



# On Unsupervised Domain Adaptation: Pseudo Label Guided Mixup for Adversarial Prompt Tuning

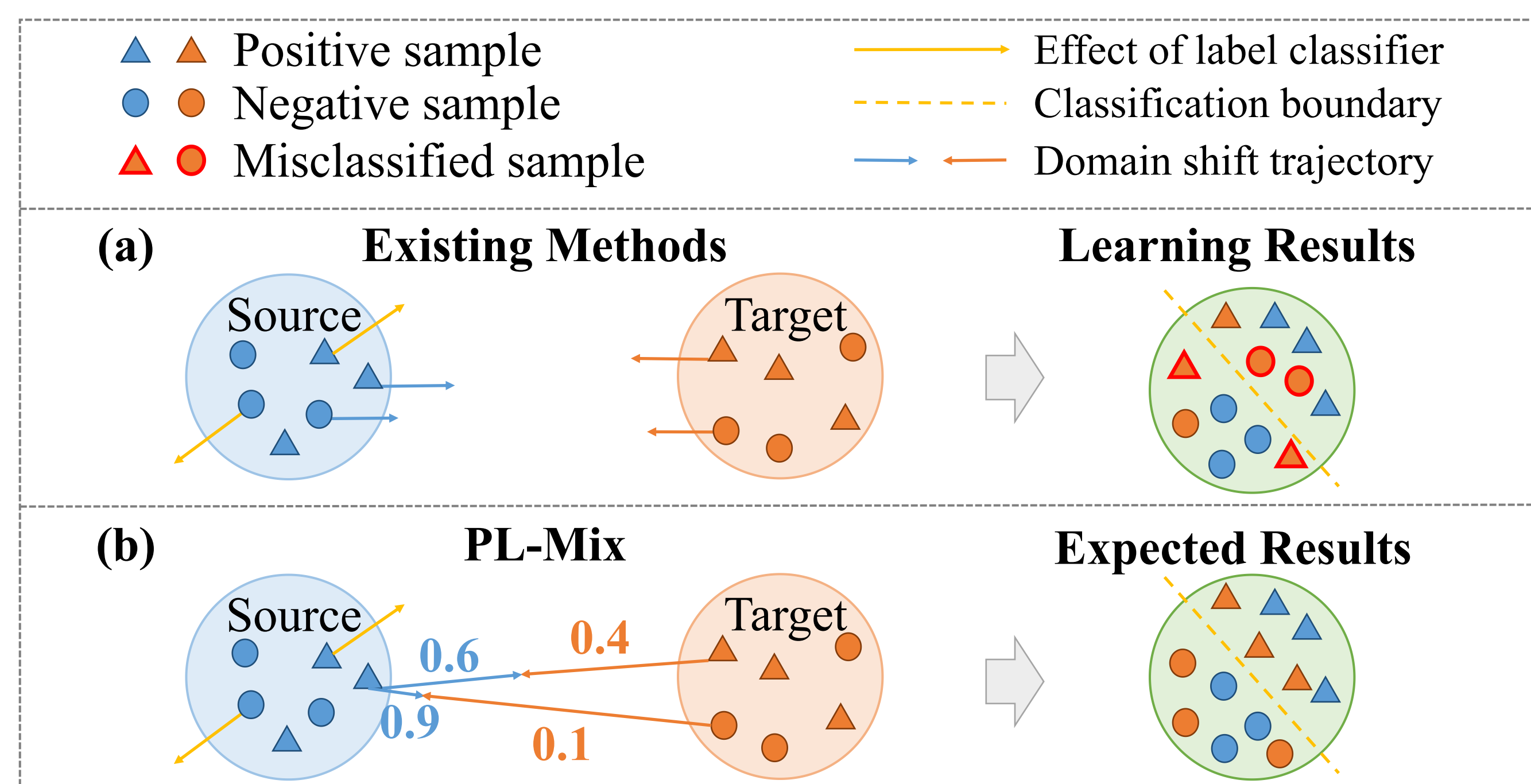
GitHub Link

<https://github.com/fskong/PL-Mix>Fanshuang Kong<sup>1</sup>, Richong Zhang<sup>1,2</sup>, Ziqiao Wang<sup>3</sup>, Yongyi Mao<sup>3</sup><sup>1</sup>SKLSDE, Beihang University, Beijing, China<sup>2</sup>Zhongguancun Laboratory, Beijing, China<sup>3</sup>School of Electrical Engineering and Computer Science, University of Ottawa, Canada

{kongfs, zhangrc}@act.buaa.edu.cn, {zwang286, ymao}@uottawa.ca



## Motivation



*No explicit mechanism that facilitates the positive (or negative) data of the source domain to be attracted towards the corresponding positive (or negative).*

- The fusion of **pseudo labels** and **Mixup** creates intermediate synthetic data between source and target data of the same class, thereby **promoting the alignment between the two domains**.
- The source data should **move more** (Mixup ratio 0.6) toward the target data with **a same and highly confident label**.
- While the **movement should be limited** (Mixup ratio 0.9) for **a distinct and highly confident label**.

## PL-Mix Improves Generalization

**Theorem 1** Let the function space of  $F$  have the finite Natarajan dimension  $d_N$ . Assume that the loss function  $\mathcal{L}_c(\cdot, \cdot; F)$  is  $R$ -subgaussian under  $P_{X^s}^s$ . Then, for any  $F$ , there exists a constant  $C > 0$  such that with probability  $1 - \delta$

$$\mathcal{E}(F) \leq C \sqrt{\frac{d_N \log |\mathcal{Y}| + \log \frac{1}{\delta}}{N_s}} + \sqrt{2R^2 D_{\text{KL}}(P_Z^t \| P_Z^s)}$$

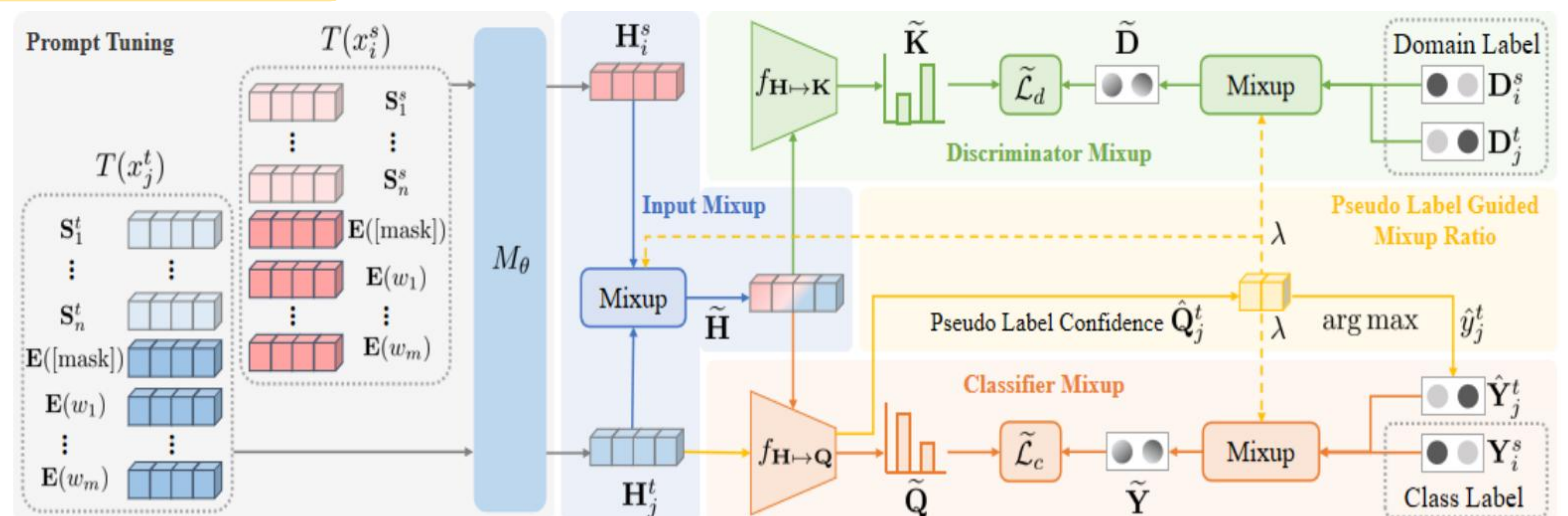
$$+ \sqrt{2R^2 D_{\text{KL}}(P_{Y|Z}^t \| P_{\hat{Y}|Z}^t)} + \sqrt{2R^2 D_{\text{KL}}(P_{\hat{Y}|Z}^t \| P_{Y|Z}^s)},$$

where  $D_{\text{KL}}(\cdot \| \cdot)$  denotes the KL divergence and  $P_{\hat{Y}|Z}^t$  is the conditional pseudo label distribution of the target data.

### □ PL-Mix can improve the generalization guarantee

- Classifier Mixup Reduces the First Term
- Discriminator Mixup Reduces the Second Term
- Classifier Mixup Controls the Last Two Terms

## PL-Mix



### □ Confidence-dependent Mixup Ratio

Pseudo label confidence:  $\hat{Q}_j^{t[\hat{y}_j^t]}$

$$\alpha = \begin{cases} \hat{Q}_j^{t[\hat{y}_j^t]}, & \text{if } \hat{y}_j^t = y_i^s \\ 1 - \hat{Q}_j^{t[\hat{y}_j^t]}, & \text{otherwise} \end{cases}$$

Mixup ratio:  $\lambda \sim \text{Beta}(\alpha, \alpha)$

$$\lambda := \begin{cases} 1, & \text{if } \hat{Q}_j^{t[\hat{y}_j^t]} < \tau \\ \max(1 - \lambda, \lambda), & \text{otherwise} \end{cases}$$

### □ Classifier Mixup

Pseudo class label Mix

$$\tilde{\mathbf{H}} = \lambda \mathbf{H}_i^s + (1 - \lambda) \mathbf{H}_j^t$$

$$\tilde{\mathbf{Y}} = \lambda \mathbf{Y}_i^s + (1 - \lambda) \hat{\mathbf{Y}}_j^t$$

Cross entropy loss of Classifier

$$\tilde{\mathcal{L}}_c = - \sum_{i=1}^{N^s + N^t} \sum_{j=1}^{|\mathcal{Y}|} \tilde{\mathbf{Y}}_i^{[j]} \log \tilde{\mathbf{Q}}_i^{[j]} \min$$

### □ Discriminator Mixup

Domain label Mix

$$\tilde{\mathbf{H}} = \lambda \mathbf{H}_i^s + (1 - \lambda) \mathbf{H}_j^t$$

$$\tilde{\mathbf{D}} = \lambda \mathbf{D}_i^s + (1 - \lambda) \mathbf{D}_j^t$$

Cross entropy loss of Discriminator

$$\max \tilde{\mathcal{L}}_d = - \sum_{i=1}^{N^s + N^t} \sum_{j=1}^{|\mathcal{D}|} \tilde{\mathbf{D}}_i^{[j]} \log \tilde{\mathbf{K}}_i^{[j]}$$

## Experiment

| Bert-base-uncased |              |              |              |              |              |              |              |              |              |              |              |              |              |
|-------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Model             | B → D        | B → E        | B → K        | D → B        | D → E        | D → K        | E → B        | E → D        | E → K        | K → B        | K → D        | K → E        | Avg.         |
| DANN              | 89.70        | 87.30        | 89.55        | 89.55        | 86.05        | 87.69        | 87.15        | 86.05        | 91.91        | 87.65        | 87.72        | 86.05        | 88.56        |
| COBE              | 90.05        | 90.45        | <b>92.90</b> | 90.98        | 90.67        | <b>92.00</b> | 87.90        | 87.87        | 93.33        | 88.38        | 87.43        | 92.58        | 90.39        |
| AdSPT             | -            | -            | -            | -            | -            | -            | -            | -            | -            | -            | -            | -            | -            |
| DANN              | 89.54        | 88.15        | 89.76        | 89.62        | 88.27        | 89.87        | 87.89        | 88.19        | 92.25        | 87.69        | 87.72        | 91.14        | 89.17        |
| COBE*             | 90.13        | 90.92        | 92.28        | 91.05        | 89.75        | 91.67        | 88.25        | <b>88.88</b> | 93.88        | 89.18        | 87.68        | <b>92.87</b> | 90.55        |
| AdSPT             | 90.10        | 90.55        | 92.25        | 90.55        | 89.40        | 90.95        | 88.35        | 87.40        | 93.75        | 88.45        | 87.80        | 92.00        | 90.13        |
| PL-Mix            | <b>90.91</b> | <b>91.04</b> | 91.82        | <b>91.19</b> | <b>91.12</b> | 91.84        | <b>88.86</b> | 88.56        | <b>93.93</b> | <b>89.25</b> | <b>88.27</b> | 92.77        | <b>90.80</b> |

| Roberta-base |              |              |              |              |              |              |       |              |              |              |              |              |              |
|--------------|--------------|--------------|--------------|--------------|--------------|--------------|-------|--------------|--------------|--------------|--------------|--------------|--------------|
| Model        | B → D        | B → E        | B → K        | D → B        | D → E        | D → K        | E → B | E → D        | E → K        | K → B        | K → D        | K → E        | Avg.         |
| DANN         | -            | -            | -            | -            | -            | -            | -     | -            | -            | -            | -            | -            | -            |
| COBE         | -            | -            | -            | -            | -            | -            | -     | -            | -            | -            | -            | -            | -            |
| AdSPT*       | 92.00        | 93.75        | 93.10        | 92.15        | 94.00        | 93.25        | 92.70 | <b>93.15</b> | 94.75        | 92.35        | <b>92.55</b> | 93.95        | 93.14        |
| DANN         | 91.79        | 92.60        | 93.12        | 92.60        | 91.58        | 93.30        | 90.48 | 90.27        | 94.24        | 91.40        | 90.15        | 93.85        | 92.11        |
| COBE         | 92.19        | 92.79        | 95.02        | 93.27        | 93.24        | 94.47        | 92.01 | 90.00        | 95.31        | 91.70        | 90.14        | 94.63        | 92.90        |
| AdSPT        | 92.86        | 93.08        | 94.45        | 93.97        | 93.16        | 94.97        | 91.75 | 89.72        | 95.43        | 91.33        | 90.76        | <b>94.70</b> | 93.02        |
| PL-Mix       | <b>93.60</b> | <b>94.22</b> | <b>95.36</b> | <b>94.19</b> | <b>94.11</b> | <b>95.29</b> | 92.77 | 92.02        | <b>95.67</b> | <b>92.50</b> | 91.71        | 94.65        | <b>93.84</b> |

### □ PL-Mix obtains convincing results

- PL-Mix **outperforms SOTA models on Avg.** both in Bert and Roberta
- PL-Mix aligns source data to target data according to their labels
- PL-Mix outperforms SOTA on **multi-source setting** (details in the paper)

