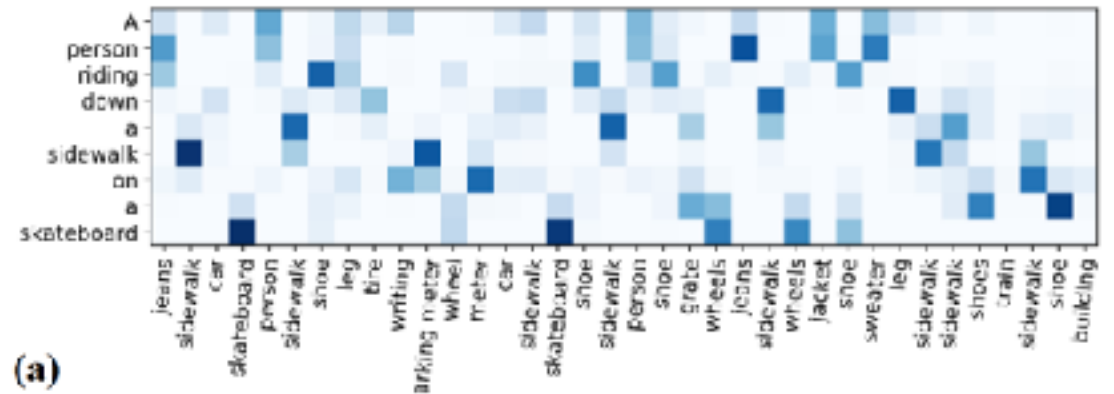


Image features

- Bottom-up features: $X \in R^{n \times s}$
- n : number of bounding boxes, s : feature dim

Text features

- Bi-LSTM/GRU for feature extraction: $Y \in R^{m \times s}$
- m : number of word tokens, s : feature dim



DENOTATION

- Entity feature vector from 2 different domain: $\tilde{\mathbf{X}} = \{\tilde{x}_i\}_{i=1}^n$, $\tilde{\mathbf{Y}} = \{\tilde{y}_j\}_{j=1}^m$
- $\mathbf{X}, \mathbf{Y} = f_{\theta}(\tilde{\mathbf{X}}, \tilde{\mathbf{Y}})$
- $\mathcal{L}(\theta) = \mathcal{L}_{\text{sup}}(\mathbf{X}, \mathbf{Y}, l)$
- $\mathcal{L}(\theta) = \mathcal{L}_{\text{sup}}(\mathbf{X}, \mathbf{Y}, l) + \alpha \cdot \mathcal{L}_{\text{CDA}}(\mathbf{X}, \mathbf{Y})$

- Probability distribution: $\mu \in \mathbf{P}(\mathbb{X})$, $\nu \in \mathbf{P}(\mathbb{Y})$, $\mu = \sum_{i=1}^n u_i \delta_{x_i}$, $\nu = \sum_{j=1}^m v_j \delta_{y_j}$
- $\Pi(\mu, \nu)$: all the joint distributions $\gamma(x, y)$ with marginals $\mu(x)$ and $\nu(y)$.
- Weight vector: $\mathbf{u} = \{u_i\}_{i=1}^n \in \Delta_n$, $\mathbf{v} = \{v_i\}_{i=1}^m \in \Delta_m$
- $\sum_{i=1}^n u_i = \sum_{j=1}^m v_j = 1$

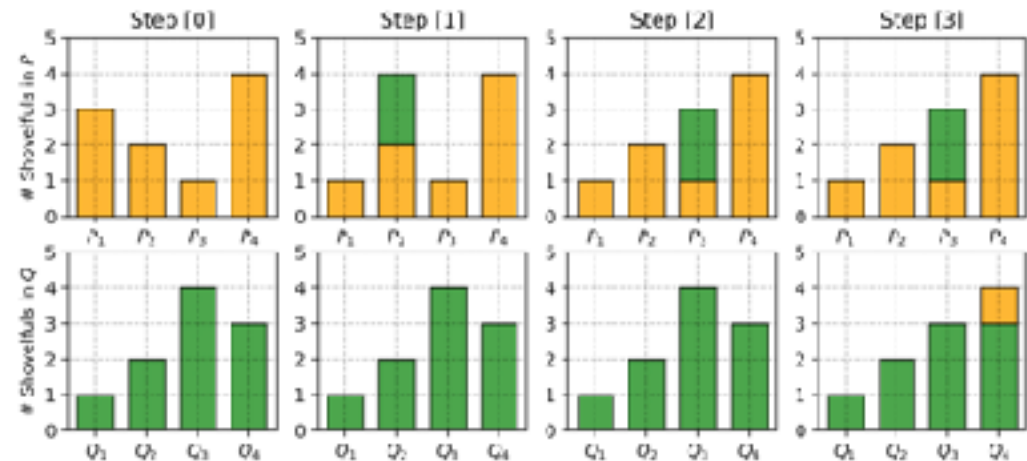


WASSERSTRIN DISTANCE

- Wasserstein Distance is a measure of the distance between two probability distributions. (Earth Mover's distance)
 - eg. the amount of dirt moved \times the moving distance
- In order to change P to look like Q

- First move 2 shovelfuls from P1 to P2 \Rightarrow (P1,Q1) match up.
- Then move 2 shovelfuls from P2 to P3 \Rightarrow (P2,Q2) match up.
- Finally move 1 shovelful from Q3 to Q4 \Rightarrow (P3,Q3) and (P4,Q4) match up.

- $\delta_{i+1} = \delta_i + P_i - Q_i$
- $\delta_0 = 0$
- $\delta_1 = 0 + 3 - 1 = 2$
- $\delta_2 = 2 + 2 - 2 = 2$
- $\delta_3 = 2 + 1 - 4 = -1$
- $\delta_4 = -1 + 4 - 3 = 0$
- Earth Mover's distance is $W = \sum |\delta_i| = 5$



WASSERSTRIN DISTANCE

- Apply optimal transport

- Wasserstrin distance

- Node matching

$$\mathcal{D}_w(\boldsymbol{\mu}, \boldsymbol{\nu}) = \inf_{\gamma \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \mathbb{E}_{(x,y) \sim \gamma} [c(x,y)]$$

$$= \min_{T \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \sum_{i=1}^n \sum_{j=1}^m T_{ij} \cdot c(x_i, y_j)$$

- T is the transport plan

- calculated by Sinkhorn algorithm

- x_i, y_j are two instances/nodes

- C is the pair-wise cost matrix:

$$C_{ij} = d(x_i, y_j)$$

- d : eg. Cosine similarity distance

- where

$$\Pi(\mathbf{u}, \mathbf{v}) = \{T \in \mathbb{R}_+^{n \times m} \mid T\mathbf{1}_m = \mathbf{u}, T^\top \mathbf{1}_n = \mathbf{v}\}$$

- $\mathbf{1}_n$: n-dimensional all-one vector

$$c(x_i, y_j) = 1 - \frac{x_i^\top y_j}{\|x_i\|_2 \|y_j\|_2}$$



GROMOV-WASSERSTEIN DISTANCE

- Gromov-Wasserstein Distance:

- Structure matching

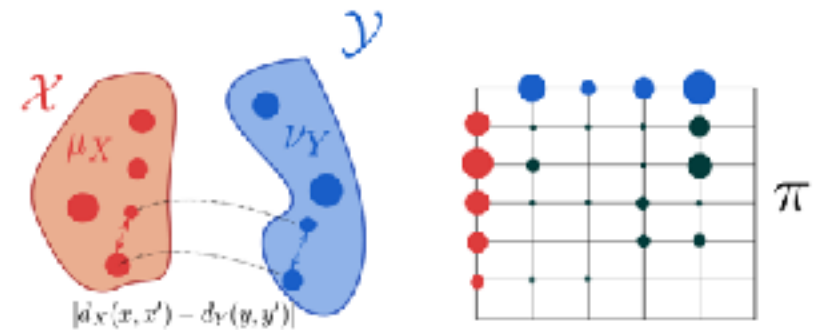
- $\mathcal{D}_{gw}(\mu, \nu) = \inf_{\gamma \in \Pi(\mu, \nu)} \mathbb{E}_{(x, y) \sim \gamma, (x', y') \sim \gamma} \left[L(x, y, x', y') \right]$

- $= \min_{\hat{T} \in \Pi(\mathbf{u}, \mathbf{v})} \sum_{i, i', j, j'} \hat{T}_{ij} \hat{T}_{i'j'} L(x_i, y_j, x_{i'}, y_{j'})$

- $L(x_i, y_j, x_{i'}, y_{j'}) = \left\| c_1(x_i, x_{i'}) - c_2(y_j, y_{j'}) \right\|$

- $L(\cdot)$: intra-graph structural similarity between two pairs of nodes

- C : node similarity within the same graph



GRAPH CONSTRUCTION

- Graph X : $\mathcal{G}_x (\mathbb{V}_x, \mathcal{E}_x)$,
 - Node $i \in \mathbb{V}_x$, feature vector x_i .
 - Edges \mathcal{E}_x : calculate the similarity between a pair of entities inside a graph
- Image graph
 - Dot-product/cosine distance between objects within the image

- $$C_x = \left\{ \cos(x_i, x_j) \right\}_{i,j} \in \mathbb{R}^{n \times n}$$

- Text graph

- $$C_y = \left\{ \cos(y_i, y_j) \right\}_{i,j} \in \mathbb{R}^{m \times m}$$

- Graph Pruning: sparse graph representation

- $$C_x = \max(C_x - \tau, 0)$$
, If $[C_x]_{ij} > 0$, an edge is added between node i and j .



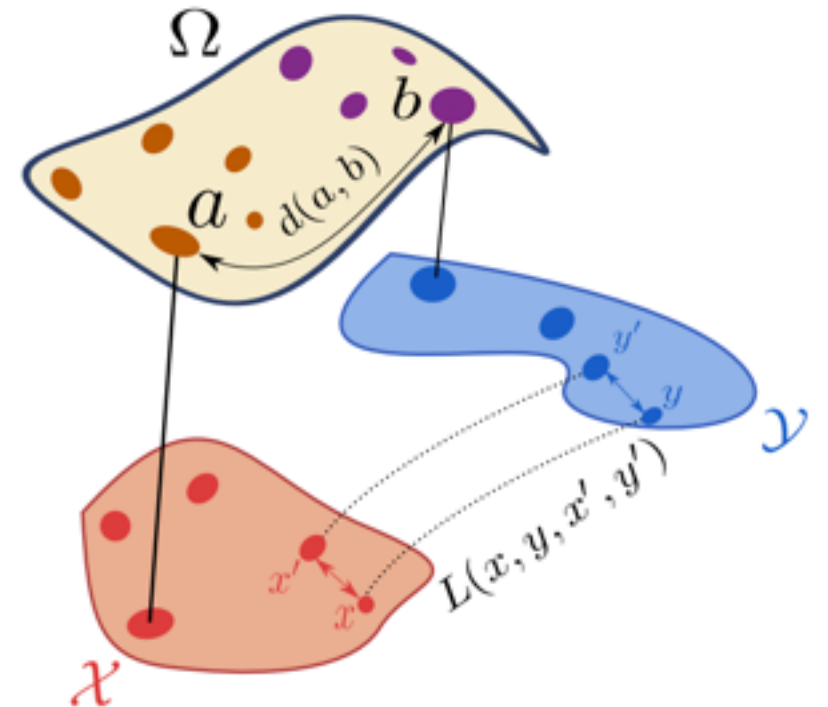
GOT (GRAPH-BASED OPTIMAL TRANSPORT)

- GOT (Graph-based Optimal transport)
- Wasserstrin distance
 - Node matching
- Gromov-Wasserstein Distance:
 - Structure matching
- Combining WD, GWD for better feature alignment

$$\mathcal{D}_{got}(\mu, \nu) = \min_{\mathbf{T} \in \Pi(\mathbf{u}, \mathbf{v})} \sum_{i, i', j, j'} \mathbf{T}_{ij} \left(\lambda c(\mathbf{x}_i, \mathbf{y}_j) + (1 - \lambda) \mathbf{T}_{i'j'} \mathcal{L}(\mathbf{x}_i, \mathbf{y}_j, \mathbf{x}_{i'}, \mathbf{y}_{j'}) \right)$$

$$L_{GOT} = \lambda_1 L_W + \lambda_2 L_{GW}$$

$$L_{\text{unified}} = \lambda c(\mathbf{x}, \mathbf{y}) + (1 - \lambda) L(\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}')$$



- To solve WD with an entropic regularizer

- $\min_{\mathbf{T} \in \Pi(\mathbf{u}, \mathbf{v})} \sum_{i=1}^n \sum_{j=1}^m \mathbf{T}_{ij} c(\mathbf{x}_i, \mathbf{y}_j) - \beta H(\mathbf{T})$

- $H(\mathbf{T}) = - \sum_{i,j} \mathbf{T}_{ij} \log \mathbf{T}_{ij}$

- Gibbs kernel for Sinkhorn:

- $P_\gamma \stackrel{\text{def}}{=} \underset{P \in U(a,b)}{\text{argmin}} \langle P, M_{XY} \rangle - \gamma E(P) \quad \mathbf{u} \in \mathbb{R}_+^n, \mathbf{v} \in \mathbb{R}_+^m$

- $P_\gamma = \text{diag}(\mathbf{u}) K \text{diag}(\mathbf{v}), \quad K \stackrel{\text{def}}{=} e^{-M_{XY}/\gamma}$

- $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a} \quad \text{and} \quad \mathbf{v} \odot (\mathbf{K}^T \mathbf{u}) = \mathbf{b}$

- $\mathbf{u}^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{a}}{\mathbf{K}\mathbf{v}^{(\ell)}} \quad \text{and} \quad \mathbf{v}^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{b}}{\mathbf{K}^T \mathbf{u}^{(\ell+1)}}$

Algorithm 1 Computing Wasserstein Distance.

- 1: **Input:** $\{\mathbf{x}_i\}_{i=1}^n, \{\mathbf{y}_j\}_{j=1}^m, \beta$
 - 2: $\sigma = \frac{1}{n} \mathbf{1}_n, \mathbf{T}^{(1)} = \mathbf{1}\mathbf{1}^T$
 - 3: $\mathbf{C}_{ij} = c(\mathbf{x}_i, \mathbf{y}_j), \mathbf{A}_{ij} = e^{-\frac{\mathbf{C}_{ij}}{\beta}}$
 - 4: **for** $t = 1, 2, 3 \dots$ **do**
 - 5: $\mathbf{Q} = \mathbf{A} \odot \mathbf{T}^{(t)}$ // \odot is Hadamard product
 - 6: **for** $k = 1, 2, 3, \dots, K$ **do**
 - 7: $\delta = \frac{1}{n\mathbf{Q}\sigma}, \sigma = \frac{1}{n\mathbf{Q}^T\delta}$
 - 8: **end for**
 - 9: $\mathbf{T}^{(t+1)} = \text{diag}(\delta)\mathbf{Q}\text{diag}(\sigma)$
 - 10: **end for**
 - 11: $\mathcal{D}_{wd} = \langle \mathbf{C}^T, \mathbf{T} \rangle$
 - 12: Return $\mathbf{T}, \mathcal{D}_w$ // $\langle \cdot, \cdot \rangle$ is the Frobenius dot-product
-



ALGORITHM2

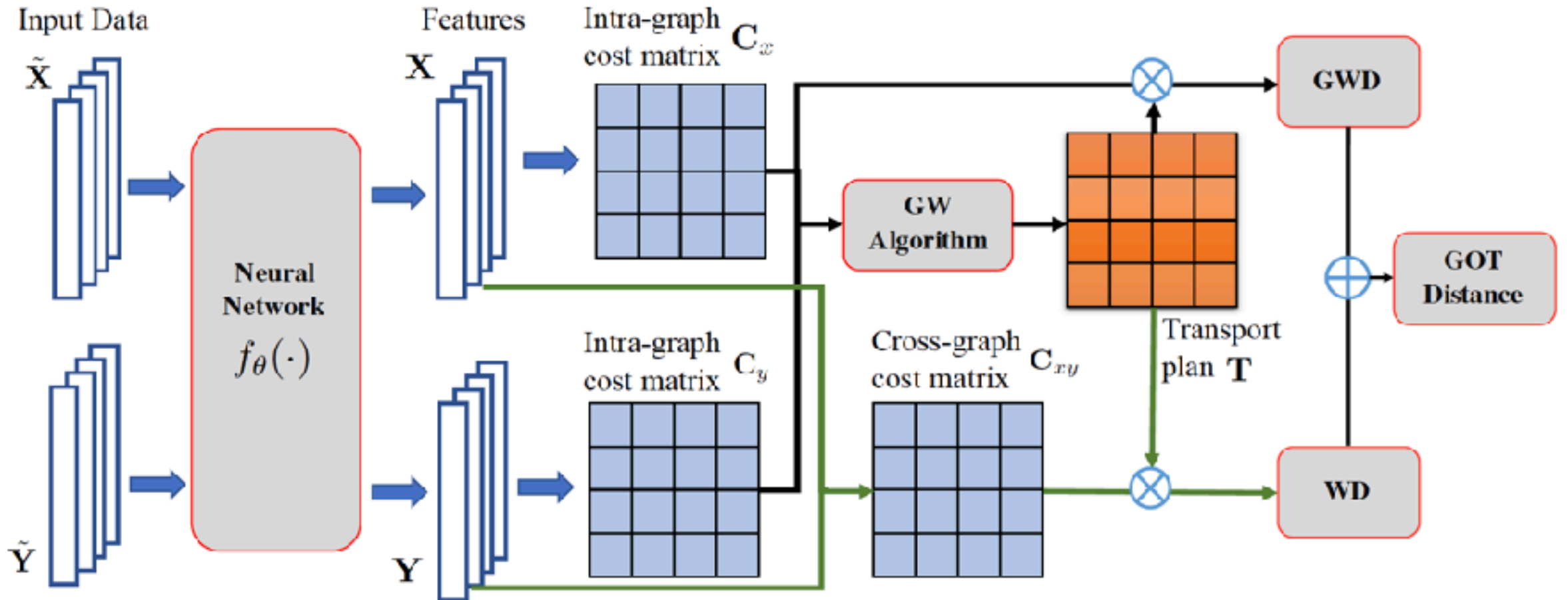
Algorithm 2 Computing Gromov-Wasserstein Distance.

1: **Input:** $\{\mathbf{x}_i\}_{i=1}^n, \{\mathbf{y}_j\}_{j=1}^m$, probability vectors \mathbf{p}, \mathbf{q}
2: Compute intra-domain similarities:
3: $[\mathbf{C}_x]_{ij} = \cos(\mathbf{x}_i, \mathbf{x}_j), [\mathbf{C}_y]_{ij} = \cos(\mathbf{y}_i, \mathbf{y}_j)$,
4: Compute cross-domain similarities:
5: $\mathbf{C}_{xy} = \mathbf{C}_x^2 \mathbf{p} \mathbf{1}_m^\top + \mathbf{C}_y \mathbf{q} (\mathbf{C}_y^2)^\top$
6: **for** $t = 1, 2, 3 \dots$ **do**
7: // Compute the pseudo-cost matrix
8: $\mathcal{L} = \mathbf{C}_{xy} - 2\mathbf{C}_x \mathbf{T} \mathbf{C}_y^\top$
9: Apply Algorithm 1 to solve transport plan \mathbf{T}
10: **end for**
11: $\mathcal{D}_{gw} = \langle \mathcal{L}^\top, \mathbf{T} \rangle$
12: Return $\mathbf{T}, \mathcal{D}_{gw}$

- $\mathcal{D}_{gw}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \inf_{\gamma \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma, (\mathbf{x}', \mathbf{y}') \sim \gamma} [L(\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}')] = \min_{\hat{\mathbf{T}} \in \Pi(\mathbf{u}, \mathbf{v})} \sum_{i, i', j, j'} \hat{\mathbf{T}}_{ij} \hat{\mathbf{T}}_{i'j'} L(\mathbf{x}_i, \mathbf{y}_j, \mathbf{x}_{i'}, \mathbf{y}_{j'})$
- $L(\mathbf{x}_i, \mathbf{y}_j, \mathbf{x}_{i'}, \mathbf{y}_{j'}) = \left\| c_1(\mathbf{x}_i, \mathbf{x}_{i'}) - c_2(\mathbf{y}_j, \mathbf{y}_{j'}) \right\|$



SCHEMATIC COMPUTATION GRAPH OF GOT



ALGORITHM3

Algorithm 3 Computing GOT Distance.

- 1: **Input:** $\{\mathbf{x}_i\}_{i=1}^n, \{\mathbf{y}_j\}_{j=1}^m$, hyper-parameter λ
- 2: Compute intra-domain similarities:
- 3: $[\mathbf{C}_x]_{ij} = \cos(\mathbf{x}_i, \mathbf{x}_j), [\mathbf{C}_y]_{ij} = \cos(\mathbf{y}_i, \mathbf{y}_j),$
- 4: $\mathbf{x}'_i = g_1(\mathbf{x}_i), \mathbf{y}'_j = g_2(\mathbf{y}_j)$ // g_1, g_2 denote two MLPs
- 5: Compute cross-domain similarities:
- 6: $\mathbf{C}_{ij} = \cos(\mathbf{x}'_i, \mathbf{y}'_j)$
- 7: **if** \mathbf{T} is shared: **then**
- 8: Update \mathcal{L} in Algorithm 2 (Line 8) with:
- 9: $\mathcal{L}_{\text{unified}} = \lambda \mathbf{C} + (1 - \lambda) \mathcal{L}$
- 10: Plug in $\mathcal{L}_{\text{unified}}$ back to Algorithm 2 and solve new \mathbf{T}
- 11: Compute \mathcal{D}_{got}
- 12: **else**
- 13: Apply Algorithm 1 to obtain \mathcal{D}_w
- 14: Apply Algorithm 2 to obtain \mathcal{D}_{gw}
- 15: $\mathcal{D}_{\text{got}} = \lambda \mathcal{D}_w + (1 - \lambda) \mathcal{D}_{gw}$
- 16: **end if**
- 17: Return \mathcal{D}_{got}

- $L_{\text{GOT}} = \lambda_1 L_W + \lambda_2 L_{GW}$
- $L_{\text{unified}} = \lambda c(\mathbf{x}, \mathbf{y}) + (1 - \lambda) L(\mathbf{x}, \mathbf{y}, \mathbf{x}', \mathbf{y}')$
- $\mathcal{D}_{\text{got}}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\mathbf{T} \in \Pi(\mathbf{u}, \mathbf{v})} \sum_{i, i', j, j'} \mathbf{T}_{ij} \left(\lambda c(\mathbf{x}_i, \mathbf{y}_j) \right. \\ \left. + (1 - \lambda) \mathbf{T}_{i'j'} \mathcal{L}(\mathbf{x}_i, \mathbf{y}_j, \mathbf{x}'_{i'}, \mathbf{y}'_{j'}) \right)$



SCAN MODEL

$$s_{ij} = \frac{v_i^T e_j}{\|v_i\| \|e_j\|}, i \in [1, k], j \in [1, n]$$

$$a_i^t = \sum_{j=1}^n \alpha_{ij} e_j \quad \alpha_{ij} = \frac{\exp(\lambda_1 \bar{s}_{ij})}{\sum_{j=1}^n \exp(\lambda_1 \bar{s}_{ij})}$$

$$R(v_i, a_i^t) = \frac{v_i^T a_i^t}{\|v_i\| \|a_i^t\|}$$

$$S_{AVG}(I, T) = \frac{\sum_{i=1}^k R(v_i, a_i^t)}{k}$$

$$l(I, T) = \sum_{\hat{T}} [\alpha - S(I, T) + S(I, \hat{T})]_+ + \sum_{\hat{I}} [\alpha - S(I, T) + S(\hat{I}, T)]_+$$



EXPERIMENTS

■ Image-Text Retrieval

- Image representation: pre-trained Faster R-CNN (36,2048)
- Textual features: bi-directional GRU
- Dataset: Flickr30K (Plummer et al.,2015), COCO (Lin et al., 2014)
- Model: SCAN

- Triple loss:
$$L_{\text{triple}}(x_i, y_i, y_j) = \max \left\{ 0, \gamma + \text{sim}(x_i, y_i) - \text{sim}(x_i, y_j) \right\}$$

- (x_i, y_i) is a positive image-text pair
- y_j is the most negative sample for x_i
- γ : hyper-parameter
- sim: function~ evaluates the similarity of image and text
- Final objective: $L_{\text{final}} = L_{\text{triple}} + L_{\text{GOT}}$



EXPERIMENTS

■ Flickr30K

Method	Sentence Retrieval			Image Retrieval			Rsum
	R@1	R@5	R@10	R@1	R@5	R@10	
SCAN (Faster R-CNN, ResNet) (Lee et al., 2018)	67.7	88.9	94.0	44.0	74.2	82.6	452.2
Ours (Faster R-CNN, ResNet):							
SCAN + WD	70.9	92.3	95.2	49.7	78.2	86.0	472.3
SCAN + GWD	69.5	91.2	95.2	48.8	78.1	85.8	468.6
SCAN + GOT	70.9	92.8	95.5	50.7	78.7	86.2	474.8

■ COCO

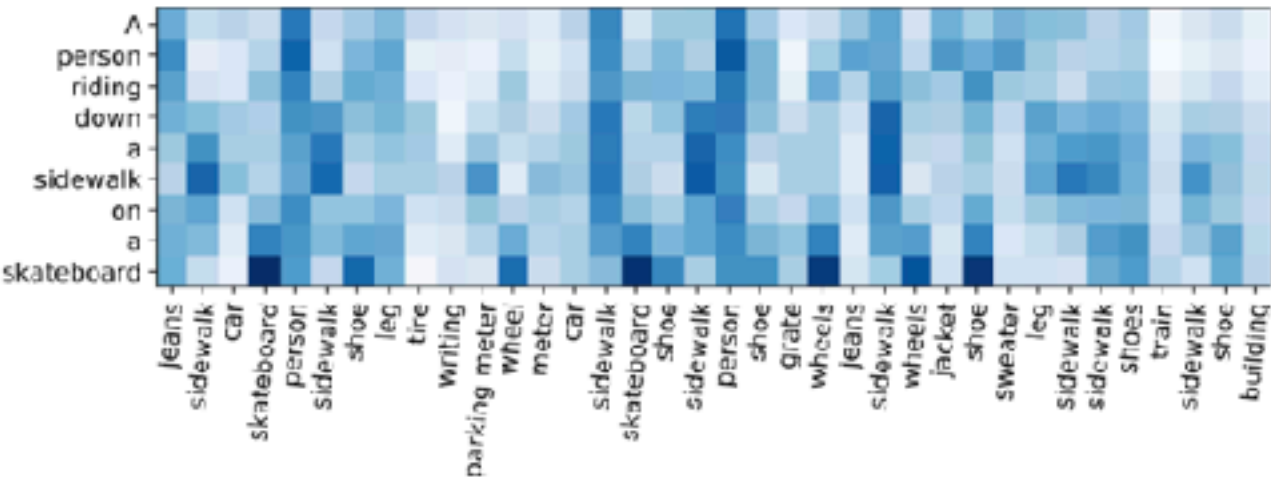
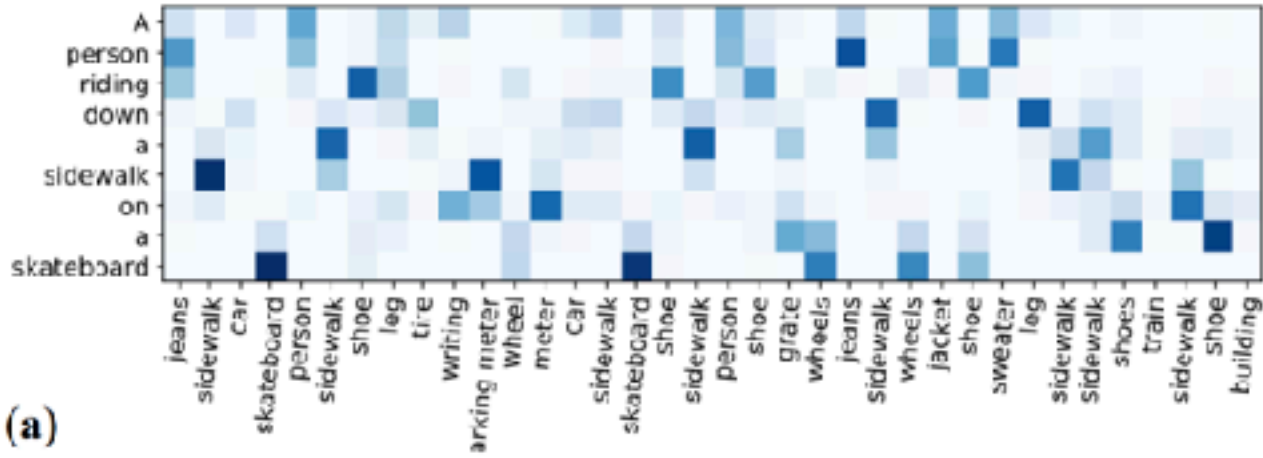
SCAN (Faster R-CNN, ResNet)(Lee et al., 2018)	46.4	77.4	87.2	34.4	63.7	75.7	384.8
Ours (Faster R-CNN, ResNet):							
SCAN + WD	50.2	80.1	89.5	37.9	66.8	78.1	402.6
SCAN + GWD	47.2	78.3	87.5	34.9	64.4	76.3	388.6
SCAN + GOT	50.5	80.2	89.8	38.1	66.8	78.5	403.9



EXPERIMENTS

word tokens

transport plan from GOT



learned attention matrix from SCAN

(b)



image regions



EXPERIMENTS

Model	BAN	BAN+GWD	BAN+WD	BAN+GOT
Score	66.00	66.21	66.26	66.44

Table 2. Results (accuracy) on VQA 2.0 validation set, using BAN (Kim et al., 2018) as baseline.

Model	BUTD	BAN-1	BAN-2	BAN-4	BAN-8
w/o GOT	63.37	65.37	65.61	65.81	66.00
w/ GOT	65.01	65.68	65.88	66.10	66.44

Table 3. Results (accuracy) of applying GOT to BUTD (Anderson et al., 2018) and BAN- m (Kim et al., 2018) on VQA 2.0. m denotes the number of glimpses.

- Visual Question Answering
- Dataset: VQA 2.0
 - human annotated QA pairs on COCO images
- Model: BAN, BUTD
- Classification:
 - cross-entropy loss L_{cls}
- Final objective:
 - $L_{\text{final}} = L_{cls} + L_{GOT}$



REFERENCE

- Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., and Zhang, L. Bottom-up and top-down attention for image captioning and visual question answering. In CVPR, 2018.
- Cuturi, M. and Peyr e, G. Computational optimal transport.2017.
- Peyr e, G., Cuturi, M., and Solomon, J. Gromov-wasserstein averaging of kernel and distance matrices. In ICML, 2016.
- Lee, K.-H. et al. Stacked cross attention for image-text matching. In ECCV, 2018.
- Faghri, F., Fleet, D. J., Kiros, J. R., and Fidler, S. Vse++: Improved visual-semantic embeddings. In BMVC, 2018.

