



HIKVISION

# Text Recognition in Real Scenarios with a Few Labeled Samples

Jinghuang Lin<sup>1</sup> Zhanzhan Cheng<sup>2</sup> Yi Niu<sup>2</sup> Shiliang Pu<sup>2</sup> Shuigeng Zhou<sup>1</sup>

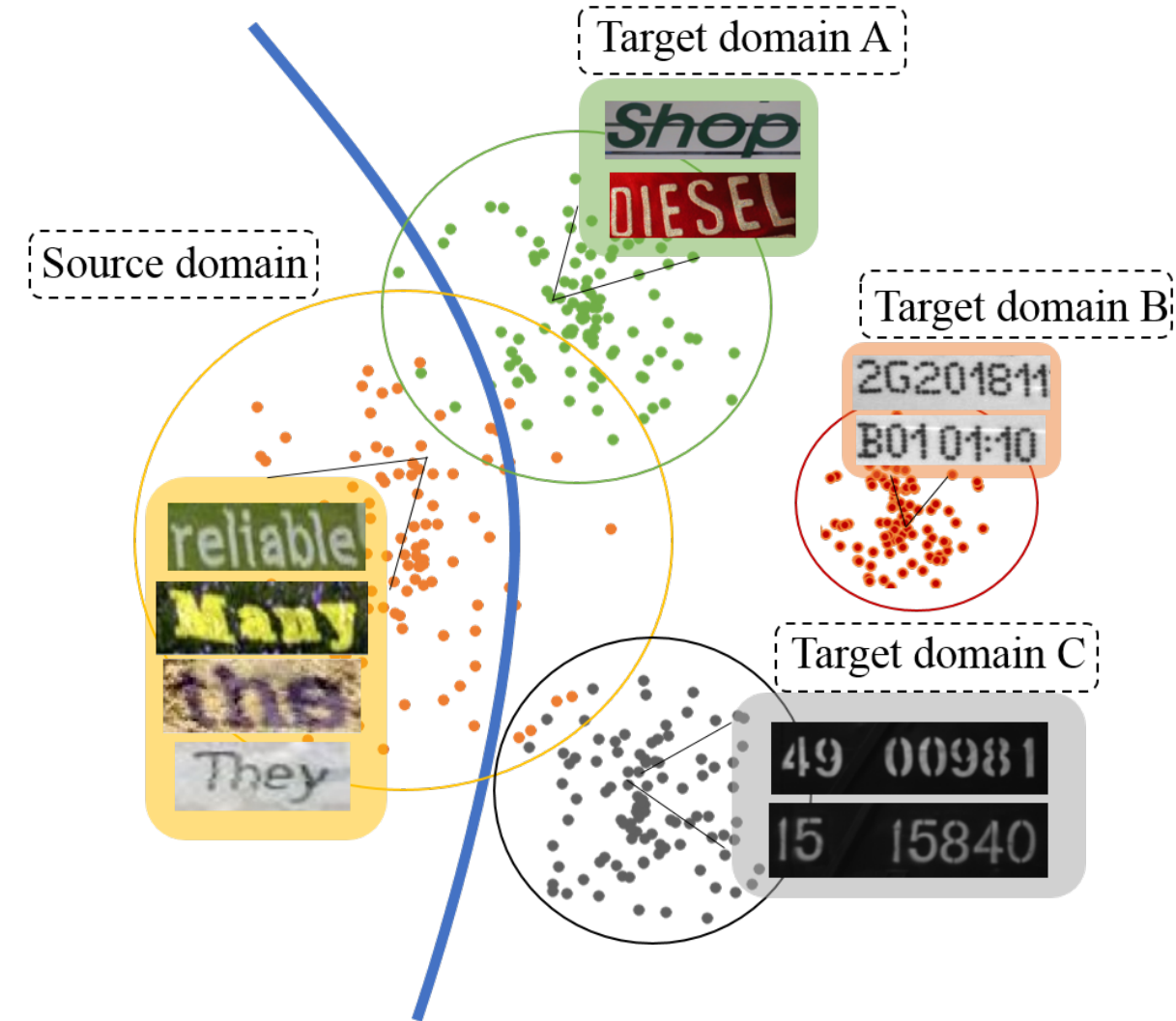
<sup>1</sup> School of Computer Science, Fudan University

<sup>2</sup> Hikvision Research Institute, China



## Motivation

- ❖ Few-shot domain adaptation techniques have shown their effectiveness in handling scenarios where labeled samples are lacked.
- ❖ Most few-shot domain adaptation techniques are focus on character-level task and can hardly handle STR problem because it is a sequence-level image classification task.

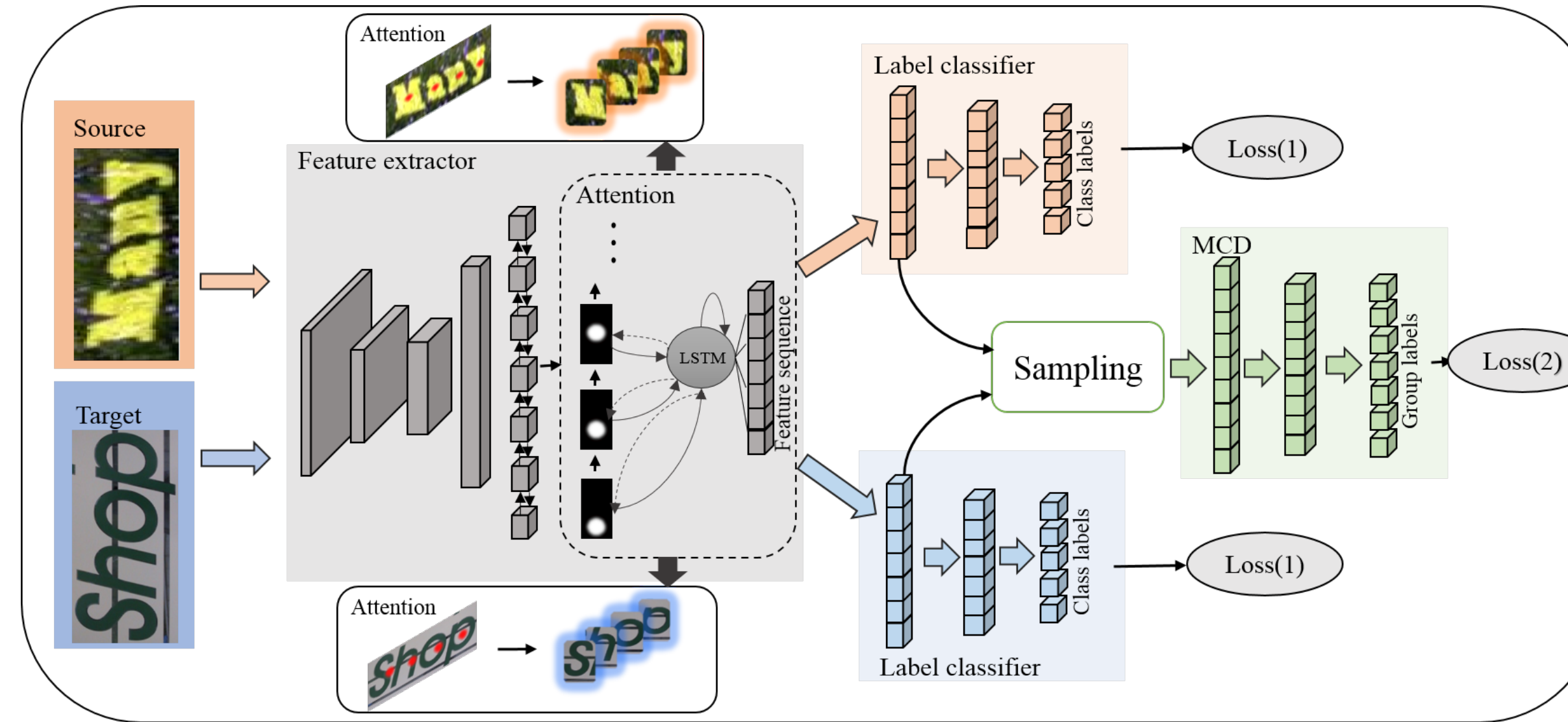


- ❖ We present a few-shot adversarial sequence domain adaptation approach to achieve sequence-level domain confusion by integrating a well-designed attention mechanism with sequence-level adversarial learning strategy into a framework.

## Few-shot Adversarial Sequence Domain Adaptation (FASDA)

- ❖ The architecture of FASDA consists of two procedures :

- 1) Weakly-supervised Character Feature Representation
  - provide “weak” character-level feature representation
- 2) Few-shot Adversarial Learning
  - maximize the character-level confusion between the source domain and the target domain



- ❖ Attention mechanism with Inclusive attending process

Define  $\alpha_{t,j}$  as the attending weight of the  $t$ -th character

$$\alpha_{t,j} = \frac{\exp(e_{t,j})}{\sum_{i=1}^M \exp(e_{t,i})}$$

$$e_{t,j} = w^T \tanh(Ws_{t-1} + Vx_j + b)$$

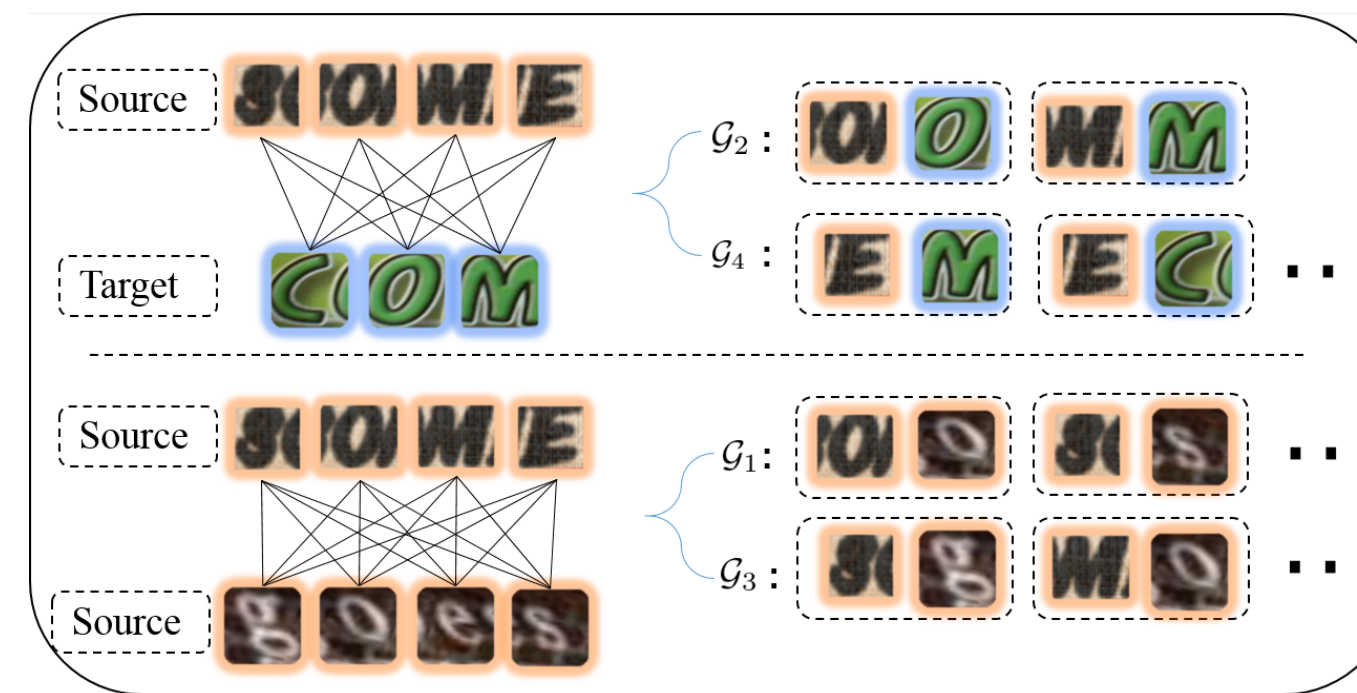
Define  $\alpha'_{t,j}$  as the re-weighted attending weight

$$\alpha'_{t,j} = \lambda \alpha_{t,j} + \frac{1-\lambda}{\eta(1+\eta)} \sum_{i=1}^{\eta} A(t, j-i)(\eta+1-i)$$

$$+ \frac{1-\lambda}{\eta(1+\eta)} \sum_{i=1}^{\eta} A(t, j+i)(\eta+1-i)$$

$$s.t. \quad A(t, j \pm i) = \begin{cases} \alpha_{t, j \pm i} & 1 \leq j \pm i \leq M \\ \alpha_{t, j} & otherwise \end{cases}$$

- ❖ Categories of character representation pair



$G_1 / G_2$  : same class, same/different domain

$G_3 / G_4$  : different class, same/different domain

- ❖ Adversarial learning target:

$$\mathcal{L}_D = - \sum_{i=1}^4 \sum_{S \in \mathcal{G}_i} y_{\mathcal{G}_i} \log(D(\phi(S))) \quad \mathcal{L}_G = - \left[ \sum_{S \in \mathcal{G}_2} y_{\mathcal{G}_1} \log(D(\phi(S))) + \sum_{S \in \mathcal{G}_4} y_{\mathcal{G}_3} \log(D(\phi(S))) \right]$$

## Experiments

- ❖ Comparison with the state-of-the-art

Method	SVT		IC03		IC13	IC15
	50	None	50	Full	None	None
Yao <i>et al.</i> (2014)[45]	75.9	-	88.5	80.3	-	-
Jaderberg <i>et al.</i> (2016)[21]	95.4	80.7	98.7	<b>98.6</b>	93.1	90.8
Shi <i>et al.</i> (2017)[46]	96.4	80.8	98.7	97.6	89.4	86.7
Lee&Osindero (2016)[3]	96.3	80.7	97.9	97.0	88.7	90.0
Cheng <i>et al.</i> (2018)[25]	96	82.8	98.5	97.1	91.5	-
Bai <i>et al.</i> (2018)[1]	96.6	87.5	98.7	97.9	94.6	<b>94.4</b>
Liu <i>et al.</i> (2018)[24]	96.8	87.1	98.1	97.5	94.7	94.0
Shi <i>et al.</i> (2018)[5]	<b>99.2</b>	<b>93.6</b>	<b>98.8</b>	98.0	94.5	91.8
Li <i>et al.</i> (2019)[47]	98.5	91.2	-	-	-	94.0
Luo <i>et al.</i> (2019)[48]	96.6	88.3	98.7	97.8	<b>95.0</b>	92.4
Zhang <i>et al.</i> (2019)[11]	-	84.5	-	-	92.1	91.8
Shi <i>et al.</i> (baseline)(2016)[26]	96.1	81.5	97.8	96.4	88.7	87.5
Cheng <i>et al.</i> (baseline)(2017)[2]	95.7	82.2	98.5	96.7	91.5	89.4
Shi <i>et al.</i> (baseline)(2018)[5]	-	<b>91.6</b>	-	-	93.6	90.5
Luo <i>et al.</i> (baseline)(2019)[48]	-	84.1	-	-	92.5	90.0
Source Only	<b>96.8</b>	85.2	99.0	97.5	92.3	91.6
FT w/ S+T	96.4	86.5	98.7	<b>97.6</b>	93.0	92.4
FASDA	96.5	88.3	<b>99.1</b>	97.5	<b>94.8</b>	<b>94.4</b>

Table 1. Results on SVT, ICDAR 2003, ICDAR 2013 and ICDAR 2015 datasets.

- ❖ Comparison with baselines

Method	SVT	IC03	IC13	IC15
Source Only	19.6	44.1	46.8	14.5
FT w/ T	23.9	46.9	49.7	15.5
FT w/ S+T	25.1	52.3	51.1	16.4
FASDA-CR	27.5	55.8	54.9	18.6
FASDA-CR <sup>+</sup>	28.8	56.8	56.6	19.1
FASDA-IA-CR <sup>+</sup>	<b>29.4</b>	<b>58.1</b>	<b>57.5</b>	<b>19.2</b>

Table 2. Results on SVT, ICDAR 2003, ICDAR 2013 and ICDAR 2015 datasets.

FASDA shows the competitive performance with the state-of-the-art STR methods. And we can also see that FASDA performs clearly better than FT w/ S+T almost on all benchmarks.

## Conclusion

We introduced FASDA to implement sequence-level domain adaptation for STR and it can maximize the character-level confusion between the source domain and the target domain to handle the scenarios that only have a few labeled samples. Our contribution can be summarized as below:

- achieve sequence-level domain confusion in STR
- the framework can be trained end-to-end with much fewer sequence-level annotations.