



# Transferable Attention for Domain Adaptation

Ximei Wang, Liang Li, Weirui Ye, Mingsheng Long,<sup>∞</sup> Jianmin Wang

School of Software, Tsinghua University, China KLiss, MOE; BNRist; Research Center for Big Data, Tsinghua University, China {wxm17,liliang17,ywr16}@mails.tsinghua.edu.cn {mingsheng,jimwang}@tsinghua.edu.cn

### Unsupervised Domain Adaptation(UDA)



• Supervised Learning



Unsupervised Domain Adaptation(UDA)



Unsupervised Domain Adaptation(UDA)

• Unsupervised Domain Adaptation(UDA)



✓ Goal: Achieving good performance on the target domain.



Background



✓ 减小最大均值差异 (Maximum Mean Discrepancy)

$$MMD(X,Y) = ||rac{1}{n}\sum_{i=1}^n \phi(x_i) - rac{1}{m}\sum_{j=1}^m \phi(y_j)||_H^2$$

✓ 对抗找到不变特征 (Domain-invariant features)



#### **Motivation**



However, it is obvious that not all regions of an image are transferable, while forcefully aligning the untransferable regions may lead to negative transfer.



□ Ignore features that are not helpful for category discrimination.

#### Transferable Attention for Domain Adaptation





- Multi-adversarial network is developed for local attention to highlight the representations of those regions with higher transferability
- Global adversarial network (green) is utilized to enhance the prediction certainty of the images more similar in the feature space across domains.

#### **Transferable Local Attention**





$$L_l = \frac{1}{Kn} \sum_{k=1}^K \sum_{\boldsymbol{x}_i \in \mathcal{D}_s \cup \mathcal{D}_t} L_d(G_d^k(\boldsymbol{f}_i^k), d_i), \ \hat{d}_i^k = G_d^k(\boldsymbol{f}_i^k)$$

#### **Transferable Global Attention**





## **Finally Objection**



$$C(\theta_{f}, \theta_{b}, \theta_{y}, \theta_{d}, \theta_{d}^{k}|_{k=1}^{K}) = L_{y} + \gamma L_{h} - \lambda(L_{g} + L_{l})$$

$$= \frac{1}{n_{s}} \sum_{\boldsymbol{x}_{i} \in \mathcal{D}_{s}} L_{y}(G_{y}(G_{b}(\boldsymbol{h}_{i})), y_{i})$$

$$- \frac{\gamma}{n} \sum_{\boldsymbol{x}_{i} \in \mathcal{D}} \sum_{j=1}^{C} m_{i} \cdot \boldsymbol{p}_{i,j} \cdot \log(\boldsymbol{p}_{i,j})$$

$$- \frac{\lambda}{n} [\sum_{\boldsymbol{x}_{i} \in \mathcal{D}} L_{d}(G_{d}(G_{b}(\boldsymbol{h}_{i}), d_{i}))$$

$$+ \frac{1}{K} \sum_{k=1}^{K} \sum_{\boldsymbol{x}_{i} \in \mathcal{D}} L_{d}(G_{d}^{k}((G_{f}(\boldsymbol{x}_{i}))^{k}), d_{i})]$$
(9)

$$\min\max\left(\hat{\theta}_{f},\hat{\theta}_{b},\hat{\theta}_{y}\right) = \underset{\theta_{f},\theta_{b},\theta_{y}}{\arg\min C}\left(\theta_{f},\theta_{b},\theta_{y},\theta_{d},\theta_{d}^{k}|_{k=1}^{K}\right),$$
$$\left(\hat{\theta}_{d},\hat{\theta}_{d}^{1},...,\hat{\theta}_{d}^{K}\right) = \underset{\theta_{d},\theta_{d}^{1},...,\theta_{d}^{K}}{\arg\max C}\left(\theta_{f},\theta_{b},\theta_{y},\theta_{d},\theta_{d}^{k}|_{k=1}^{K}\right).$$
(10)

#### Experiments



Table 1: Accuracy (%) on *Office-31* for unsupervised domain adaption (ResNet)

Method	$A {\rightarrow} W$	$D {\rightarrow} W$	$W {\rightarrow} D$	$A \rightarrow D$	D→A	$W {\rightarrow} A$	Avg
ResNet-50 (He et al. 2016)	$68.4 \pm 0.2$	$96.7\pm0.1$	$99.3 \pm 0.1$	$68.9\pm0.2$	$62.5\pm0.3$	$60.7\pm0.3$	76.1
TCA (Pan et al. 2011)	$72.7\pm0.0$	$96.7\pm0.0$	$99.6\pm0.0$	$74.1\pm0.0$	$61.7\pm0.0$	$60.9\pm0.0$	77.6
GFK (Gong et al. 2012)	$72.8\pm0.0$	$95.0\pm0.0$	$98.2\pm0.0$	$74.5\pm0.0$	$63.4\pm0.0$	$61.0\pm0.0$	77.5
DAN (Long et al. 2015)	$80.5\pm0.4$	$97.1\pm0.2$	$99.6 \pm 0.1$	$78.6\pm0.2$	$63.6\pm0.3$	$62.8\pm0.2$	80.4
RTN (Long et al. 2016)	$84.5\pm0.2$	$96.8\pm0.1$	$99.4\pm0.1$	$77.5\pm0.3$	$66.2\pm0.2$	$64.8\pm0.3$	81.6
DANN (Ganin et al. 2016)	$82.0 \pm 0.4$	$96.9\pm0.2$	$99.1\pm0.1$	$79.7\pm0.4$	$68.2 \pm 0.4$	$67.4 \pm 0.5$	82.2
ADDA (Tzeng et al. 2017)	$86.2\pm0.5$	$96.2\pm0.3$	$98.4\pm0.3$	$77.8\pm0.3$	$69.5\pm0.4$	$68.9\pm0.5$	82.9
JAN (Long et al. 2017)	$85.4\pm0.3$	$97.4\pm0.2$	$\textbf{99.8}\pm0.2$	$84.7\pm0.3$	$68.6\pm0.3$	$70.0\pm0.4$	84.3
MADA (Pei et al. 2018)	$90.0\pm0.1$	$97.4\pm0.1$	$99.6\pm0.1$	$87.8\pm0.2$	$70.3\pm0.3$	$66.4\pm0.3$	85.2
SimNet (Pinheiro 2018)	$88.6\pm0.5$	$98.2\pm0.2$	$99.7\pm0.2$	$85.3\pm0.3$	$73.4\pm0.8$	$71.6\pm0.6$	86.2
GTA (Sankaranarayanan et al. 2018)	$89.5\pm0.5$	$97.9\pm0.3$	$99.8\pm0.4$	$87.7\pm0.5$	$72.8\pm0.3$	$71.4\pm0.4$	86.5
TADA (local)	$89.4\pm0.4$	<b>98.7</b> ± 0.2	$99.8\pm0.2$	$87.2\pm0.2$	$66.4\pm0.2$	$65.3\pm0.3$	84.5
TADA (global)	$92.9\pm0.4$	$98.2\pm0.2$	$99.8\pm0.2$	$88.9\pm0.2$	$69.6\pm0.2$	$71.0\pm0.3$	86.7
TADA (local+global)	$\textbf{94.3}\pm0.3$	<b>98.7</b> $\pm$ 0.1	<b>99.8</b> $\pm$ 0.2	$\textbf{91.6}\pm0.3$	$72.9\pm0.2$	<b>73.0</b> $\pm$ 0.3	88.4

Table 2: Accuracy (%) on *Office-Home* for unsupervised domain adaption (ResNet)

Method	Ar→Cl	Ar→Pr	Ar→Rw	Cl→Ar	Cl→Pr	Cl→Rw	Pr→Ar	Pr→Cl	Pr→Rw	$Rw \rightarrow Ar$	Rw→Cl	$Rw {\rightarrow} Pr$	Avg
ResNet-50 (He et al. 2016)	34.9	50.0	58.0	37.4	41.9	46.2	38.5	31.2	60.4	53.9	41.2	59.9	46.1
DAN (Long et al. 2015)	43.6	57.0	67.9	45.8	56.5	60.4	44.0	43.6	67.7	63.1	51.5	74.3	56.3
DANN (Ganin et al. 2016)	45.6	59.3	70.1	47.0	58.5	60.9	46.1	43.7	68.5	63.2	51.8	76.8	57.6
JAN (Long et al. 2017)	45.9	61.2	68.9	50.4	59.7	61.0	45.8	43.4	70.3	63.9	52.4	76.8	58.3
TADA (local)	47.3	69.1	75.2	56.9	66.4	69.1	55.9	46.9	75.7	68.2	56.2	80.4	63.9
TADA (global)	51.3	66.0	76.5	58.6	69.3	70.3	58.3	52.0	77.1	70.2	57.0	81.5	65.7
TADA (local+global)	<b>53.1</b>	<b>72.3</b>	<b>77.2</b>	<b>59.1</b>	<b>71.2</b>	<b>72.1</b>	<b>59.7</b>	<b>53.1</b>	<b>78.4</b>	<b>72.4</b>	<b>60.0</b>	<b>82.9</b>	<b>67.6</b>

Experiments





Figure 2: The t-SNE visualization of features learned by (a) ResNet, (b) DANN, (c) MADA, and (d) TADA (red: A; blue: W).

Experiments





Figure 3: Attention visualization of the last convolutional layer of ResNet on *Office-Home*. The images on the left are randomly sampled from source domain (Ar) while the right from target domain (Rw). In each group of images, the original input images, the corresponding attentions and the attentions shown in the original input images are illustrated from left to right respectively.



# Thanks