

# Context-Aware Mixup for Domain Adaptive Semantic Segmentation

Qianyu Zhou, Zhengyang Feng, Qiqi Gu, Jiangmiao Pang,  
Guangliang Cheng, Xuequan Lu, Jianping Shi, Lizhuang Ma

**Abstract**—Unsupervised domain adaptation (UDA) aims to adapt a model of the labeled source domain to an unlabeled target domain. Although the domain shifts may exist in various dimensions such as appearance, textures, etc, the contextual dependency, which is generally shared across different domains, is neglected by recent methods. In this paper, we utilize this important clue as explicit prior knowledge and propose end-to-end Context-Aware Mixup (CAMix) for domain adaptive semantic segmentation. Firstly, we design a contextual mask generation strategy by leveraging accumulated spatial distributions and contextual relationships. The generated contextual mask is critical in this work and will guide the domain mixup. In addition, we define the significance mask to indicate where the pixels are credible. To alleviate the over-alignment (e.g., early performance degradation), the source and target significance masks are mixed based on the contextual mask into the mixed significance mask, and we introduce a significance-reweighted consistency loss on it. Experimental results show that the proposed method outperforms the state-of-the-art methods by a large margin on two widely-used domain adaptation benchmarks, i.e., GTAV  $\rightarrow$  Cityscapes and SYNTHIA  $\rightarrow$  Cityscapes.

**Index Terms**—Domain Adaptation, Semantic Segmentation, Domain Mixup, Autonomous Driving.

## I. INTRODUCTION

SEMANTIC segmentation aims to assign a semantic label to each pixel for a given image. Over the past few years, researchers have made great efforts to explore a variety of CNN methods trained on a large-scale segmentation dataset [1]–[3] to tackle this problem [4]–[8]. However, building such a large annotated dataset is both cost-expensive and time-consuming due to the process of annotating pixel-wise labels [3]. A natural idea to overcome this bottleneck is using synthetic data [9], [10] to supervise the segmentation model instead of real data. However, the existing domain gap between the synthetic images [9], [10] and real images [3] often leads to a significant performance drop when the learned source models are directly applied to the unlabelled target data.

Manuscript received xx xx, 2021. revised xx xx 2021.

Q. Zhou, F. Zheng, Q. Gu are with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: {zhouqianyu, zfyfeng97, miemie}@sjtu.edu.cn).

J. Pang is with the Multimedia Laboratory, The Chinese University of Hong Kong, China (e-mail: pangjiangmiao@gmail.com).

G. Cheng and J. Shi are with SenseTime Research, Beijing, China (e-mail: guangliangcheng2014@gmail.com, shijianping@sensetime.com).

X. Lu is with the School of Information Technology, Deakin University, Victoria 3216, Australia (e-mail: xuequan.lu@deakin.edu.au).

L. Ma is with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai 200240, China, and also with the School of Computer Science and Technology, East China Normal University, Shanghai 200062, China (e-mail: ma-lz@cs.sjtu.edu.cn).

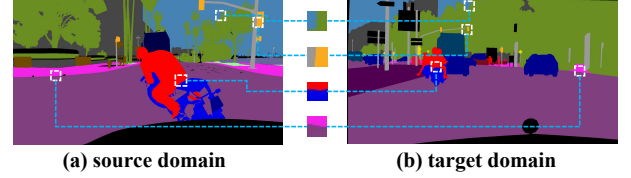


Fig. 1. Previous domain adaptation methods neglect the shared context dependency across different domains and could result in severe negative transfer and training instability. We observe that exploiting contexts as explicit prior knowledge is essential when adapting from the source domain to the target domain.

To address this issue, various unsupervised domain adaptation (UDA) techniques have been proposed to reduce the domain gap in pixel level [11]–[18], feature level [19], [20], [20]–[27] and output level [28]–[34]. Among them, the most common practices are based on adversarial learning [25], [28]–[33], self-training [13], [21], [22], [35]–[37], consistency regularization [38]–[43], entropy minimization [12], [44], [45], etc.

Previous works mainly focused on utilizing common prior knowledge, e.g., appearances, scales, textures, weather, etc., to narrow down the domain gap. Nevertheless, context dependency across different domains has been very sparsely exploited so far in UDA, and how to transfer such cross-domain context still remains under-explored. As shown in Figure 1, we observe that the source and target images usually share similar semantic contexts, e.g., rider is over the bicycle or motorcycle, sidewalk is beside the road, and such context knowledge is crucial particularly when adapting from the source domain to the target domain. The lack of context will lead to severe negative transfer, e.g., early performance degradation during the adaptation process. In addition, most state-of-the-art approaches cannot be trained end-to-end. They heavily depend on the adversarial learning, image-to-image translation or pseudo labeling, and most of them need to fine-tune the models in many offline stages.

In this paper, we attempt to identify context dependency across domains as explicit prior domain knowledge when adapting from the source domain to the target domain. We propose context-aware domain mixup (CAMix) to explore and transfer cross-domain contexts for domain adaptation. Our whole framework is fully end-to-end. The proposed CAMix consists of two key components: contextual mask generation (CMG) and significance-reweighted consistency loss (SRC).

To be specific, the CMG firstly generates a contextual mask

by selectively leveraging the accumulated spatial distribution of the source domain and the contextual relationship of the target domain. This mask is critical in our work and will guide the domain mixup. Guided by it, context-aware domain mixup is performed in three different levels, *i.e.*, input level, output level and significance mask level. Notice that the significance mask is a mask that we define to indicate where the pixels are credible. This contextual mask respectively mixes the input images, the labels and the corresponding significance-masks to narrow down the domain gap.

In addition, we introduce a SRC loss on the significance mask level to alleviate the over-alignment, *e.g.*, early performance degradation, during the adaptation process. In particular, we calculate a significance mask with the help of the target predictive entropy and its dynamic threshold. Then, we mix the target and the source significance masks using the context knowledge as supervisory signals, and utilize the mixed significance mask to reweigh the consistency loss.

To sum up, we propose a *context-aware mixup* architecture for domain adaptive semantic segmentation, which is a fully end-to-end framework. Our contributions are summarized as follows.

- We present a *contextual mask generation* strategy, which leverages the spatial distribution of the source domain and the contextual relationship of the target domain. It acts as prior knowledge for guiding the context-aware domain mixup on three different levels.
- We introduce a *significance-reweighted consistency loss*, which alleviates the adverse impacts of the adaptation procedure, *e.g.*, early performance degradation and training instability, under the guidance of context.
- Extensive experimental results show that we outperform state-of-the-art methods by a large margin on two challenging UDA benchmarks. We achieve 55.2% mIoU in GTAV [9]  $\rightarrow$  Cityscapes [3], and 59.7% mIoU in SYNTHIA [10]  $\rightarrow$  Cityscapes [3], respectively.

## II. RELATED WORK

The current mainstream approaches for cross-domain semantic segmentation [11]–[14], [16], [17], [19], [21]–[30], [32]–[34], [38], [46]–[49] include adversarial learning [25], [28]–[33], consistency regularization [38]–[43] and self-training [13], [21], [22], [35]–[37]. As our work is mostly relevant to the latter two categories, we mainly focus on reviewing them.

**Domain mixup:** Mixup has been well-studied in other communities to improve the robustness of models, *e.g.*, semi-supervised learning [50], [51], and point cloud classification [52], [53]. A few works [54]–[56] studied cross-domain mixup in UDA. Nevertheless, these methods work well on simple and small classification datasets (*e.g.* MNIST [57] and SVHN [58]), but can hardly be applied to more challenging tasks, *e.g.*, domain adaptive semantic segmentation. DACS [40] is designed for segmentation, while little attention has been paid to exploiting contexts as prior knowledge to mitigate the domain gaps.

**Consistency regularization:** The key idea of consistency regularization is that the target prediction of the student model

and that of the teacher model should be invariant under different perturbations. The teacher model is an exponential moving average (EMA) of the student model, and then the teacher model could transfer the learned knowledge to the student. Consistency regularization typically appears in Semi-supervised Learning (SSL) [59] and is recently applied to UDA recently [38]–[43], [60]–[62]. For simplicity, we choose [59] as a base framework to realize end-to-end learning.

**Self-training:** Self training [21], [22], [63], [64] aims to generate pseudo labels for the unlabeled target domain, and then fine-tuned the segmentation model on the pseudo labels iteratively in an offline way. Mei et al. [36] concentrated on the quality of pseudo labels and designed an instance adaptive self-training. Li et al. presented a self-supervised learning [13], which alternately trained the image translation model and the self-supervised segmentation adaptation model. In addition, CBST-BNN [65] and ESL [66] both leveraged predictive entropy rather than the maximum softmax predicted probabilities to refine the pseudo labels during the offline self-training. *Our method differs from these methods in several aspects.* Firstly, in contrast to previous offline self-training that generates pseudo labels and fine-tunes the segmentation model iteratively in many stages, our approach can be trained end-to-end in an online manner. Secondly, instead of using a probability-based mask in common self-training, *e.g.*, [21], [22], we calculate an entropy-guided mask with a novel significance-reweighted loss. Thirdly, different from [65], [66] to refine the pseudo labels, our significance mask is calculated based on the prior knowledge of context information.

**Uncertainty estimation:** The idea of exploiting model prediction uncertainty has been utilized in domain adaptation for classification, *e.g.*, Bayesian classifier [67] and Bayesian discriminator [68]. These methods always require an extra discriminator in adversarial training, and can work well on simple and small classification datasets. *Our method differs from these methods in several aspects.* At first, we tackle the more challenging task of semantic segmentation rather than image classification, where the uncertainty of dense pixel-wise predictions instead of image-wise prediction needs to be decreased. Secondly, we avoid using adversarial adaptation in uncertainty estimation which tends to be unstable and inaccurate. Thirdly, in comparison with the aforementioned approaches, we design significance mask level domain mixup between the target significance mask and the source mask, which enables a more informative entropy-guided mask during the domain mixup.

## III. METHODOLOGY

Following the UDA protocols [21], [28], [44], we have access to the source images  $X_S \in S$  with their corresponding labels  $Y_S$ . For the target domain  $T$ , only unlabeled images  $X_T \in T$  are available. Unlike existing UDA methods that overlook the shared context knowledge across domains, we propose a novel context-aware domain mixup (CAMix) to exploit and transfer such cross-domain contexts.

Figure 2 shows the overview of our proposed architecture. Firstly, we present a contextual mask generation (CMG)

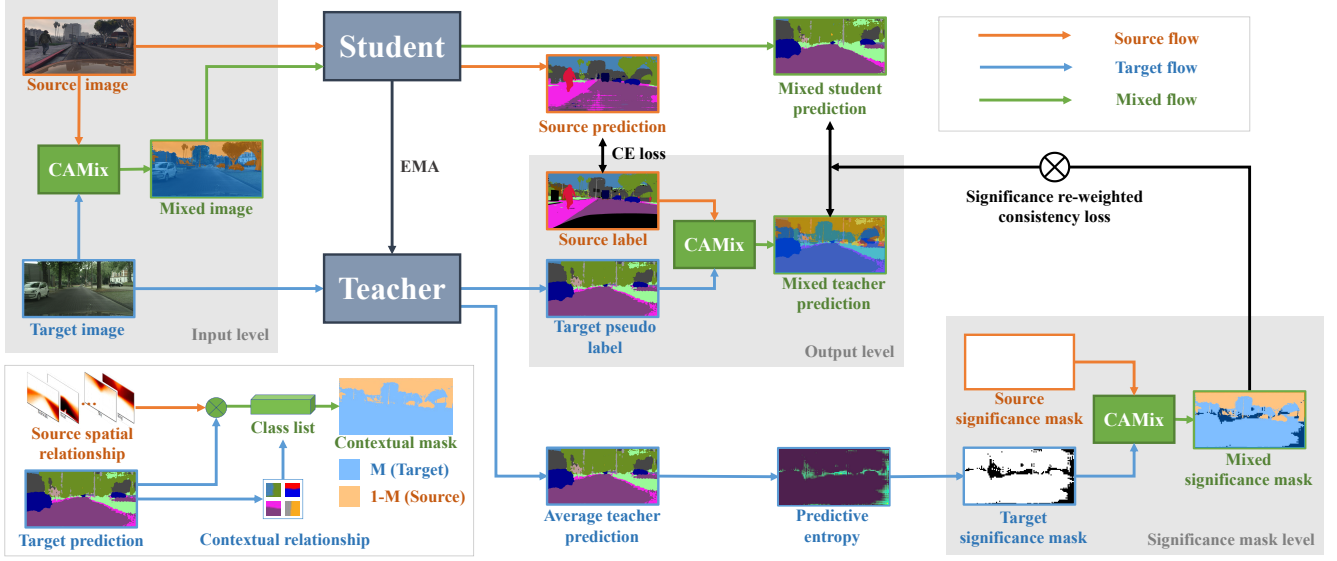


Fig. 2. Overview of the proposed architecture. Firstly, we generate a contextual mask (CMG) by leveraging the spatial distribution of the source domain and the contextual relationship of the target domain. Guided by this mask  $M$ , we perform context-aware mixup (CAMix) in three levels, *i.e.*, input level, output level and significance mask level. Provided the context knowledge, we design a significance re-weighted consistency (SRC) loss to ease the over-alignment between the mixed student and teacher prediction.

strategy for mining the prior spatial distribution of the source domain and the contextual relationships of the target domain, thus generating a mask  $M$ . Guided by this mask, we perform an efficient CAMix on three levels, *i.e.*, input level, output level, and significance mask level (a mask we define to indicate where the pixels are credible.). In particular, the teacher model  $f_{\theta'}$  is an exponential moving average (EMA) of the student model  $f_{\theta}$ . In other words, the proposed CAMix uses the labeled source samples and unlabelled target samples to synthesize the mixed images, the mixed pseudo labels (Section III-B), and the corresponding mixed significance masks (Section III-C). We introduce a significance-reweighted consistency loss (SRC) on the significance mask (SigMask) level to alleviate the over-alignment during the online adaptation procedure.

#### A. Contextual Mask Generation

Intuitively, the source and the target domain share similar context dependency between domains. With this in mind, we identify two kinds of semantic contexts as explicit prior domain knowledge for guiding the domain adaptation procedure. The former is prior spatial contexts of the source domain, and the latter is contextual relationships of the categories in the target domain.

The scenes often have their intrinsic spatial structures, *e.g.*, sky tends to appear on the top of the image while roads are more likely to appear on the bottom. It is intuitive to explore the spatial relationships of the source domain. Thus, we generate a spatial prior tensor  $Q$  with the shape of  $C * H * W$  by counting the class frequencies in the source domain. Each spatial location of  $Q$  is a class distribution, and we treat it as prior knowledge to regularize the target prediction:  $\hat{f}_{\theta'} \leftarrow Q \odot f_{\theta'}(X_T)$ , where  $f_{\theta'}(T)$  is the target prediction of the teacher model.

To exploit the contextual relationship, *e.g.*, the traffic sign should be beside the pole, our core idea is to find the semantic-related categories of the current class presented in the image. In other words, these classes that have contextual relationships to each other can be treated as a meta-class, and then we copy them together from the target images and paste them onto the source images, which prevents certain semantic categories hanging on an inappropriate context.

Specifically, we first get the spatially-modulated pseudo label:  $\tilde{Y}_T \leftarrow \arg \max_{c'} \hat{f}_{\theta'}(i, j, c')$ . Next, we randomly select half of the classes present in the argmax prediction  $\tilde{Y}_T$ , namely  $c$ . After that, we judge whether each category  $k \in c$  presented in  $\tilde{Y}_T$  is in the meta-class list  $m$  or not. The meta-class list involves several groups of heuristic meta-classes, *e.g.*, pole, traffic sign, traffic light are in one group, and bicycle, motorcycle and rider are in another group, etc. This list is chosen from the prior distribution of the source domain and it is shared in all experiments. If  $k \in c$ , we append the semantic-related classes  $\tilde{k}$  of class  $k$  to the current list  $c$ . The current list  $c$  is dynamic in each forwarding, and we observe it has only about 50% ~ 60% of all classes of the target image.

A binary contextual mask  $M$  is generated by setting the pixels from the final class list  $c$  to value 1 in  $M$ , and all others to value 0.

$$M(i, j) = \begin{cases} 1, & \text{if } \tilde{Y}_T(i, j) \in c \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

where  $i \in h, j \in w$ . We iterate each spatial location to generate the mask. This mask  $M$  is then utilized as prior knowledge to mix the images in the input level, the labels in output level (Section III-B), and the significance mask on the significance mask level (Section III-C) between the source domain and the target domain. The whole algorithm of contextual mask generation is described in Algorithm 1.

---

**Algorithm 1: Contextual Mask Generation Algorithm**


---

**Input:** teacher model  $f_{\theta'}$ , target image  $X_T$ , spatial matrix  $Q$ , a meta-class list  $m$ .

**Output:** contextual mask  $M$  for CAMix.

```

1  $\hat{f}_{\theta'} \leftarrow Q \odot f_{\theta'}(X_T)$ ;
2  $\tilde{Y}_T \leftarrow \arg \max_{c'} \hat{f}_{\theta'}(i, j, c')$ ;
3  $C \leftarrow$  Set of the classes present in  $\tilde{Y}_T$ ;
4  $c \leftarrow$  Randomly select  $|C|/2$  classes in  $C$ ;
5 for each  $k \in c$  do
6   if  $k \in c$  and  $k \in m$  then
7      $\tilde{k} \leftarrow$  the semantic-related classes of  $k$ ;
8     if  $\tilde{k} \in C$  then
9        $c.append(\tilde{k})$ ;
10    end
11  end
12 end
13 for each  $i, j$  do
14    $M(i, j) = \begin{cases} 1, & \text{if } \tilde{Y}_T(i, j) \in c \\ 0, & \text{otherwise} \end{cases}$ 
15 end
16 return  $M$ ;
```

---

### B. Input-level and Output-level Domain Mixup

In the *input level*, the image  $X_S$  and  $X_T$  sampled from the source domain and target domain are synthesized into  $X_M$ :

$$X_M = M \odot X_T + (1 - M) \odot X_S, \quad (2)$$

where  $\odot$  denotes element-wise multiplication.

The weights  $\Phi'_t$  of the teacher model at training step  $t$  are updated by the student's weights  $\Phi_t$  with a smoothing coefficient  $\alpha \in [0, 1]$ , which can be formulated as follows:

$$\Phi'_t = \alpha \cdot \Phi'_{t-1} + (1 - \alpha) \cdot \Phi_t, \quad (3)$$

where  $\alpha$  is the EMA decay that controls the updating rate.

Regarding the *output level*, the label of source domain  $Y_S$  and the pseudo label of target domain  $\hat{Y}_T = f_{\theta'}(X_T)$  are mixed as:

$$Y_M = M \odot \hat{Y}_T + (1 - M) \odot Y_S. \quad (4)$$

Different from [40], [62], we mix the images and the corresponding labels in a target-to-source direction rather than the source-to-target direction. In other words, we copy some categories from the target domain and paste them onto the source domain, where we can add our consideration of both spatial relationship and contextual relationship in such a direction.

### C. Significance-mask Level Domain Mixup

In the significance-mask (SigMask) level domain mixup, we decrease the uncertainties of the mixed teacher prediction with the guidance of contextual mask  $M$  as additional supervisory signals. As a result, we are able to alleviate the adverse impact, *e.g.*, training instability and early performance degradation, and transfer more reasonable knowledge from the teacher to the student.

**Stochastic forward passes.** In particular, we repeat each target image for  $N$  copies and inject a random Gaussian noise for each target sample copy. Then, given a set of pixel-wise predicted class scores  $\{P_i^{(h,w,c)}(x_t)\}_{i=1}^N$  of target samples, we can get the mean of the predictive probability  $\hat{P}_c$  of the  $c$ -th class:

$$\hat{P}_c = \frac{1}{N} \sum_{i=1}^N P_i^{(h,w,c)}(X_T). \quad (5)$$

Note that we do not use any dropout layers during stochastic forward passes. The predictive entropy  $\zeta$  is calculated as:

$$\zeta^{(h,w)} = - \sum_{c=1}^C \hat{P}_c \cdot \log(\hat{P}_c), \quad (6)$$

where all volumes of pixel-wise entropy forms a set  $K = \{\zeta\}_{i=1}^N$ .

**Dynamic threshold.** A dynamic threshold  $H$  is then determined by the predictive entropy rather than the softmax probabilities to filter out the unreliable pixel-wise predictions:

$$H = \beta + (1 - \beta) \cdot e^{\gamma(1-t/t_{max})^2} \cdot K_{sup}, \quad (7)$$

where  $t$  denotes the current training step and  $t_{max}$  is the maximum training step.  $K_{sup}$  means the upper-bound of the volumes' self-information, which is denoted as:  $K_{sup} = \sup\{\zeta\}_{i=1}^N$ . we use the same  $\beta, \gamma$  by default in all experiments. **Significance mask.** We denote the SigMask  $U_T = I(\zeta < H)$ , where  $I$  is an indicator function. Note that although the predictive entropy  $\zeta^{(h,w)}$  is similar to ADVENT [44], we do not perform entropy minimization at all, and our SigMask  $U_T$  is calculated from Eq. (5) to Eq. (7) in a completely different way, with the help of target predictive entropy  $\zeta$  and its dynamic threshold  $H$ .

Given the contextual mask  $M$  as additional supervisory signals, we perform **SigMask level** domain mixup. The significance mask of the source domain  $U_S$  and the target domain  $U_T$  are mixed into  $U_M$ :

$$U_M = M \odot U_T + (1 - M) \odot U_S, \quad (8)$$

where  $U_S$  is a tensor full of 1, because the source labels are provided without uncertainties. And these certain areas do not need to reweigh the consistency loss. Only the uncertain areas in the target  $U_T$  which is below the dynamic threshold  $H$ , are set to 0 to reweigh the consistency loss.

**Significance-reweighted consistency loss.** To encourage the teacher model to transfer more credible knowledge to the student model, we define a SRC loss with the guidance of  $U_M$ :

$$\mathcal{L}_{con}(f_{\theta'}, f_{\theta}) = \frac{\sum_j (U_M \cdot CE(f_{\theta}(X_M), Y_M))}{\sum_j U_M}, \quad (9)$$

where  $f_{\theta'}$  and  $f_{\theta}$  are the teacher model and the student model, respectively.  $CE$  is the abbreviation of the cross-entropy loss. The pixel-wise SigMask  $U_M$  is used to reweigh the consistency loss in a weighted averaging manner.



---

**Algorithm 2:** Context Aware Mixup Algorithm
 

---

**Input:** student model  $f_\theta$ , teacher model  $f_{\theta'}$ , source domain  $D_S$ , target domain  $D_T$ , total iterations  $N$ .

**Output:** teacher model  $f_{\theta'}$ .

- 1 Initialize network parameters  $\theta$  randomly;
- 2  $S_S \leftarrow f_{\theta'}(T)$ ;
- 3  $\tilde{Y}_T \leftarrow Q \odot f_{\theta'}(X_T)$ ;
- 4 Initialize network parameters  $\theta$  randomly. ;
- 5 **for**  $i=1$  **to**  $N$  **do**
- 6    $X_S, Y_S \sim D_S$ ;
- 7    $X_T \sim D_T$ ;
- 8    $\tilde{Y}_T \leftarrow f_{\theta'}(X_T)$ ;
- 9    $X_M \leftarrow$  Input-level mixup by Eq.(2);
- 10    $\tilde{Y}_S \leftarrow f_\theta(X_S), \tilde{Y}_M \leftarrow f_\theta(X_M)$ ;
- 11    $Y_M \leftarrow$  Output-level mixup by Eq.(4);
- 12    $U_T \leftarrow$  Target SigMask by Eq.(5)~Eq.(7) ;
- 13    $U_M \leftarrow$  SigMask-level mixup by Eq.(8);
- 14    $\mathcal{L}_{total} \leftarrow$  Total loss by Eq.(11);
- 15   Compute  $\nabla_\theta \mathcal{L}_{total}$  by backpropagation;
- 16   Perform stochastic gradient descent on  $\theta$ ;
- 17 **end**
- 18 **return**  $(\hat{F}_{enc}, \hat{F}_{seg})$

---

#### D. End-to-end Training

**Segmentation loss.** The segmentation loss  $L_{seg}$  is a cross-entropy loss for optimizing the images from the source domain:

$$\mathcal{L}_{seg} = - \sum_{h=1}^H \sum_{w=1}^W \sum_{c=1}^C Y_S^{(h,w,c)} \log(P_S^{(h,w,c)}), \quad (10)$$

where  $Y_S$  is the ground truth for source images and  $P_S = f_\theta(X_S)^{(h,w,c)}$  is the segmentation output of source images.

**Total loss.** During training, all models on three different levels are jointly trained in an end-to-end manner. The whole framework is optimized by integrating all the aforementioned loss functions:

$$\mathcal{L}_{total} = \mathcal{L}_{seg} + \lambda_{con} \mathcal{L}_{con}, \quad (11)$$

where  $\lambda_{con}$  is the weight of consistency loss. Note that our model is insensitive to hyper-parameters in Eq. 7 and Eq. 11, and we use the same weights and same hyper-parameters across all datasets and experiments. For example, we set the same  $\beta = 0.75, \gamma = -5$  for Eq. 7 and we use the same adaptive schedule for the weight  $\lambda_{con}$  as [40] in all experiments.

The algorithm of CAMix for the whole training process is illustrated in Algorithm 2.

## IV. EXPERIMENTS

Following common UDA protocols [28], [74], we treat the labeled synthetic dataset, *i.e.*, GTAV [9] and SYNTHIA [10], as the source domain, and the unlabeled real dataset *i.e.*, Cityscapes [3] as the target domain.

#### A. Datasets

**Cityscapes** [3] is a dataset focused on autonomous driving, which consists of 2,975 images in the training set and 500 images in the validation set. The images have a fixed spatial resolution of  $2048 \times 1024$  pixels. Following common practice, we trained the model on the unlabelled training set and report our results on the validation set.

**GTAV** [9] is a synthetic dataset including 24,966 photo-realistic images rendered by the gaming engine Grand Theft Auto V (GTAV). The semantic categories are compatible between the two datasets. We used all the 19 official training classes in our experiments.

**SYNTHIA** [10] is another synthetic dataset composed of 9,400 annotated synthetic images with the resolution of  $1280 \times 960$ . Like GTAV, it has semantically compatible annotations with Cityscapes. Following the prior works [20], [75], [76], we use the SYNTHIA-RAND-CITYSCAPES subset [10] as our training set.

#### B. Implementation Details

In our implementation, we employ DeepLab-v2 [4] with ResNet 101 backbone [77]. The backbone is pre-trained on ImageNet [78]. For the DeepLab-v2 network, we use Adam as the optimizer. The initial learning rate is  $2.5 \times 10^{-4}$  which is then decreased using polynomial decay with an exponent of 0.9. The weight decay is  $5 \times 10^{-5}$  and the momentum is 0.9. Following the common UDA protocol [13], [29], when the source domain is GTAV, we resize all images to  $1280 \times 720$ ; when the source domain is SYNTHIA, we resize all images to  $1280 \times 760$ . Then, both the source and target images are randomly cropped to  $512 \times 512$ . We use the same data augmentation as DACS [40], *i.e.*, color jittering and Gaussian blurring. In our SigMask-level CAMix, we perform  $N = 8$  times of stochastic forward passes. Following the previous consistency regularization works, we use the same adaptive schedule as CutMix [50] and DACS [40] for the consistency weight  $\lambda_{con}$ . Our method is implemented in Pytorch [79] on a single NVIDIA Tesla V100, and our proposed method requires about 22GB memory for training the model. Each mini-batch contains two images, one from the source domain and the other sampled from the target domain. We totally train the model for 250k iterations. Following the common UDA protocol [12], [28]–[30], we use the early-stopping.

#### C. Comparison with the State-of-the-Art Methods

Table I and Table III present the comparison results with the state-of-the-arts on two challenging tasks: “GTAV  $\rightarrow$  Cityscapes” and “SYNTHIA  $\rightarrow$  Cityscapes”. Our proposed method significantly outperforms the state-of-the-art techniques by 5%  $\sim$  10% on GTAV  $\rightarrow$  Cityscapes and 6%  $\sim$  12% on SYNTHIA  $\rightarrow$  Cityscapes. Also, it is superior to the non-adaptive baselines by around 22% and 30% on two benchmarks, respectively.

Most of the state-of-the-art approaches perform the adversarial learning, *e.g.*, APODA [34], IntraDA [33], WLabel [69], MRNet [80], FADA [24] and DADA [26], and they need to

TABLE I  
COMPARISON RESULTS (mIoU) FROM GTAV TO CITYSCAPES.

Method	Venue	road	sidewalk	building	wall	fence	pole	light	sign	vegetation	terrain	sky	person	rider	car	truck	bus	train	motorcycle	bike	mIoU
Source Only	-	63.3	15.7	59.4	8.6	15.2	18.3	26.9	15.0	80.5	15.3	73.0	51.0	17.7	59.7	28.2	33.1	3.5	23.2	16.7	32.9
BDL [13]	CVPR'19	91.0	44.7	84.2	34.6	27.6	30.2	36.0	36.0	85.0	43.6	83.0	58.6	31.6	83.3	35.3	49.7	3.3	28.8	35.6	48.5
APODA [34]	AAAI'20	85.6	32.8	79.0	29.5	25.5	26.8	34.6	19.9	83.7	40.6	77.9	59.2	28.3	84.6	34.6	49.2	8.0	32.6	39.6	45.9
STAR [23]	CVPR'20	88.4	27.9	80.8	27.3	25.6	26.9	31.6	20.8	83.5	34.1	76.6	60.5	27.2	84.2	32.9	38.2	1.0	30.2	31.2	43.6
IntraDA [33]	CVPR'20	90.6	37.1	82.6	30.1	19.1	29.5	32.4	20.6	85.7	40.5	79.7	58.7	31.1	86.3	31.5	48.3	0.0	30.2	35.8	46.3
SIM [30]	CVPR'20	90.6	44.7	84.8	34.3	28.7	31.6	35.0	37.6	84.7	43.3	85.3	57.0	31.5	83.8	42.6	48.5	1.9	30.4	39.0	49.2
LTIR [14]	CVPR'20	92.9	55.0	85.3	34.2	31.1	34.9	40.7	34.0	85.2	40.1	87.1	61.0	31.1	82.5	32.3	42.9	0.3	36.4	46.1	50.2
FDA [12]	CVPR'20	92.5	53.3	82.4	26.5	27.6	36.4	40.6	38.9	82.3	39.8	78.0	62.6	34.4	84.9	34.1	53.1	16.9	27.7	46.4	50.5
PCEDA [17]	CVPR'20	91.0	49.2	85.6	37.2	29.7	33.7	38.1	39.2	85.4	35.4	85.1	61.1	32.8	84.1	45.6	46.9	0.0	34.2	44.5	50.5
LSE [37]	ECCV'20	90.2	40.0	83.5	31.9	26.4	32.6	38.7	37.5	81.0	34.2	84.6	61.6	33.4	82.5	32.8	45.9	6.7	29.1	30.6	47.5
WLabel [69]	ECCV'20	91.6	47.4	84.0	30.4	28.3	31.4	37.4	35.4	83.9	38.3	83.9	61.2	28.2	83.7	28.8	41.3	8.8	24.7	46.4	48.2
CrCDA [70]	ECCV'20	92.4	55.3	82.3	31.2	29.1	32.5	33.2	35.6	83.5	34.8	84.2	58.9	32.2	84.7	40.6	46.1	2.1	31.1	32.7	48.6
FADA [24]	ECCV'20	92.5	47.5	85.1	37.6	32.8	33.4	33.8	18.4	85.3	37.7	83.5	63.2	39.7	87.5	32.9	47.8	1.6	34.9	39.5	49.2
LDR [18]	ECCV'20	90.8	41.4	84.7	35.1	27.5	31.2	38.0	32.8	85.6	42.1	84.9	59.6	34.4	85.0	42.8	52.7	3.4	30.9	38.1	49.5
CCM [71]	ECCV'20	93.5	57.6	84.6	39.3	24.1	25.2	35.0	17.3	85.0	40.6	86.5	58.7	28.7	85.8	49.0	56.4	5.4	31.9	43.2	49.9
CD-SAM [72]	WACV'21	91.3	46.0	84.5	34.4	29.7	32.6	35.8	36.4	84.5	43.2	83.0	60.0	32.2	83.2	35.0	46.7	0.0	33.7	42.2	49.2
ASA [46]	TIP'21	89.2	27.8	81.3	25.3	22.7	28.7	36.5	19.6	83.8	31.4	77.1	59.2	29.8	84.3	33.2	45.6	16.9	34.5	30.8	45.1
CLAN [32]	TPAMI'21	88.7	35.5	80.3	27.5	25.0	29.3	36.4	28.1	84.5	37.0	76.6	58.4	29.7	81.2	38.8	40.9	5.6	32.9	28.8	45.5
DAST [73]	AAAI'21	92.2	49.0	84.3	36.5	28.9	33.9	38.8	28.4	84.9	41.6	83.2	60.0	28.7	87.2	45.0	45.3	7.4	33.8	32.8	49.6
Ours	-	<b>93.3</b>	<b>58.2</b>	<b>86.5</b>	36.8	<b>31.5</b>	<b>36.4</b>	35.0	<b>43.5</b>	<b>87.2</b>	<b>44.6</b>	<b>88.1</b>	<b>65.0</b>	24.7	<b>89.7</b>	<b>46.9</b>	<b>56.8</b>	<b>27.5</b>	<b>41.1</b>	<b>56.0</b>	<b>55.2</b>

carefully tune the optimization procedure for min-max problems through a domain discriminator. However, such domain discriminators tend to be unstable and inaccurate. Instead, our method does not require to maintain an extra discriminator during the domain adaptation process, and we outperform these approaches by 6% ~ 10% in mIoU.

In contrast to the offline self-training methods that need to fine-tune the models in many rounds, *e.g.*, CRST [22], LSE [37], CCM [71] and TPLD [35], our whole framework can be trained in a fully end-to-end manner. Benefiting from the online consistency regularization by our proposed components CMG and SRC, our approach significantly outperforms the self-training methods by around 5% ~ 9%.

Compared to the methods which require an image-to-image (I2I) translation or style transfer algorithm to filter out the domain-specific texture or style information, *e.g.*, BDL [13], LDR [18], LTIR [14], FDA [12] and PCEDA [17], our context-aware domain mixup does not require any style/spectral transfer algorithms or deep neural networks for I2I translation. Our domain mixup algorithm is simple and works very well, and it surpasses the translation-based methods by around 5% ~ 8%.

CrCDA [70] learned and enforced the prototypical local contextual-relations in the feature space, and similarly CD-SAM [72] exploits contexts implicitly in the feature space, while the visual cues of context knowledge tend to be lost. Moreover, both of the learning [70], [72] does not *explicitly* exploit the cross-domain contexts in the image space and cannot be trained end-to-end. In contrast, our CAMix explicitly explores the contexts in the image space rather than the feature space, and our architecture can be trained end-to-end. Our approach outperforms the CrCDA [70] by 6.6% and 9.7% in two benchmarks, respectively.

Taking a closer look at per-category performance in Table I and Table III, our approach achieves the highest IoU on most categories, *e.g.*, motorcycle, bicycle, traffic sign, etc. This phenomenon reveals the effectiveness of CAMix among different classes during the adaptation process.

TABLE II  
COMPARISONS WITH EXISTING DOMAIN MIXUP METHODS.

method	mIoU (%)	Gain (%)
Mean Teacher	43.1	-
+ CowMix [51]	48.3	+5.2
+ CutMix [50]	48.7	+5.6
+ DACS [40]	52.1	+9.0
+ iDACS [40]	51.5	+8.4
+ CAMix	<b>55.2</b>	<b>+12.1</b>

#### D. Comparison with the Other Domain Mixup

As shown in Table II, we present the adaptation results of our method and the existing domain mixup algorithms on GTAV → Cityscapes. We choose the Mean Teacher architecture [59] as our baseline in this experiment. The existing domain mixup algorithms are implemented under the same settings. CowMix [51], CutMix [50] are proposed for semi-supervised learning (SSL), and we adapt them to the UDA task, which mixes the source domain image and the target domain image. Besides, we implement the existing cross-domain mixup method, *e.g.*, DACS [40] and inverse DACS. The former DACS means using ClassMix to copy the source categories and paste them onto the target. Inverse DACS (iDACS) [40] uses a target-to-source direction. We re-implement the iDACS and the results are based on experiments.

We analyze that using CowMix [51] results in the occurrence of partial objects in the mixed images, which are harder to learn in the training process. Besides, CutMix [50], DACS [40] and iDACS tend to result in severe label contamination and category confusion when generating the mixed results, thus leading to negative transfer. The main reason is that they neglect the context dependency as prior knowledge for facilitating the domain adaptation. The results shown in Table II demonstrate the superiority of our proposed CAMix to other domain mixup algorithms.

TABLE III  
COMPARISON RESULTS (mIoU) FROM SYNTHIA TO CITYSCAPES.

Method	Venue	road	sidewalk	building	light	sign	vegetation	sky	person	rider	car	bus	motorcycle	bike	mIoU <sub>13</sub>
Source Only	-	36.3	14.6	68.8	5.6	9.1	69.0	79.4	52.5	11.3	49.8	9.5	11.0	20.7	29.5
BDL [13]	CVPR'19	86.0	46.7	80.3	14.1	11.6	79.2	81.3	54.1	27.9	73.7	42.2	25.7	45.3	51.4
DADA [26]	ICCV'19	89.2	44.8	81.4	8.6	11.1	81.8	84.0	54.7	19.3	79.7	40.7	14.0	38.8	49.8
APODA [34]	AAAI'20	86.4	41.3	79.3	22.6	17.3	80.3	81.6	56.9	21.0	84.1	49.1	24.6	45.7	53.1
STAR [23]	CVPR'20	82.6	36.2	81.1	12.2	8.7	78.4	82.2	59.0	22.5	76.3	33.6	11.9	40.8	48.1
IntraDA [33]	CVPR'20	84.3	37.7	79.5	9.2	8.4	80.0	84.1	57.2	23.0	78.0	38.1	20.3	36.5	48.9
LTIR [14]	CVPR'20	92.6	53.2	79.2	1.6	7.5	78.6	84.4	52.6	20.0	82.1	34.8	14.6	39.4	49.3
SIM [30]	CVPR'20	83.0	44.0	80.3	17.1	15.8	80.5	81.8	59.9	33.1	70.2	37.3	28.5	45.8	52.1
FDA [12]	CVPR'20	79.3	35.0	73.2	19.9	24.0	61.7	82.6	61.4	31.1	83.9	40.8	38.4	51.1	52.5
LSE [37]	ECCV'20	82.9	43.1	78.1	9.1	14.4	77.0	83.5	58.1	25.9	71.9	38.0	29.4	31.2	49.4
CrCDA [70]	ECCV'20	86.2	44.9	79.5	9.4	11.8	78.6	86.5	57.2	26.1	76.8	39.9	21.5	32.1	50.0
WLabel [69]	ECCV'20	92.0	53.5	80.9	3.8	6.0	81.6	84.4	60.8	24.4	80.5	39.0	26.0	41.7	51.9
CCM [71]	ECCV'20	79.6	36.4	80.6	22.4	14.9	81.8	77.4	56.8	25.9	80.7	45.3	29.9	52.0	52.9
LDR [18]	ECCV'20	85.1	44.5	81.0	16.4	15.2	80.1	84.8	59.4	31.9	73.2	41.0	32.6	44.7	53.1
CD-SAM [72]	WACV'21	82.5	42.2	81.3	18.3	15.9	80.6	83.5	61.4	33.2	72.9	39.3	26.6	43.9	52.4
CLAN [32]	TPAMI'21	82.7	37.2	81.5	17.1	13.1	81.2	83.3	55.5	22.1	76.6	30.1	23.5	30.7	48.8
ASA [46]	TIP'21	91.2	48.5	80.4	5.5	5.2	79.5	83.6	56.4	21.9	80.3	36.2	20.0	32.9	49.3
DAST [73]	AAAI'21	87.1	44.5	82.3	13.9	13.1	81.6	86.0	60.3	25.1	83.1	40.1	24.4	40.5	52.5
Ours	-	<b>91.8</b>	<b>54.9</b>	<b>83.6</b>	<b>23.0</b>	<b>29.0</b>	<b>83.8</b>	<b>87.1</b>	<b>65.0</b>	26.4	<b>85.5</b>	<b>55.1</b>	36.8	<b>54.1</b>	<b>59.7</b>

TABLE IV  
ABLATION STUDY OF EACH COMPONENT IN CAMIX.

iDACS [40]	SP	CR	SRC	mIoU
✓				51.5
✓	✓			53.1
✓	✓	✓		54.1
✓	✓	✓	✓	55.2

TABLE V  
ABLATION STUDY OF THE SRC LOSS

Baseline	Mixup	$\mathcal{L}_{con}$	mIoU	$\Delta$
iDACS [40]	CMG	SRC	55.2	-
	CMG	MSE	44.5	↓ 9.7
	CMG	CE	54.2	↓ 1.0

### E. Ablation Studies

In this section, we study the effectiveness of each component (Table IV) and each level (Table VI) in our approach and investigate how they contribute to the final performance when adapting from the GTAV [9] to Cityscapes [3].

**Effectiveness of CMG:** The CMG strategy is a fundamental component of our framework, which is designed to capture the shared context dependency across domains for CAMix. *Spatial prior (SP)* and *contextual relationship (CR)* are two key components of CMG. The ablation studies of each component in CAMix are reported in Table IV. Compared to the baseline (iDACS) [40], SP and CR could successfully bring 1.6% and 2.6% of improvements, achieving 53.1% and 54.1% on the former two levels, respectively. By adding the SRC loss on the

TABLE VI  
ABLATION STUDY OF EACH LEVEL IN CAMIX.

MT	SigMask	In-Out	mIoU (GTAV)	mIoU <sub>13</sub> (SYN)
✓			43.1	45.9
✓	✓		44.6	47.1
✓		✓	54.5	59.0
✓	✓	✓	55.2	59.7

SigMask level, we can achieve an even higher performance of 55.2%.

**Effectiveness of SRC:** Table V shows the contribution of the SRC loss on the GTAV → Cityscapes benchmark. The full CAMix with all three levels and SRC loss achieves 55.2%. If we directly replace the SRC loss with a normal *mean square error (MSE)*, the result is even worse and only reaches 44.5%. Using the *cross-entropy (CE)* as the consistency loss boosts the mIoU to 54.2%, which is still 1.0% worse than our SRC loss in Eq. (9). The main benefits of the SRC loss are reflected as follows. The SigMask-level domain mixup with the SRC loss could further decrease the uncertainty of the teacher model, and promote the teacher model to transfer reasonable knowledge to the student, thus improving the performance. As such, our approach tends to be more stable and effectively ease these negative impacts, *i.e.*, training instability and early performance degradation, during the adaptation process.

**Effectiveness of different levels:** Table VI lists the impacts of different levels on the above two settings, *i.e.*, taking GTAV and SYNTHIA as the source domains, respectively. The Mean Teacher (MT) baseline achieves 43.1% and 45.9% on

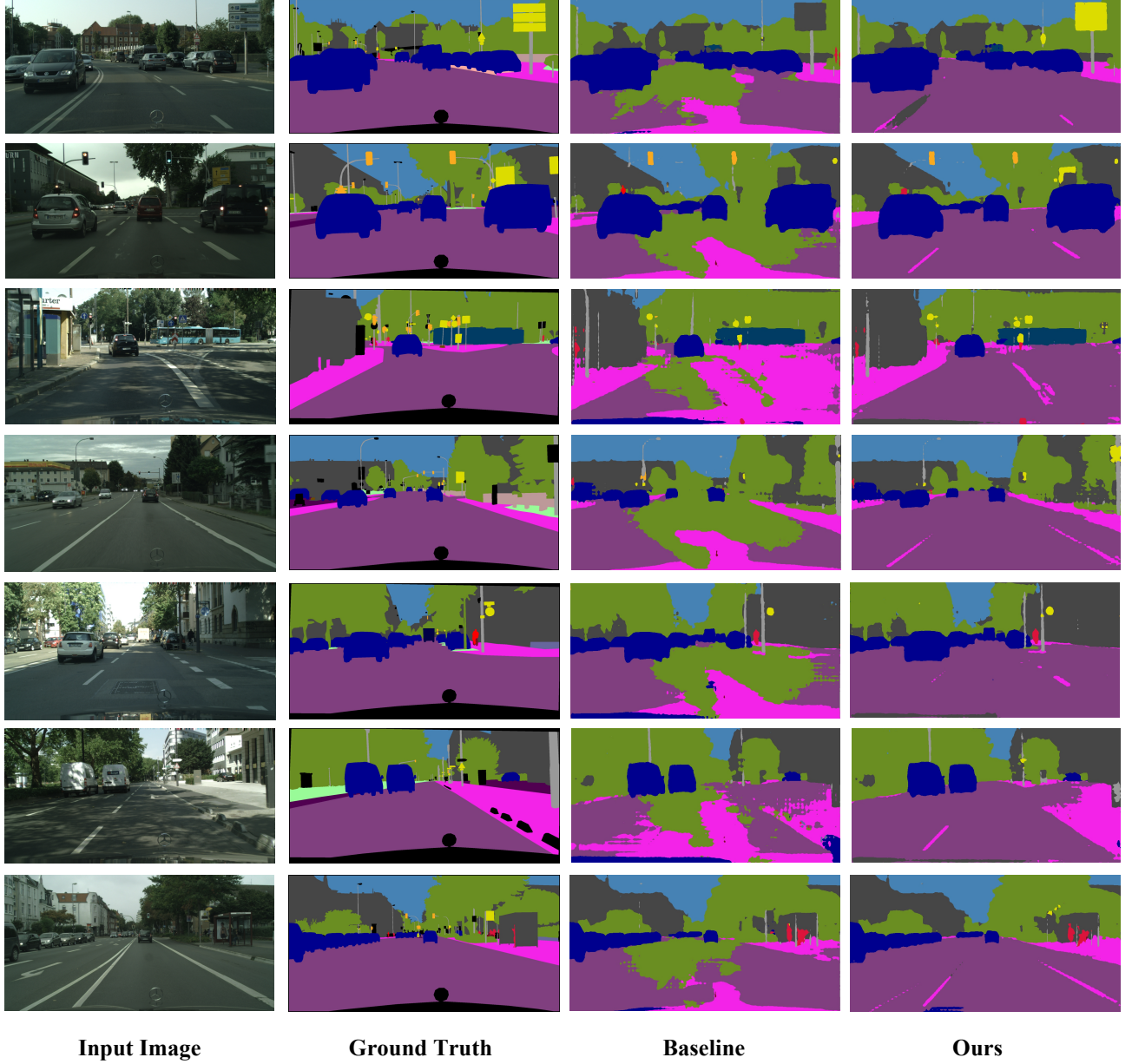


Fig. 3. Qualitative segmentation results in the SYNTHIA  $\rightarrow$  Cityscapes setup. The four columns plot (a) RGB input image, (b) ground-truth, (c) the predictions of DACS [40] and (d) the predictions of our CAMix.

two benchmarks, respectively. By adding the SigMask-level branch, our method respectively brings +1.5% and +1.2% improvements. In Table VI, In-Out means using both the input and output level mixup. By performing CAMix in the input and output level, we can successfully boost the mIoU by an additional +12.4% and +13.1%, reaching 54.5% and 59.0%, respectively. By integrating the CAMix on three levels together, we finally achieve 55.2% and 59.7% mIoU, respectively.

#### F. Visualization

**Qualitative segmentation results.** Figure 3 visualizes some segmentation results in the SYNTHIA  $\rightarrow$  Cityscapes (16

classes) set-up. The four columns plot (a) RGB input images, (b) ground truth, (c) DACS baseline outputs [40] and (d) the predictions of CAMix. As we can see from the figure, due to the lack of context dependency, DACS [40] tends to produce noisy segmentation predictions on some large categories, e.g., ‘road’, ‘sidewalk’, ‘truck’, etc, and incorrectly classifies some large categories e.g., the road as sidewalk or terrain, and produces some false predictions on some sophisticated classes, e.g., traffic sign. With the help of our proposed CAMix and SRC loss, our model manages to produce correct predictions at a high level of confidence. Fig. 3 shows that CAMix enables good performance on ‘road’, ‘sidewalk’, ‘bus’, ‘car’, ‘truck’, ‘motorcycle’, ‘bicycle’, ‘building’ and ‘terrain’ classes. Our

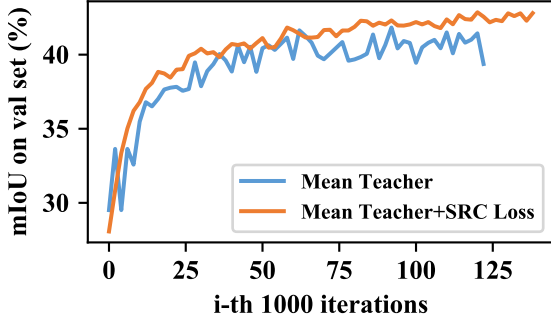


Fig. 4. Comparison on adapting from GTA5 [9] dataset to Cityscapes [3] dataset. The blue line corresponds to the conventional consistency regularization strategy [38]. The orange line indicates the consistency-based adaptation with our SRC loss. Our method can ease the issue of training instability and early performance drop.

proposed method is capable of outputting high confidence predictions compared to the previous work.

**Performance curve of adaptation.** Figure 4 plots the performance curves to show the effectiveness of SRC loss when adapting from GTAV [9] to Cityscapes [3]. Previous methods, *e.g.*, Mean Teacher [38], neglect the context knowledge shared by different domains and perform a rough distribution matching, resulting in training instability and early performance degradation. Instead, we effectively ease these negative impacts and decrease the uncertainty of segmentation model, by introducing the SRC loss.

## V. CONCLUSION

In this paper, we proposed a novel context-aware domain mixup (CAMix) framework for domain adaptive semantic segmentation. We present a contextual mask generation (CMG) strategy, which is critical for guiding the whole pipeline on three different levels, *i.e.*, input level, output level and significance mask level. Our approach can explicitly explore and transfer the shared context dependency across domains, thus narrowing down the domain gap. We also introduce a significance-reweighted consistency loss (SRC) to penalize the inconsistency between the mixed student prediction and the mixed teacher prediction, which effectively eases the adverse impacts of the adaptation, *e.g.*, training instability and early performance degradation. The extensive experiments with ablation studies demonstrate that our approach soundly outperforms the state-of-the-art methods in domain adaptive semantic segmentation.

## REFERENCES

- [1] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [2] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.
- [3] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proc. CVPR*, 2016, pp. 3213–3223.
- [4] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, 2018.
- [5] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.
- [6] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 801–818.
- [7] C. Liu, L. Chen, F. Schroff, H. Adam, W. Hua, A. L. Yuille, and F. Li, "Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 82–92.
- [8] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2881–2890.
- [9] S. R. Richter, V. Vineet, S. Roth, and V. Koltun, "Playing for data: Ground truth from computer games," in *European conference on computer vision*. Springer, 2016, pp. 102–118.
- [10] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, and A. M. Lopez, "The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3234–3243.
- [11] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, and T. Darrell, "CyCADA: Cycle-consistent adversarial domain adaptation," in *International conference on machine learning*, 2018, pp. 1989–1998.
- [12] Y. Yang and S. Soatto, "Fda: Fourier domain adaptation for semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4085–4095.
- [13] Y. Li, L. Yuan, and N. Vasconcelos, "Bidirectional learning for domain adaptation of semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6936–6945.
- [14] M. Kim and H. Byun, "Learning texture invariant representation for domain adaptation of semantic segmentation," in *Proc. CVPR*, 2020, pp. 12 975–12 984.
- [15] Y.-C. Chen, Y.-Y. Lin, M.-H. Yang, and J.-B. Huang, "Crdoco: Pixel-level domain transfer with cross-domain consistency," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1791–1800.
- [16] S. Guo, Q. Zhou, Y. Zhou, Q. Gu, J. Tang, Z. Feng, and L. Ma, "Label-free regional consistency for image-to-image translation," in *2021 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2021, pp. 1–6.
- [17] Y. Yang, D. Lao, G. Sundaramoorthi, and S. Soatto, "Phase consistent ecological domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9011–9020.
- [18] J. Yang, W. An, S. Wang, X. Zhu, C. Yan, and J. Huang, "Label-driven reconstruction for domain adaptation in semantic segmentation," in *European conference on computer vision*, vol. 12372. Springer, 2020, pp. 480–498.
- [19] Y. Luo, P. Liu, T. Guan, J. Yu, and Y. Yang, "Significance-aware information bottleneck for domain adaptive semantic segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6778–6787.
- [20] Y.-H. Chen, W.-Y. Chen, Y.-T. Chen, B.-C. Tsai, Y.-C. Frank Wang, and M. Sun, "No more discrimination: Cross city adaptation of road scene segmenters," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 1992–2001.
- [21] Y. Zou, Z. Yu, B. Kumar, and J. Wang, "Unsupervised domain adaptation for semantic segmentation via class-balanced self-training," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 289–305.
- [22] Y. Zou, Z. Yu, X. Liu, B. Kumar, and J. Wang, "Confidence regularized self-training," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 5982–5991.
- [23] Z. Lu, Y. Yang, X. Zhu, C. Liu, Y.-Z. Song, and T. Xiang, "Stochastic classifiers for unsupervised domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9111–9120.
- [24] H. Wang, T. Shen, W. Zhang, L. Duan, and T. Mei, "Classes matter: A fine-grained adversarial approach to cross-domain semantic segmentation," Springer, 2020, pp. 642–659.



- [25] Y.-H. Tsai, K. Sohn, S. Schuster, and M. Chandraker, "Domain adaptation for structured output via discriminative patch representations," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1456–1465.
- [26] T. Vu, H. Jain, M. Bucher, M. Cord, and P. Pérez, "DADA: depth-aware domain adaptation in semantic segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 7363–7372.
- [27] W.-L. Chang, H.-P. Wang, W.-H. Peng, and W.-C. Chiu, "All about structure: Adapting structural information across domains for boosting semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1900–1909.
- [28] Y.-H. Tsai, W.-C. Hung, S. Schuster, K. Sohn, M.-H. Yang, and M. Chandraker, "Learning to adapt structured output space for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7472–7481.
- [29] Y. Luo, L. Zheng, T. Guan, J. Yu, and Y. Yang, "Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2507–2516.
- [30] Z. Wang, M. Yu, Y. Wei, R. Feris, J. Xiong, W.-m. Hwu, T. S. Huang, and H. Shi, "Differential treatment for stuff and things: A simple unsupervised domain adaptation method for semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 635–12 644.
- [31] Y. Chen, W. Li, X. Chen, and L. V. Gool, "Learning semantic segmentation from synthetic data: A geometrically guided input-output adaptation approach," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1841–1850.
- [32] Y. Luo, P. Liu, L. Zheng, T. Guan, J. Yu, and Y. Yang, "Category-level adversarial adaptation for semantic segmentation using purified features," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2021.
- [33] F. Pan, I. Shin, F. Rameau, S. Lee, and I. S. Kweon, "Unsupervised intra-domain adaptation for semantic segmentation through self-supervision," in *Unsupervised Intra-domain Adaptation for Semantic Segmentation through Self-Supervision*, 2020, pp. 3764–3773.
- [34] J. Yang, R. Xu, R. Li, X. Qi, X. Shen, G. Li, and L. Lin, "An adversarial perturbation oriented domain adaptation approach for semantic segmentation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 12 613–12 620.
- [35] I. Shin, S. Woo, F. Pan, and I. S. Kweon, "Two-phase pseudo label densification for self-training based domain adaptation," in *European conference on computer vision*. Springer, 2020, pp. 532–548.
- [36] K. Mei, C. Zhu, J. Zou, and S. Zhang, "Instance adaptive self-training for unsupervised domain adaptation," in *European conference on computer vision*. Springer, 2020, pp. 415–430.
- [37] M. Naseer Subhani and M. Ali, "Learning from scale-invariant examples for domain adaptation in semantic segmentation," in *European conference on computer vision*. Springer, 2020, pp. 290–306.
- [38] J. Choi, T. Kim, and C. Kim, "Self-ensembling with gan-based data augmentation for domain adaptation in semantic segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6830–6840.
- [39] Y. Xu, B. Du, L. Zhang, Q. Zhang, G. Wang, and L. Zhang, "Self-ensembling attention networks: Addressing domain shift for semantic segmentation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 5581–5588.
- [40] W. Tranheden, V. Olsson, J. Pinto, and L. Svensson, "Dacs: Domain adaptation via cross-domain mixed sampling," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 1379–1389.
- [41] Q. Zhou, Z. Feng, Q. Gu, G. Cheng, X. Lu, J. Shi, and L. Ma, "Uncertainty-aware consistency regularization for cross-domain semantic segmentation," *arXiv preprint arXiv:2004.08878*, 2020.
- [42] L. Melas-Kyriazi and A. K. Manrai, "Pixmatch: Unsupervised domain adaptation via pixelwise consistency training," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12 435–12 445.
- [43] N. Araslanov and S. Roth, "Self-supervised augmentation consistency for adapting semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 384–15 394.
- [44] T.-H. Vu, H. Jain, M. Bucher, M. Cord, and P. Pérez, "Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2517–2526.
- [45] M. Chen, H. Xue, and D. Cai, "Domain adaptation for semantic segmentation with maximum squares loss," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2090–2099.
- [46] W. Zhou, Y. Wang, J. Chu, J. Yang, X. Bai, and Y. Xu, "Affinity space adaptation for semantic segmentation across domains," *IEEE Transactions on Image Processing*, vol. 30, pp. 2549–2561, 2020.
- [47] X. Zhu, H. Zhou, C. Yang, J. Shi, and D. Lin, "Penalizing top performers: Conservative loss for semantic segmentation adaptation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 568–583.
- [48] Q. Gu, Q. Zhou, M. Xu, Z. Feng, G. Cheng, X. Lu, J. Shi, and L. Ma, "Pit: Position-invariant transform for cross-fov domain adaptation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021.
- [49] Q. Zhou, Q. Gu, J. Pang, Z. Feng, G. Cheng, X. Lu, J. Shi, and L. Ma, "Self-adversarial disentangling for specific domain adaptation," *arXiv preprint arXiv:2108.03553*, 2021.
- [50] G. French, S. Laine, T. Aila, S. Laine, M. Mackiewicz, and G. Finlayson, "Semi-supervised semantic segmentation needs strong, varied perturbations," in *British Machine Vision Conference*, 2020.
- [51] G. French, A. Oliver, and T. Salimans, "Milking cowmask for semi-supervised image classification," *arXiv preprint arXiv:2003.12022*, 2020.
- [52] J. Zhang, L. Chen, B. Ouyang, B. Liu, J. Zhu, Y. Chen, Y. Meng, and D. Wu, "Pointcutmix: Regularization strategy for point cloud classification," *arXiv preprint arXiv:2101.01461*, 2021.
- [53] Y. Chen, V. T. Hu, E. Gavves, T. Mensink, P. Mettes, P. Yang, and C. G. Snoek, "Pointmixup: Augmentation for point clouds," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*. Springer, 2020, pp. 330–345.
- [54] M. Xu, J. Zhang, B. Ni, T. Li, C. Wang, Q. Tian, and W. Zhang, "Adversarial domain adaptation with domain mixup," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 6502–6509.
- [55] Y. Wu, D. Inkpen, and A. El-Roby, "Dual mixup regularized learning for adversarial domain adaptation," in *European Conference on Computer Vision*. Springer, 2020, pp. 540–555.
- [56] X. Mao, Y. Ma, Z. Yang, Y. Chen, and Q. Li, "Virtual mixup training for unsupervised domain adaptation," *arXiv preprint arXiv:1905.04215*, 2019.
- [57] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *IEEE Proc.*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [58] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," in *NeurIPS workshop*, 2011.
- [59] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *Advances in Neural Information Processing Systems* 30, 2017, pp. 1195–1204.
- [60] G. French, M. Mackiewicz, and M. Fisher, "Self-ensembling for visual domain adaptation," in *Proceedings of the International Conference on Learning Representations*, 2018.
- [61] C. S. Perone, P. Ballester, R. C. Barros, and J. Cohen-Adad, "Unsupervised domain adaptation for medical imaging segmentation with self-ensembling," *NeuroImage*, vol. 194, pp. 1–11, 2019.
- [62] V. Olsson, W. Tranheden, J. Pinto, and L. Svensson, "Classmix: Segmentation-based data augmentation for semi-supervised learning," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 1369–1378.
- [63] Z. Feng, Q. Zhou, Q. Gu, X. Tan, G. Cheng, X. Lu, J. Shi, and L. Ma, "Dmt: Dynamic mutual training for semi-supervised learning," *arXiv preprint arXiv:2004.08514*, 2020.
- [64] Z. Feng, Q. Zhou, G. Cheng, X. Tan, J. Shi, and L. Ma, "Semi-supervised semantic segmentation via dynamic self-training and classbalanced curriculum," *arXiv preprint arXiv:2004.08514*, vol. 1, no. 2, p. 5, 2020.
- [65] L. Han, Y. Zou, R. Gao, L. Wang, and D. Metaxas, "Unsupervised domain adaptation via calibrating uncertainties," in *CVPR Workshops*, 2019.
- [66] A. Saporta, T.-H. Vu, M. Cord, and P. Pérez, "Esl: Entropy-guided self-supervised learning for domain adaptation in semantic segmentation," in *CVPR Workshops*, 2020.
- [67] J. Wen, N. Zheng, J. Yuan, Z. Gong, and C. Chen, "Bayesian uncertainty matching for unsupervised domain adaptation," *arXiv preprint arXiv:1906.09693*, 2019.

- [68] V. K. Kurmi, S. Kumar, and V. P. Namboodiri, "Attending to discriminative certainty for domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 491–500.
- [69] S. Paul, Y. Tsai, S. Schuler, A. K. Roy-Chowdhury, and M. Chandraker, "Domain adaptive semantic segmentation using weak labels," in *European conference on computer vision*, vol. 12354. Springer, 2020, pp. 571–587.
- [70] J. Huang, S. Lu, D. Guan, and X. Zhang, "Contextual-relation consistent domain adaptation for semantic segmentation," in *European conference on computer vision*, vol. 12360. Springer, 2020, pp. 705–722.
- [71] G. Li, G. Kang, W. Liu, Y. Wei, and Y. Yang, "Content-consistent matching for domain adaptive semantic segmentation," in *European conference on computer vision*, vol. 12359. Springer, 2020, pp. 440–456.
- [72] J. Yang, W. An, C. Yan, P. Zhao, and J. Huang, "Context-aware domain adaptation in semantic segmentation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 514–524.
- [73] F. Yu, M. Zhang, H. Dong, S. Hu, B. Dong, and L. Zhang, "Dast: Unsupervised domain adaptation in semantic segmentation based on discriminator attention and self-training," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 12, 2021, pp. 10754–10762.
- [74] J. Hoffman, D. Wang, F. Yu, and T. Darrell, "Fcns in the wild: Pixel-level adversarial and constraint-based adaptation," *CoRR*, vol. abs/1612.02649, 2016.
- [75] Y. Zhang, P. David, and B. Gong, "Curriculum domain adaptation for semantic segmentation of urban scenes," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2020–2030.
- [76] Y. Chen, W. Li, and L. Van Gool, "Road: Reality oriented adaptation for semantic segmentation of urban scenes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7892–7901.
- [77] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [78] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [79] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimeshine, L. Antiga, A. Desmaison, A. Köpf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Proc. NeurIPS*, 2019, pp. 8024–8035.
- [80] Z. Zheng and Y. Yang, "Unsupervised scene adaptation with memory regularization in vivo," in *International Joint Conference on Artificial Intelligence*, 2020, pp. 1076–1082.



**Qianyu Zhou** is currently pursuing his Ph.D. degree in the Department of Computer Science and Engineering, Shanghai Jiao Tong University. Before that, he received a B.Sc. degree in Jilin University in 2019. His current research interests focus on computer vision, scene understanding, domain adaptation.



**Zhengyang Feng** is currently pursuing his M.Sc. degree in the Department of Computer Science and Engineering, Shanghai Jiao Tong University. Before that, he received a B.E. degree in information security from Harbin Institute of Technology, Weihai, China, in 2020. His current research interests focus on pattern recognition with limited human supervision.



**Qiqi Gu** received a B.E. degree in Shanghai Jiao Tong University in 2015 and is now a second-year master student in Department of Computer Science and Engineering, Shanghai Jiao Tong University. Her current research interests focus on domain adaptation of object detection and semantic segmentation.



**Jiangmiao Pang** is currently a Postdoctoral Research Fellow at Multimedia Laboratory, the Chinese University of Hong Kong. He obtained his Ph.D. degree from Zhejiang University in 2021. His research interests include computer vision and robotics, especially their applications in autonomous driving.



**Guangliang Cheng** is currently a Senior Research Manager in SenseTime. Before that, he was a Post-doc researcher in the Institute of Remote Sensing and Digital Earth, Chinese Academy of Sciences, China, and he received his Ph.D. degree with national laboratory of pattern recognition (NLPR) from the Institute of Automation, Chinese Academy of Sciences, Beijing. His research interests include autonomous driving, scene understanding, domain adaptation and remote sensing image processing.



**Xuequan Lu** is an Assistant Professor at the School of Information Technology, Deakin University, Australia. He spent more than two years as a Research Fellow in Singapore. Prior to that, he earned his Ph.D at Zhejiang University (China) in June 2016. His research interests mainly fall into the category of visual computing, for example, geometry modeling, processing and analysis, animation/simulation, 2D data processing and analysis. More information can be found at <http://www.xuequanlu.com>.



**Jianping Shi** is an Executive Research Director at SenseTime. Currently her team works on developing algorithms for autonomous driving, scene understanding, remote sensing, etc. She got her Ph.D. degree in Computer Science and Engineering Department in the Chinese University of Hong Kong in 2015 under the supervision of Prof. Jiaya Jia. Before that, she received the B. Eng degree from Zhejiang University in 2011. She has served regularly on the organization committees of numerous conferences, such as Area Chair of CVPR 2020, ICCV 2019, etc.



**Lizhuang Ma** received his B.S. and Ph.D. degrees from the Zhejiang University, China in 1985 and 1991, respectively. He is now a Distinguished Professor, Ph.D. Tutor, and the Head of the Digital Media and Computer Vision Laboratory at the Department of Computer Science and Engineering, Shanghai Jiao Tong University, China. He was a Visiting Professor at the Frounhofer IGD, Darmstadt, Germany in 1998, and was a Visiting Professor at the Center for Advanced Media Technology, Nanyang Technological University, Singapore from 1999 to 2000. He has published more than 200 academic research papers in both domestic and international journals. His research interests include computer aided geometric design, computer graphics, computer vision, scientific data visualization, computer animation, digital media technology, and theory and applications for computer graphics, CAD/CAM. He serves as the reviewer of IEEE Transactions on Pattern Analysis and Machine Intelligence, IEEE Transactions on Image Processing, IEEE Transactions on Multimedia, etc.