

# Macro-Micro Adversarial Network for Human Parsing

Yawei Luo<sup>1,2</sup>, Zhedong Zheng<sup>2</sup>, Liang Zheng<sup>2,3</sup>, Tao Guan<sup>1</sup>, Junqing Yu<sup>1</sup>, and  
Yi Yang<sup>2</sup>

<sup>1</sup> School of Computer Science and Technology,  
Huazhong University of Science and Technology

{royalvane,qd.gt,yjqing}@hust.edu.cn

<sup>2</sup> CAI, University of Technology Sydney

<sup>3</sup> Singapore University of Technology and Design  
{zdzheng12,liangzheng06,yee.i.yang}@gmail.com

**Abstract.** In human parsing, the pixel-wise classification loss has drawbacks in its low-level local inconsistency and high-level semantic inconsistency. The introduction of the adversarial network tackles the two problems using a single discriminator. However, the two types of parsing inconsistency are generated by distinct mechanisms, so it is difficult for a single discriminator to solve them both. To address the two kinds of inconsistencies, this paper proposes the Macro-Micro Adversarial Net (MMAN). It has two discriminators. One discriminator, Macro  $D$ , acts on the low-resolution label map and penalizes semantic inconsistency, *e.g.*, misplaced body parts. The other discriminator, Micro  $D$ , focuses on multiple patches of the high-resolution label map to address the local inconsistency, *e.g.*, blur and hole. Compared with traditional adversarial networks, MMAN not only enforces local and semantic consistency explicitly, but also avoids the poor convergence problem of adversarial networks when handling high resolution images. In our experiment, we validate that the two discriminators are complementary to each other in improving the human parsing accuracy. The proposed framework is capable of producing competitive parsing performance compared with the state-of-the-art methods, *i.e.*, mIoU=46.81% and 59.91% on LIP and PASCAL-Person-Part, respectively. On a relatively small dataset PPSS, our pre-trained model demonstrates impressive generalization ability. The code is publicly available at <https://github.com/RoyalVane/MMAN>.

**Keywords:** Human parsing, Adversarial network, Inconsistency, Macro-Micro

## 1 Introduction

Human parsing aims to segment a human image into multiple semantic parts. It is a pixel-level prediction task which requires to understand human images in both the global level and the local level. Human parsing can be widely applied to human behavior analysis [9], pose estimation [34] and fashion synthesis [40].

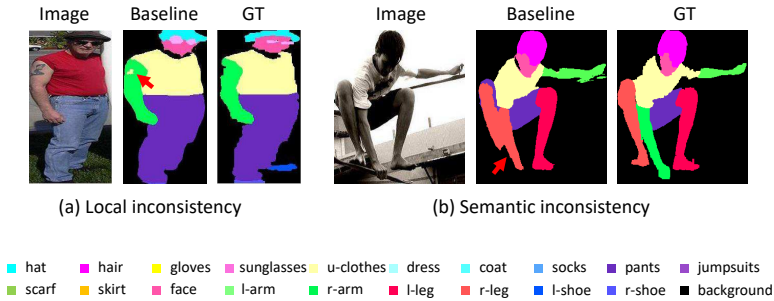


Fig. 1: Drawbacks of the pixel-wise classification loss. (a) Local inconsistency, which leads to a hole on the arm. (b) Semantic inconsistency, which causes unreasonable human poses. The inconsistencies are indicated by red arrows.

Recent advances in human parsing and semantic segmentation [19,34,10,23,37,36] mostly explore the potential of the convolutional neural network (CNN).

Based on CNN architecture, the *pixel-wise classification loss* is usually used [19,34,10] which punishes the classification error for each pixel. Despite providing an effective baseline, the pixel-wise classification loss which is designed for per-pixel category prediction, has two drawbacks. First, the pixel-wise classification loss may lead to *local inconsistency*, such as holes and blur. The reason is that it merely penalizes the false prediction on every pixel without explicitly considering the correlation among the adjacent pixels. For illustration, we train a baseline model (see Section 3.2) with the pixel-wise classification loss. As shown in Fig. 1(a), some pixels which belongs to “arm” are incorrectly predicted as “upper-clothes” by the baseline. This is undesirable but is the consequence of local inconsistency of the baseline loss. Second, pixel-wise classification loss may lead to *semantic inconsistency* in the overall segmentation map, such as unreasonable human poses and incorrect spatial relationship of body parts. Compared to the local inconsistency, the semantic inconsistency is generated from deeper layers. When only looking at a local region, the learned model does not have an overall sense of the topology of body parts. As shown in Fig. 1(b), the “arm” is merged with an adjacent “leg”, indicating incorrect part topology (three legs). Therefore, the pixel-wise classification loss does not explicitly consider the semantic consistency, so that long-range dependency may not be well captured.

In the attempt to address the inconsistency problems, the conditional random fields (CRFs) [17] can be employed as a post processing method. However, CRFs usually handle inconsistency in very limited scope (locally) due to the pairwise potentials, and may even generate worse label maps given poor initial segmentation result. As an alternative to CRFs, a recent work proposes the use of adversarial network [24]. Since the adversarial loss assesses whether a label map is real or fake by joint configuration of many label variables, it can enforce higher-level consistency, which cannot be achieved with pairwise terms or the per-pixel classification loss. Now, an increasing number of works adopt the routine

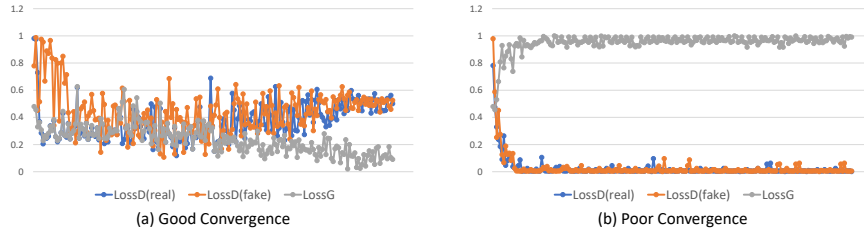


Fig. 2: Two types of convergence in adversarial network training.  $LossD(real)$  and  $LossD(fake)$  denote the adversarial losses of discriminator on real and fake image respectively, and  $LossG$  denotes the loss of generator. **(a)** Good convergence, where  $LossD(real)$  and  $LossD(fake)$  converge to 0.5 and  $LossG$  converges to 0. It indicates a successful adversarial network training, where  $G$  is able to fool  $D$ . **(b)** Poor convergence, where  $LossD(real)$  and  $LossD(fake)$  converge to 0 and  $LossG$  converges to 1. It stands for an unbalanced adversarial network training, where  $D$  can easily distinguish generated images from real images.

of combining the cross entropy loss with an adversarial loss to produce label maps closer to the ground truth [5,27,12].

Nevertheless, the previous adversarial network also has its limitations. First, the single discriminator back propagates only one adversarial loss to the generator. However, the local inconsistency is generated from top layers and the semantic inconsistency is generated from deep layers. The two targeted layers can not be discretely trained with only one adversarial loss. Second, a single discriminator has to look at overall high-resolution image (or a large part of it) in order to supervise the global consistency. As mentioned by numbers of literatures [7,14], it is very difficult for a generator to fool the discriminator on a high-resolution image. As a result, the single discriminator back propagates a maximum adversarial loss invariably, which makes the training unbalanced. We call it *poor convergence problem*, as shown in Fig. 2.

In this paper, the basic objective is to improve the local and semantic consistency of label maps in human parsing. We adopt the idea of adversarial training and at the same time aim to address its limitations, *i.e.*, the inferior ability in improving parsing consistency with a single adversarial loss and the poor convergence problem. Specifically, we introduce the Macro-Micro Adversarial Nets (MMAN). MMAN consists of a dual-output generator ( $G$ ) and two discriminators ( $D$ ), named Macro  $D$  and Micro  $D$ . The three modules constitute two adversarial networks (Macro AN, Micro AN), addressing the semantic consistency and the local consistency, respectively. Given an input human image, the CNN-based generator outputs two segmentation maps with different resolution levels, *i.e.*, low resolution and high resolution. The input of Macro  $D$  is a low-resolution segmentation map, and the output is the confidence score of semantic consistency. The input of Micro  $D$  is the high-resolution segmentation result, and its outputs

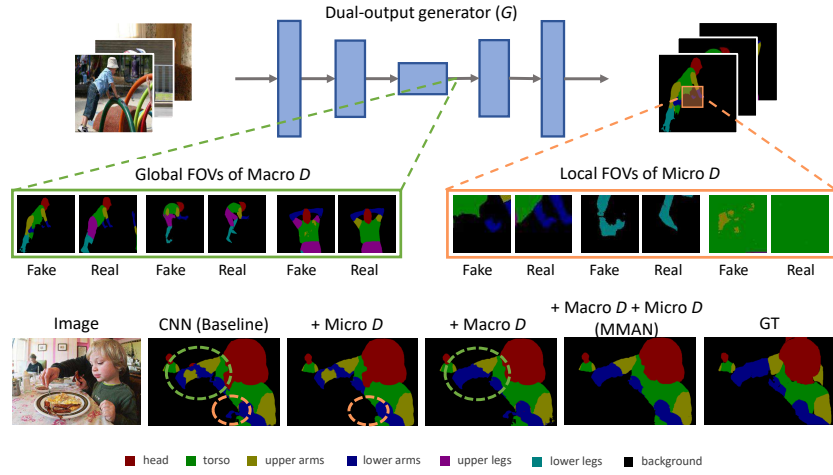


Fig. 3: **Top:** A brief pipeline of MMAN. Two discriminators are attached to a CNN-based generator ( $G$ ). The Macro  $D$  works on the low-resolution label map and has a global receptive field, focusing on semantic consistency. Micro  $D$  focuses on multiple patches and has small receptive fields on high-resolution label map, thus supervising the local consistency. The Macro (Micro) discriminator yields “fake” if semantic (local) inconsistency is observed, otherwise it gives “real”. **Bottom:** qualitative results of using Macro  $D$ , Micro  $D$  and MMAN, respectively. We observe that Macro  $D$  and Micro  $D$  correct semantic inconsistency (green dashed circle) and local inconsistency (orange dashed circle), respectively, and that MMAN possesses the merits of both.

is the confidence score of local consistency. A brief pipeline of the proposed framework is shown in Fig. 3. It is in two critical aspects that MMAN departs from previous works. First, our method explicitly copes with the local inconsistency and semantic inconsistency problem using two task-specific adversarial networks individually. Second, our method does not use large-sized FOVs on high-resolution image, so we can avoid the poor convergence problem. More detailed description of the merits of the proposed network is provided in Section 3.5.

Our contributions are summarized as follows:

- We propose a new framework called Macro-Micro Adversarial Network (MMAN) for human parsing. The Macro  $AN$  and Micro  $AN$  focus on semantic and local inconsistency respectively, and work in complementary way to improve the parsing quality.
- The two discriminators in our framework achieve local and global supervision on the label maps with small field of views (FOVs), which avoids the poor convergence problem caused by high-resolution images.

- The proposed adversarial net achieves very competitive mIoU on the LIP and PASCAL-Person-Part datasets, and can be well generalized on a relatively small dataset PPSS.

## 2 Related works

Our review focuses on three lines of literature most relevant to our work, *i.e.*, CNN-based human parsing, the conditional random fields (CRFs) and the adversarial networks.

**Human parsing.** Recent progress in human parsing has been due to the two factors: 1) the available of the large-scale datasets [10,19,25,4]. Comparing to the small datasets, the large-scale datasets contain the common visual variance of people and provide a comprehensive evaluation. 2) the end-to-end learned model. Human parsing demands understanding the person on the pixel level. The recent works apply the convolutional neural network (CNN) to learn the segmentation result in an end-to-end manner. In [34], human poses are extracted in advance and utilized as strong structural cues to guide the parsing. In [21], four human-related contexts are integrated into a unified network. A novel human-related grammar is presented by [29] which infers human body pose and human part segmentation jointly.

**Conditional random fields** Using the pixel-wise classification loss, CNN usually ignores the micro context between pixels and the macro context between semantic parts. Conditional random fields (CRFs) [17,22,18] are one of the common methods to enforce spatial contiguity in the output label maps. Served as a post-process procedure for image segmentation, CRFs further fine-tune the output map. However, the most common used CRFs are with pair-wise potentials [2,26], which has very limited parameters and handles low-level inconsistencies with a small scope. Higher-order potentials [16,18] have also been observed to be effective in enforcing the semantic validity, but the corresponding energy pattern and the clique form are usually difficult to design. In summary, the utilization of context in CNN remains an open problem.

**Adversarial networks.** Adversarial networks have demonstrated the effectiveness in image synthesis [13,28,30,39,38]. By minimizing the adversarial loss, the discriminator leads the generator to produce high-fidelity images. In [24], Luc *et al.* add the adversarial loss for training semantic segmentation and yield the competitive results. Similar idea then has been applied in street scene segmentation [12] and medical image segmentation [5,27]. Contemporarily, an increasing body of literature [7,14] report the difficulty of training the adversarial networks on the high-resolution images. Discriminator can easily recognize the fake high-resolution image, which leads to the training unbalance. The generator and discriminator are prone to stuck in a local minimum.

The main difference between MMAN and the adversarial learning methods above is that the we explicitly endow adversarial training with the macro and micro subtasks. We observe that the two subtasks are complementary to each

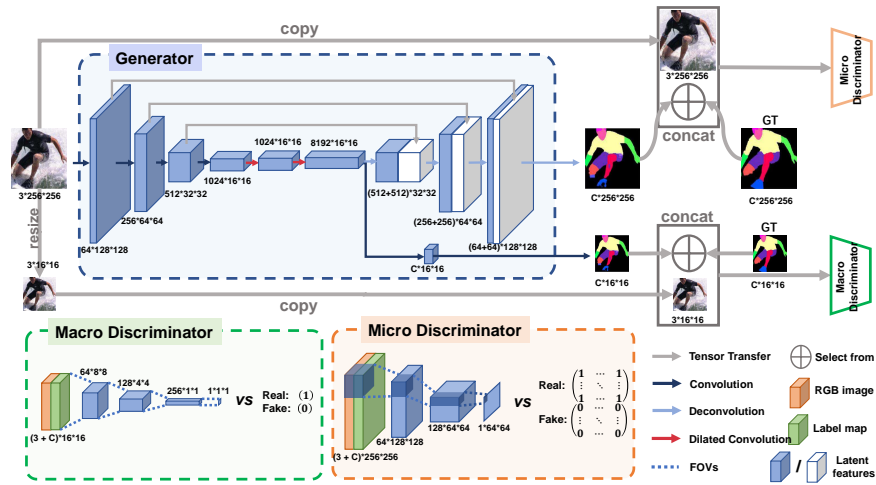


Fig. 4: MMAN has three components: a dual-output generator (blue dashed box), a Macro discriminator (green dashed box) and a Micro discriminator (orange dashed box). Given an input image of size  $3 \times 256 \times 256$ , the generator  $G$  first produces a low-resolution ( $8192 \times 16 \times 16$ ) tensor, from which a low-resolution label map ( $C \times 16 \times 16$ ) and a high-resolution label map ( $C \times 256 \times 256$ ) are generated, where  $C$  is the number of classes. Finally, for the each label map (sized  $C \times 16 \times 16$ , for example), we concatenate it with an RGB image (sized  $3 \times 16 \times 16$ ) along the 1st axis (number of channels), which is fed into the corresponding discriminator.

other to achieve superior parsing accuracy to the baseline with a single adversarial loss and are able to reduce the risk of the training unbalance.

### 3 Macro-Micro Adversarial Network

Figure 4 illustrates the architecture of the proposed Macro-Micro Adversarial Network. The network consists of three components, *i.e.*, a dual-output generator ( $G$ ) and two task-specific discriminators ( $D_{Ma}$  and  $D_{Mi}$ ). Given an input image of size  $3 \times 256 \times 256$ ,  $G$  outputs two label maps of size  $C \times 16 \times 16$  and  $C \times 256 \times 256$ , respectively.  $D_{Ma}$  supervises the entire label map of  $C \times 16 \times 16$  and  $D_{Mi}$  focuses on patches of the label map of size  $C \times 256 \times 256$ , respectively, so that global and local inconsistencies are penalized. In Section 3.1, we illustrate the training objectives, followed by the structure illustration in Section 3.2, 3.3 and 3.4. The merits of the proposed network are discussed in Section 3.5.

#### 3.1 Training Objectives

Given a human image  $x$  of shape  $3 \times H \times W$  and a target label map  $y$  of shape  $C \times H \times W$  where  $C$  is the number of classes including the background, the

traditional pixel-wise classification loss (multi-class cross-entropy loss) can be formulated as:

$$\mathcal{L}_{mce}(G) = \sum_{i=1}^{H \times W} \sum_{c=1}^C -y_{ic} \log \hat{y}_{ic}, \quad (1)$$

where  $\hat{y}_{ic}$  denotes the predicted probability of the class  $c$  on the  $i$ -th pixel. The  $y_{ic}$  denotes the ground truth probability of the class  $c$  on the  $i$ -th pixel. If the  $i$ -th pixel belongs to class  $c$ ,  $y_{ic} = 1$ , else  $y_{ic} = 0$ .

To enforce the spatial consistency, we combine the pixel-wise classification loss with the adversarial loss. It can be formulated as:

$$\mathcal{L}_{mix}(G, D) = \mathcal{L}_{mce}(G) + \lambda \mathcal{L}_{adver}(G, D), \quad (2)$$

where  $\lambda$  controls the relative importance of the pixel-wise classification loss and the adversarial loss. Specifically, the adversarial loss  $\mathcal{L}_{adver}(G, D)$  is:

$$\mathcal{L}_{adver}(G, D) = \mathbb{E}_{x,y} [\log D(x, y)] + \mathbb{E}_x [\log(1 - D(x, G(x)))]. \quad (3)$$

As shown in Fig. 4, the proposed MMAN employs the “cross-entropy loss + adversarial loss” to supervise both the bottom and top output from the generator  $G$ :

$$\mathcal{L}_{MMAN}(G, D_{Ma}, D_{Mi}) = \mathcal{L}_{adver}(G, D_{Ma}) + \lambda_1 \mathcal{L}_{mce_l}(G) + \lambda_2 \mathcal{L}_{adver}(G, D_{Mi}) + \lambda_3 \mathcal{L}_{mce_h}(G), \quad (4)$$

where  $\mathcal{L}_{mce_l}(G)$  donates the cross-entropy loss between the low-resolution output and the small-sized target label map, while the  $\mathcal{L}_{mce_h}(G)$  refers to the cross-entropy loss between the high-resolution output and the original ground-truth label map. Similarly,  $\mathcal{L}_{adver}(G, D_{Ma})$  is the adversarial loss focusing on the low-resolution map, and  $\mathcal{L}_{adver}(G, D_{Mi})$  is based on the high-resolution map. The hyper parameters  $\lambda_1$ ,  $\lambda_2$  and  $\lambda_3$  control the relative importance of the four losses. The training objective of MMAN is:

$$G^*, D_{Ma}^*, D_{Mi}^* = \arg \min_G \max_{D_{Ma}, D_{Mi}} \mathcal{L}_{MMAN}(G, D_{Ma}, D_{Mi}). \quad (5)$$

We solve Eq. 5 by alternate between optimizing  $G$ ,  $D_{Ma}$  and  $D_{Mi}$  until  $\mathcal{L}_{MMAN}(G, D_{Ma}, D_{Mi})$  converges.

### 3.2 Dual-output Generator

For the generator ( $G$ ), we utilize DeepLab-ASPP [2] framework with ResNet-101 [11] model pre-trained on the ImageNet dataset [6] as our starting point due to its simplicity and effectiveness. We augment DeepLab-ASPP architecture with cascaded upsampling layers and skip connect them with early layers, which is

similar with U-net [31]. Furthermore, we add a bypass to output the deep feature tensor from the bottom layers and transfer it to a label map with a convolution layer. The small-sized label map serves as the second output in parallel with the original sized label map from the top layer. We refer to the augmented dual-output architecture as Do-DeepLab-ASPP and adopt it as our baseline. For the dual output, we supervise the cross-entropy loss from top layers with ground truth label maps of original size, since it can retain visual details. Besides, we supervise the cross-entropy loss of bottom layers with a resized label map, *i.e.*, 1/16 times of the original size. The shrunken label map pays more attentions to the coarse-grained human structure. The same strategy is applied to adversarial loss. We concatenated the respect label map with RGB image of corresponding size along class channel as a strong condition to discriminators.

### 3.3 Macro Discriminator

Macro discriminator ( $D_{Ma}$ ) aims to lead the generator to produce realistic label map that consist with high-level human characteristics, such as reasonable human poses and correct spatial relationship of body parts.  $D_{Ma}$  is attached to the bottom layer of  $G$  and focuses on an overall low-resolution label map. It consists of 4 convolution layers with kernel size of  $4 \times 4$  and stride of 2. Each convolution layer follows by one instance-norm layer and one LeakyRelu function. Given a output label map from  $G$ ,  $D_{Ma}$  downsamples it to  $1 \times 1$  to achieve the global supervision on it. The output of  $D_{Ma}$  is the confidence score of semantic consistency.

### 3.4 Micro Discriminator

Micro discriminator ( $D_{Mi}$ ) is designed to enforce the local consistency in label maps. We follow the idea of ‘‘PatchGAN’’ [13] in designing the  $D_{Mi}$ . Different from  $D_{Ma}$  that has a global receptive field on the (shrunken) label map,  $D_{Mi}$  only penalizes local error at the scale of image patches. The kernel size of  $D_{Mi}$  is  $4 \times 4$  and the stride is 2. Micro  $D$  has a shallow structure of 3 convolution layers, each convolution layer follows by one instance-norm layer and one LeakyRelu function.  $D_{Mi}$  aims to classify if each  $22 \times 22$  patch in an high-resolution image is real or fake, which is suitable for enforcing the local consistency. After running  $D_{Mi}$  convolutionally across the label map, we will obtain multiple response from every receptive field. We finally averages all responses to provide the ultimate output of  $D_{Mi}$ .

### 3.5 Discussions

In CNN-based human parsing, convolution layers go deep to extract part-level features, and deconvolution layers bring the in-depth features back to pixel-level locations. It seems intuitive to arrange the Macro  $D$  to deeper layers to supervise high-level semantic features and Micro  $D$  to top layers, focusing on low-level visual features. Besides the intuitive motivation, however, we can benefit more from such arrangement. The merits of MMAN are summarized in four aspects.



**Functional specialization of Macro  $D$  and Micro  $D$ .** Compared with the single discriminator which attempts to solve two levels of inconsistency alone, Macro  $D$  and Micro  $D$  are specified in addressing one of the two consistency problems. Take Macro  $D$  as an example. First, Macro  $D$  is attached to the deep layer of  $G$ . Because the semantic inconsistency is originally generated from the deep layers, a such designed Macro  $D$  allows the loss to back propagated to  $G$  more directly. Second, Macro  $D$  acts on a low-resolution label map that retains the semantic-level human structure while filtering out the pixel-level details. It enforces Macro  $D$  to focus on the global inconsistency without disturbing by local errors. The same reasoning applies to Micro  $D$ . In section 4.5, we validate that MMAN consistently outperforms the adversarial networks with a single adversarial loss [24,5].

**Functional complementarity of Macro  $D$  and Micro  $D$ .** As mentioned in [35], supervising classification loss in early deep layers can offer a good coarse-grained initialization for later top layers. Correspondingly, decreasing the loss in top layers can remedy the coarse semantic feature with fine-grained visual details. We assume that the adversarial loss has the same characteristic to work in complementary pattern. We clarify our hypothesis in Section 4.4.

**Small FOVs to avoid poor convergence problem.** Reported by increasing literatures [7,14], the existing adversarial networks have drawbacks in coping with complex high-resolution images. In our framework, Macro  $D$  acts on a low-resolution label map and Micro  $D$  has multiple but small FOVs on a high-resolution label map. As a result, both Macro  $D$  and Micro  $D$  avoid using large FOVs as the actual input, which effectively reduce the convergence risk caused by high resolution. We show this benefit in Section 4.5.

**Efficiency.** Comparing with the single adversarial network [24,5], MMAN achieves the supervision across the overall images with two shallower discriminators, which have fewer parameters. It also owing to the small FOVs of the discriminators. The efficiency of MMAN is showed in variant study in Section 4.5.

## 4 Experiment

### 4.1 Dataset

**LIP** [10] is a recently introduced large-scale dataset, challenging in the severe pose complexity, heavy occlusions and body truncation. It contains 50,462 images in total, including 30,362 for training, 10,000 for testing and 10,000 for validation. LIP defines 19 human part (clothes) labels, including hat, hair, sunglasses, upper-clothes, dress, coat, socks, pants, gloves, scarf, skirt, jumpsuits, face, right arm, left arm, right leg, left leg, right shoe and left shoe, and a background class.

**PASCAL-Person-Part** [4] annotates the human part segmentation labels and is a subset of PASCAL-VOC 2010 [8]. PASCAL-Person-Part includes 1,716 images for training and 1,817 for testing. In this dataset, an image may contain multiple persons with unconstrained poses and environment. Six human body part classes and the background class are annotated.

**PPSS** [25] includes 3,673 annotated samples, which are divided into a training set of 1,781 images and a testing set of 1,892 images. It defines seven human parts and a background class. Collected from 171 surveillance videos, the dataset can reflect the occlusion and illumination variation in real scene.

**Evaluation metric.** The human parsing accuracy of each class is measured in terms of pixel intersection-over-union (IoU). The mean intersection-over-union (mIoU) is computed by averaging the IoU across all classes. We use both IoU for each class and mIoU as evaluation metrics for each dataset.

## 4.2 Implementation Details

In our implementation, input images are resized so that its shorter side is fixed to 288. A  $256 \times 256$  crop is randomly sampled from the image or its horizontal flipped version. The per-pixel mean is subtracted from the cropped image. We adopt instance normalization [32] after each convolution. For the hyperparameters in Eq.4, we set  $\lambda_1 = 25$ ,  $\lambda_2 = 1$  and  $\lambda_3 = 100$ . For the down-sampling network of the generator, we use the ImageNet [6] pretrained network as initialization. The weights of the rest of the network are initialized from scratch using Gaussian distribution with standard deviation as 0.001. We use Adam optimizer [15] with a mini-batch size of 1. We set  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $weightdecay = 0.0001$ . Learning rate starts from 0.0002. On the LIP dataset, learning rate is divided by 10 after 15 epochs, and the models are trained for 30 epochs. On the Pascal-Person-Part dataset, learning rate is divided by 10 after 25 epochs, and the models are trained for 50 epochs. We use dropout in the deconvolution layers, following the practice in [13]. We alternately optimize the  $D$  and  $G$ . During testing, we average the per-pixel classification scores at multiple scales, *i.e.*, testing images are resized to  $\{0.8, 1, 1.2\}$  times of their original size.

## 4.3 Comparison with the State-of-the-Art Methods

In this section, we compare our result with the state-of-the-art methods on the three datasets. First, on the **LIP dataset**, we compare MMAN with five state-of-the-art methods in Table 1. The proposed MMAN yields an mIoU of 46.65%, while the mIoU of the five competing methods is 18.17% [1], 28.29% [23], 42.92% [3], 44.13% [2] and 44.73% [10], respectively. For a fair comparison, we further implement ASN [24] and SSL [10] on our baseline, *i.e.*, Do-Deeplab-ASPP. On the same baseline, MMAN outperforms ASN [24] and SSL [10] by +1.40% and +0.62% in terms of mIoU, respectively. It clearly indicates that our method outperforms the state of the art. The comparison of per-class IoU indicates that improvement is mainly from classes which are closely related to human pose, such as arms, legs and shoes. In particular, MMAN is capable of distinguishing between “left” and “right”, which gives a huge boost in following human parts: more than +2.5% improvement in left/right arm, more than +10% improvement in left/right leg and more than +5% improvement in left/right shoe. The comparison implies that MMAN is capable of enforcing the consistency of semantic-level features, *i.e.*, human pose.

Table 1: Method comparison of per-class IoU and mIoU on LIP validation set.

Method	hat	hair	glov	sung	clot	dress	coat	sock	pant	suit	scarf	skirt	face	l-arm	r-arm	l-leg	r-leg	l-sh	r-sh	bkg	avg
SegNet[1]	26.60	44.01	0.01	0.00	34.46	0.00	15.97	3.59	33.56	0.01	0.00	0.00	52.38	15.30	24.23	13.82	13.17	9.26	6.47	70.62	18.17
PCN-8s[23]	39.79	58.96	5.32	3.08	49.08	12.36	26.82	15.66	49.41	6.48	0.00	2.16	62.65	29.78	36.63	28.12	26.05	17.76	17.70	78.02	28.29
Attention[3]	58.87	66.78	23.32	19.48	63.20	29.63	49.70	35.23	66.04	24.73	12.84	20.41	70.58	50.17	54.03	38.35	37.70	26.20	27.09	84.00	42.92
DeepLab-ASPP[2]	56.48	65.33	29.98	19.67	62.44	30.33	51.03	40.51	69.00	22.38	11.29	20.56	70.11	49.25	52.88	42.37	35.78	33.81	32.89	84.53	44.03
Attention+SSL[10]	<b>59.75</b>	<b>67.25</b>	28.95	21.57	<b>65.30</b>	29.49	51.92	38.52	68.02	24.48	<b>14.92</b>	<b>24.32</b>	71.01	52.64	55.79	40.23	38.80	28.08	29.03	84.56	44.73
Do-DeepLab-ASPP	56.16	65.28	28.53	20.16	62.54	29.04	51.22	38.00	69.82	22.62	10.63	19.94	69.88	51.83	53.01	45.68	46.08	35.82	34.72	83.47	44.72
Macro AN	57.24	65.28	28.87	19.56	64.02	27.51	51.39	38.13	70.11	22.81	9.05	19.35	68.60	54.19	56.29	50.57	51.22	37.15	37.42	83.25	45.60
Micro AN	57.47	65.05	28.66	16.93	63.95	<b>31.45</b>	51.11	39.64	70.85	25.58	6.87	18.96	68.89	53.62	56.69	49.81	49.42	35.35	35.65	84.46	45.52
ASN [24]	56.92	64.34	28.07	17.78	64.90	30.85	51.90	39.75	<b>71.78</b>	25.57	7.97	17.63	70.77	53.53	56.70	49.58	48.21	34.57	33.31	84.01	45.41
SSL [10]	58.21	67.17	<b>31.20</b>	<b>23.65</b>	63.66	28.31	<b>52.35</b>	39.58	69.40	<b>28.61</b>	13.70	22.52	<b>74.84</b>	52.83	55.67	48.22	47.49	31.80	29.97	84.64	46.19
MMAN	57.66	65.63	30.07	20.02	64.15	28.39	51.98	<b>41.46</b>	71.03	23.61	9.65	23.20	69.54	<b>55.30</b>	<b>58.13</b>	<b>51.90</b>	<b>52.17</b>	<b>38.58</b>	<b>39.05</b>	<b>84.75</b>	<b>46.81</b>

Table 2: Performance comparison in terms of per-class IoU with five state-of-the-art methods on the PASCAL-Person-Part test set.

Method	head	torso	u-arms	l-arms	u-legs	l-legs	bkg	avg
Deeplab-ASPP [2]	81.33	60.06	41.16	40.95	37.49	32.56	92.81	55.19
HAZN [33]	80.79	59.11	43.05	42.76	38.99	34.46	93.59	56.11
Attention [3]	81.47	59.06	44.15	42.50	38.28	35.62	93.65	56.39
LG-LSTM [20]	82.72	60.99	45.40	<b>47.76</b>	42.33	37.96	88.63	57.97
Attention + SSL [10]	<b>83.26</b>	62.40	47.80	45.58	42.32	39.48	94.68	59.36
Do-Deeplab-ASPP	81.82	59.53	44.80	42.79	38.32	36.38	93.91	56.79
Macro AN	82.01	61.19	45.24	44.30	39.73	36.75	93.89	57.58
Micro AN	82.44	61.35	44.79	43.68	38.41	36.05	93.93	57.23
MMAN	82.46	61.41	46.05	45.17	40.93	38.83	94.30	58.45
Attention + MMAN	82.58	<b>62.83</b>	<b>48.49</b>	47.37	<b>42.80</b>	<b>40.40</b>	<b>94.92</b>	<b>59.91</b>

Second, on **PASCAL-Person-Part**, the comparison is shown in Table 2. We apply the same model structure used on the LIP dataset to train the PASCAL-Person-Part dataset. Our model yields an mIoU of 58.45% on the test set. It is higher than most of the compared methods and is only slightly inferior to “Attention+SSL” [10] by 0.91%. This is probably due to the human scale variance in this dataset, which can be addressed by the attention algorithm proposed in [3] and applied in [10].

Therefore, we add a plug-and-play module to our model, *i.e.*, attention network [3]. In particular, we employ multi-scale input and use the attention network to merge the results. The final model “Attention+MMAN” improves mIoU to 59.91%, which is higher than the current state-of-the-art method [10] by +0.55%. When we look into the per-class IoU scores, we have similar observations to the those on LIP. The largest improvement can be observed in arms and legs. The improvement over the state-of-the-art methods [10,20,3] is over +0.6% in upper arms, over +1.8% in lower arms, over +0.4% in upper legs and over +0.9% in lower legs, respectively. The comparisons indicate that our method is very competitive.

Third, we deploy the model trained on LIP to the testing set of the **PPSS dataset** without any fine-tuning. We aim to evaluate the generalization ability of the proposed model.

To make the labels in the LIP and PPSS datasets consistent, we merge the fine-grained labels of LIP into coarse-grained human part labels defined in PPSS.

Table 3: Comparison of human parsing accuracy on the PPSS dataset [25]. Best performance is highlighted in blue.

Method	head	face	up-cloth	arms	lo-cloth	legs	bkg	avg
DL [25]	22.0	29.1	57.3	10.6	46.1	12.9	68.6	35.2
DDN [25]	35.5	44.1	68.4	17.0	<b>61.7</b>	<b>23.8</b>	80.0	47.2
ASN [24]	51.7	<b>51.0</b>	65.9	<b>29.5</b>	52.8	20.3	83.8	50.7
MMAN	<b>53.1</b>	50.2	<b>69.0</b>	29.4	55.9	21.4	<b>85.7</b>	<b>52.1</b>

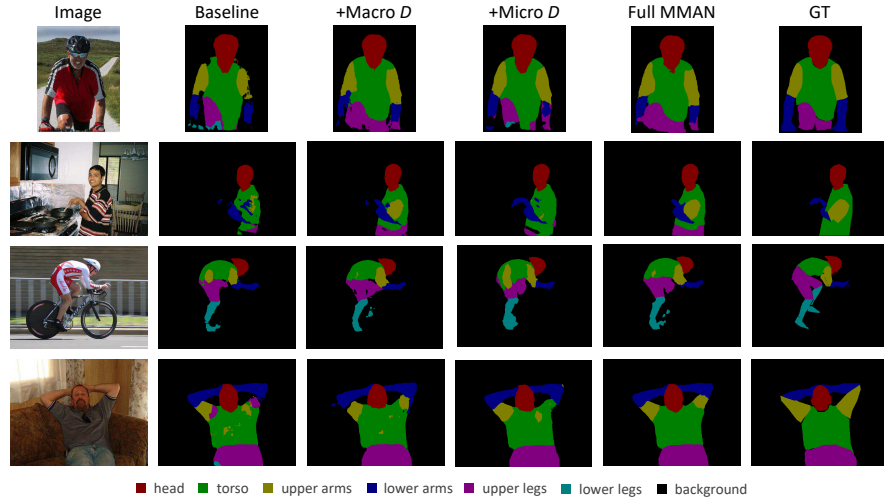


Fig. 5: Qualitative parsing results on the Pascal-Person-Part dataset.

The evaluation result is reported in Table 3. MMAN yields an mIoU of 52.11%, which significantly outperforms DL [25] DDN [25] and ASN [24] by +16.9%, +4.9% and +1.4%, respectively. Therefore, when directly tested on another dataset with different image styles, our model still yields good performance.

In Fig. 5, we provide some segmentation examples obtained by Baseline (Do-Deeplab-ASPP), Baseline+Macro  $D$ , Baseline+Micro  $D$  and full MMAN, respectively. The ground truth label maps are also shown. We observe that Baseline+Micro  $D$  reduces the blur and noise significantly and aids to generate sharp boundaries, and that Baseline+Macro  $D$  corrects the unreasonable human poses. The full MMAN method integrates the advantages of both Macro AN and Micro AN and achieves higher parsing accuracy. We also present qualitative results on the PPSS dataset in Fig. 6.

#### 4.4 Ablation Study

This section presents ablation studies of our method. Since two components are involved, *i.e.*, Macro  $D$  and Micro  $D$ , we remove them one at a time to evaluate their contributions respectively. Results on LIP and PASCAL-Person-Part datasets are shown in Table 1 and Table 2, respectively.



Fig. 6: Qualitative parsing results on the PPSS dataset. RGB image and the label map are showed in pairs.

On the LIP dataset, when removing Macro  $D$  or Micro  $D$  from the system, mIoU will drop 1.21% and 1.29%, respectively, compared with the full MMAN system. Meanwhile, when compared with the baseline approach, employing Macro  $D$  or Micro  $D$  alone brings +0.88% and +0.80% improvement in mean IoU. Similar observations can be made on the PASCAL-Person-Part dataset as well.

To further evaluate the respective function of the two different discriminators, we add two external experiments: 1) For Macro  $D$ , we calculate another mIoU using the low-resolution segmentation maps, which filter out pixel-wise details and retain high-level human structures. So this new mIoU is more suitable for evaluating Macro  $D$ . 2) For Micro  $D$ , we count the “isolated pixels” in high-resolution segmentation maps, which reflects local inconsistency such as “holes”. The “isolated pixel rate” (IPR) can be viewed as a better indicator for evaluating Micro  $D$ . We see from Table 4 that Macro  $D$  is better than Micro  $D$  at improving “mIoU (low-reso.)”, proving that Macro  $D$  *specializes in preserving high-level human structures*. We also see that Micro  $D$  is better than Macro  $D$  at decreasing IPR, suggesting that Micro  $D$  *specializes in improving local consistency* of the result.

#### 4.5 Variant Study


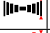


We further evaluate three different variants of MMAN, *i.e.*, Single AN, Double AN, and Multiple AN, on the LIP dataset. Table 5 details the number of parameter, global FOV (g.FOV) and local FOV (l.FOV) sizes, as well as the architecture sketch of each variant. The result of original MMAN is also presented for a clear comparison.

Single AN refers to the traditional adversarial network with only one discriminator. The discriminator is attached to the top layer and has a global receptive field on a  $256 \times 256$  label map. As the result shows, Single AN yields 45.23% in mean IoU, which is slightly higher than the baseline but lower than MMAN. This result suggests that employing Macro  $D$  and Micro  $D$  outperforms the single

Table 4: Comparison in IPR and mIOUs

method	IPR	mIoU (low-reso.)	mIoU (high-reso.)
baseline	5.62	50.66	44.72
+macro $D$	4.23	55.79	45.60
+micro $D$	2.81	53.60	45.52
+CRF	1.53	52.77	45.45
MMAN	2.47	56.95	46.81

Table 5: Variant study of MMAN.

variant	arch.	g.FOV	l.FOV	#par	pc.	mIoU
sAN		$256 \times 256$	-	3.2M	✓	45.23
dAN		$256 \times 256$	$22 \times 22$	3.8M	✓	46.15
mAN		$16 \times 16$	$22 \times 22$	1.8M	-	46.97
MMAN		$16 \times 16$	$22 \times 22$	1.2M	-	46.81

discriminator, which proves the correctness of the analysis in Section 3.5. What is more, we observe the poor convergence (pc) problem when training the Single AN. It is due to the employment of large FOVs on the high-resolution label map.

Double AN has the same number of discriminators with MMAN. The difference lies in that the Double AN attaches the Macro  $D$  to the top layer. Compared to Double AN, MMAN significantly improves the result by 0.82%. The result illustrates the complementary effects of Macro  $D$  and Micro  $D$ : Macro  $D$  acts on deep layers and offers a good coarse-grained initialization for later top layers and Micro  $D$  helps to remedies the coarse semantic feature with fine-grained visual details.

Multiple AN is designed to evaluate the parsing accuracy when employing more than two discriminators. To this end, we attach an extra discriminator to the 3rd deconvolution layer of  $G$ . In particular, the discriminator has the same architecture with micro  $D$  and focuses on  $22 \times 22$  patches on a  $64 \times 64$  label map. As the result shows in Table 5, employing three discriminators brings very slightly improvement (0.16%) in mean IoU, but results in more complex architecture and more parameters.

## 5 Conclusions

In this paper, we introduce a novel Macro-Micro adversarial network (MMAN) for human parsing, which significantly reduces the semantic inconsistency, *e.g.*, misplaced human parts, and the local inconsistency, *e.g.*, blur and holes, in the parsing results. Our model achieves comparative parsing accuracy with the state-of-the-art methods on two challenge human parsing datasets and has a good generalization ability on other datasets. The two adversarial losses are complementary and outperform previous methods that employ a single adversarial loss. Furthermore, MMAN achieves both global and local supervisions with small receptive fields, which effectively avoids the poor convergence problem of adversarial network in handling high-resolution images.

**Acknowledgment.** This work is partially supported by the National Natural Science Foundation of China (No. 61572211). We acknowledge the Data to Decisions CRC (D2D CRC) and the Cooperative Research Centers Programme for funding this research.

## References

1. Badrinarayanan, V., Kendall, A., Cipolla, R.: Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence* **39**(12), 2481–2495 (2017)
2. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *arXiv preprint arXiv:1606.00915* (2016)
3. Chen, L.C., Yang, Y., Wang, J., Xu, W., Yuille, A.L.: Attention to scale: Scale-aware semantic image segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 3640–3649 (2016)
4. Chen, X., Mottaghi, R., Liu, X., Fidler, S., Urtasun, R., Yuille, A.: Detect what you can: Detecting and representing objects using holistic models and body parts. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 1971–1978 (2014)
5. Dai, W., Doyle, J., Liang, X., Zhang, H., Dong, N., Li, Y., Xing, E.P.: Scan: Structure correcting adversarial network for chest x-rays organ segmentation. *arXiv preprint arXiv:1703.08770* (2017)
6. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. pp. 248–255. IEEE (2009)
7. Denton, E.L., Chintala, S., Fergus, R., et al.: Deep generative image models using a laplacian pyramid of adversarial networks. In: *Advances in neural information processing systems*. pp. 1486–1494 (2015)
8. Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The PASCAL Visual Object Classes Challenge 2010 (VOC2010) Results. <http://www.pascal-network.org/challenges/VOC/voc2010/workshop/index.html>
9. Gan, C., Lin, M., Yang, Y., de Melo, G., Hauptmann, A.G.: Concepts not alone: Exploring pairwise relationships for zero-shot video activity recognition. In: *AAAI*. p. 3487 (2016)
10. Gong, K., Liang, X., Shen, X., Lin, L.: Look into person: Self-supervised structure-sensitive learning and a new benchmark for human parsing. *arXiv preprint arXiv:1703.05446* (2017)
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 770–778 (2016)
12. Hung, W.C., Tsai, Y.H., Liou, Y.T., Lin, Y.Y., Yang, M.H.: Adversarial learning for semi-supervised semantic segmentation. *arXiv preprint arXiv:1802.07934* (2018)
13. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. *arXiv preprint* (2017)
14. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. *ICLR* (2018)
15. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
16. Kohli, P., Torr, P.H., et al.: Robust higher order potentials for enforcing label consistency. *International Journal of Computer Vision* **82**(3), 302–324 (2009)
17. Krähenbühl, P., Koltun, V.: Efficient inference in fully connected crfs with gaussian edge potentials. In: *Advances in neural information processing systems*. pp. 109–117 (2011)

18. Li, Q., Arnab, A., Torr, P.H.: Holistic, instance-level human parsing. arXiv preprint arXiv:1709.03612 (2017)
19. Liang, X., Liu, S., Shen, X., Yang, J., Liu, L., Dong, J., Lin, L., Yan, S.: Deep human parsing with active template regression. *IEEE transactions on pattern analysis and machine intelligence* **37**(12), 2402–2414 (2015)
20. Liang, X., Shen, X., Xiang, D., Feng, J., Lin, L., Yan, S.: Semantic object parsing with local-global long short-term memory. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 3185–3193 (2016)
21. Liang, X., Xu, C., Shen, X., Yang, J., Liu, S., Tang, J., Lin, L., Yan, S.: Human parsing with contextualized convolutional neural network. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 1386–1394 (2015)
22. Liu, Z., Li, X., Luo, P., Loy, C.C., Tang, X.: Semantic image segmentation via deep parsing network. In: *Computer Vision (ICCV), 2015 IEEE International Conference on*. pp. 1377–1385. IEEE (2015)
23. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 3431–3440 (2015)
24. Luc, P., Couprie, C., Chintala, S., Verbeek, J.: Semantic segmentation using adversarial networks. arXiv preprint arXiv:1611.08408 (2016)
25. Luo, P., Wang, X., Tang, X.: Pedestrian parsing via deep decompositional network. In: *Computer Vision (ICCV), 2013 IEEE International Conference on*. pp. 2648–2655. IEEE (2013)
26. Luo, Y., Guan, T., Pan, H., Wang, Y., Yu, J.: Accurate localization for mobile device using a multi-planar city model. In: *Pattern Recognition (ICPR), 2016 23rd International Conference on*. pp. 3733–3738. IEEE (2016)
27. Moeskops, P., Veta, M., Lafarge, M.W., Eppenhof, K.A., Pluim, J.P.: Adversarial training and dilated convolutions for brain mri segmentation. In: *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, pp. 56–64. Springer (2017)
28. Odena, A., Olah, C., Shlens, J.: Conditional image synthesis with auxiliary classifier gans. arXiv preprint arXiv:1610.09585 (2016)
29. Park, S., Nie, X., Zhu, S.C.: Attribute and-or grammar for joint parsing of human pose, parts and attributes. *IEEE transactions on pattern analysis and machine intelligence* (2017)
30. Reed, S.E., Akata, Z., Mohan, S., Tenka, S., Schiele, B., Lee, H.: Learning what and where to draw. In: *Advances in Neural Information Processing Systems*. pp. 217–225 (2016)
31. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: *International Conference on Medical image computing and computer-assisted intervention*. pp. 234–241. Springer (2015)
32. Ulyanov, D., Vedaldi, A., Lempitsky, V.S.: Instance normalization: The missing ingredient for fast stylization. *CoRR* **abs/1607.08022** (2016), <http://arxiv.org/abs/1607.08022>
33. Xia, F., Wang, P., Chen, L.C., Yuille, A.L.: Zoom better to see clearer: Human and object parsing with hierarchical auto-zoom net. In: *European Conference on Computer Vision*. pp. 648–663. Springer (2016)
34. Xia, F., Zhu, J., Wang, P., Yuille, A.L.: Pose-guided human parsing by an and/or graph using pose-context features. In: *AAAI*. pp. 3632–3640 (2016)
35. Xue, Y., Xu, T., Zhang, H., Long, R., Huang, X.: Segan: Adversarial network with multi-scale  $L_1$  loss for medical image segmentation. arXiv preprint arXiv:1706.01805 (2017)



36. Zhang, X., Kang, G., Wei, Y., Yang, Y., Huang, T.: Self-produced guidance for weakly-supervised object localization. In: European Conference on Computer Vision. Springer (2018)
37. Zhang, X., Wei, Y., Feng, J., Yang, Y., Huang, T.: Adversarial complementary learning for weakly supervised object localization. In: IEEE CVPR (2018)
38. Zhong, Z., Zheng, L., Li, S., Yang, Y.: Generalizing a person retrieval model hetero- and homogeneously. In: ECCV (2018)
39. Zhong, Z., Zheng, L., Zheng, Z., Li, S., Yang, Y.: Camera style adaptation for person re-identification. In: CVPR (2018)
40. Zhu, S., Fidler, S., Urtasun, R., Lin, D., Loy, C.C.: Be your own prada: Fashion synthesis with structural coherence. In: International Conference on Computer Vision (ICCV) (2017)