

CCNet: Criss-Cross Attention for Semantic Segmentation

Zilong Huang, Xinggang Wang, *Member, IEEE*, Yunchao Wei, Lichao Huang, Humphrey Shi, *Member, IEEE*, Wenyu Liu, *Senior Member, IEEE*, and Thomas Huang, *Life Fellow, IEEE*

Abstract—Contextual information is vital in visual understanding problems, such as semantic segmentation and object detection. We propose a Criss-Cross Network (CCNet) for obtaining full-image contextual information in a very effective and efficient way. Concretely, for each pixel, a novel criss-cross attention module harvests the contextual information of all the pixels on its criss-cross path. By taking a further recurrent operation, each pixel can finally capture the full-image dependencies. Besides, a category consistent loss is proposed to enforce the criss-cross attention module to produce more discriminative features. Overall, CCNet is with the following merits: 1) GPU memory friendly. Compared with the non-local block, the proposed recurrent criss-cross attention module requires $11 \times$ less GPU memory usage. 2) High computational efficiency. The recurrent criss-cross attention significantly reduces FLOPs by about 85% of the non-local block. 3) The state-of-the-art performance. We conduct extensive experiments on semantic segmentation benchmarks including Cityscapes, ADE20K, human parsing benchmark LIP, instance segmentation benchmark COCO, video segmentation benchmark CamVid. In particular, our CCNet achieves the mIoU scores of 81.9%, 45.76% and 55.47% on the Cityscapes test set, the ADE20K validation set and the LIP validation set respectively, which are the new state-of-the-art results. The source codes are available at <https://github.com/speedinghzl/CCNet>.

Index Terms—Semantic Segmentation, Graph Attention, Criss-Cross Network, Context Modeling

1 INTRODUCTION

SEMANTIC segmentation, which is a fundamental problem in the computer vision community, aims at assigning semantic class labels to each pixel in a given image. It has been extensively and actively studied in many recent works and is also critical for various significant applications such as autonomous driving [1], augmented reality [2], and image editing [3]. Specifically, current state-of-the-art semantic segmentation approaches based on the fully convolutional network (FCN) [4] have made remarkable progress. However, due to the fixed geometric structures, the conventional FCN is inherently limited to local receptive fields that only provide short-range contextual information. The limitation of insufficient contextual information imposes a great adverse effect on its segmentation accuracy.

To make up for the above deficiency of FCN, some works have been proposed to introduce useful contextual information to benefit the semantic segmentation task. Specifically, Chen *et al.* [6] proposed atrous spatial pyramid pooling module with multi-scale dilation convolutions for contextual information aggregation. Zhao *et al.* [7] further intro-

- Z. Huang, X. Wang and W. Liu are with the School of Electronic Information and Communications, Huazhong University of Science and Technology, Wuhan 430074, China (e-mail: hzl@hust.edu.cn, xg-wang@hust.edu.cn, liuwy@hust.edu.cn).
 - Y. Wei is with the Centre for Artificial Intelligence, Faculty of Engineering and Information Technology, University of Technology Sydney, Ultimo, NSW 2007, Australia. (e-mail: yunchao.wei@uts.edu.au).
 - L. Huang is with Horizon Robotics. (e-mail: lichao.huang@horizon.ai).
 - H. Shi is with the University of Oregon and the University of Illinois at Urbana-Champaign. (e-mail: shihonghui3@gmail.com).
 - T. S. Huang was with the University of Illinois at Urbana-Champaign. (e-mail: t-huang1@illinois.edu).

Corresponding author: Xinggang Wang. Zilong Huang and Xinggang Wang contributed equally to this work.

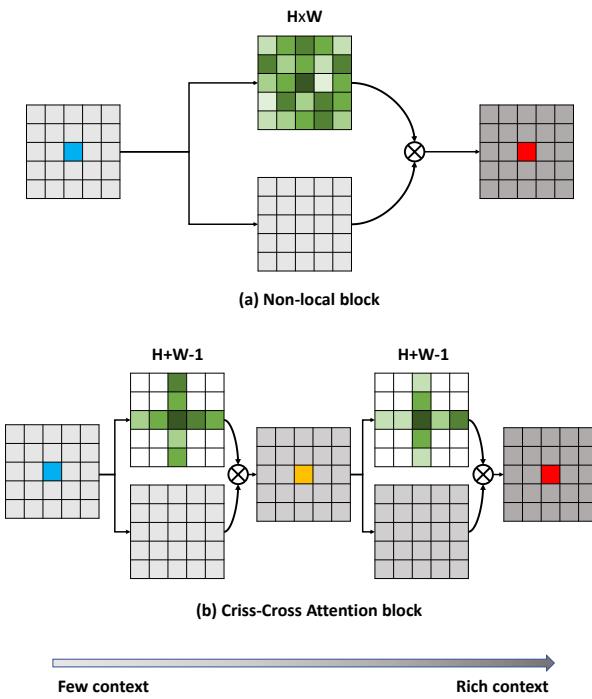


Fig. 1. Diagrams of two attention-based context aggregation methods. (a) For each position (e.g., blue), the Non-local module [5] generates a dense attention map which has N weights (in green). (b) For each position (e.g., blue), the criss-cross attention module generates a sparse attention map which only has about $2\sqrt{N}$ weights. After the recurrent operation, each position (e.g., red) in the final output feature maps can collect information from all pixels. For clear display, residual connections are ignored.

duced PSPNet with pyramid pooling module to capture contextual information. However, the dilated convolution based methods [6], [8], [9] collect information from a few surrounding pixels and cannot generate dense contextual information actually. Meanwhile, the pooling based methods [7], [10] aggregate contextual information in a non-adaptive manner and the homogeneous context extraction procedure is adopted by all image pixels, which does not satisfy the requirement that different pixels need different contextual dependencies.

To incorporate dense and pixel-wise contextual information, some fully-connected graph neural network (GNN) [11] methods were proposed to augment traditional convolutional features with an estimated full-image context representation. PSANet [12] learns to aggregate contextual information for each position via a predicted attention map. Non-local Networks [5] utilizes a self-attention mechanism [13], [14], which enables a single feature from any position to perceive features of all the other positions, thus harvesting full-image contextual information, see Fig. 1 (a). These non-local operations could be viewed as a densely-connected GNN module based on attention mechanism [14]. This feature augmentation method allows a flexible way to represent non-local relations between features and has led to significant improvements in several vision recognition tasks. However, these GNN-based non-local neural networks need to generate huge attention maps to measure the relationships for each pixel-pair, leading to a very high complexity of $\mathcal{O}(N^2)$ for both time and space, where N is the number of input features. Since the dense prediction tasks, such as semantic segmentation, inherently require high resolution feature maps, the non-local based methods will often with high computation complexity and occupy a huge number of GPU memory. Thus, is there an alternative solution to achieve such a target in a more efficient way?

To address the above mentioned issue, our motivation is to replace the common single densely-connected graph with several consecutive sparsely-connected graphs, which usually require much lower computational resources. Without loss of generality, we use two consecutive criss-cross attention modules, in which each one only has sparse connections (about \sqrt{N}) for each position in the feature map. For each pixel/position, the criss-cross attention module aggregates contextual information in its horizontal and vertical directions. By serially stacking two criss-cross attention modules, each position can collect contextual information from all pixels in the given image. The above decomposition strategy will greatly reduce the complexities of both time and space from $\mathcal{O}(N^2)$ to $\mathcal{O}(N\sqrt{N})$.

We compare the differences between the non-local module [5] and our criss-cross attention module in Fig. 1. Concretely, both non-local module and criss-cross attention module feed the input feature map to generate an attention map for each position and transform the input feature map into an adapted feature map. Then, a weighted sum is adopted to collecting contextual information from other positions in the adapted feature map based on the attention maps. Different from the dense connections adopted by the non-local module, each position (e.g., blue) in the feature map is sparsely connected with other ones which are in the same row and the same column in our criss-cross attention

module, leading to the predicted attention map only has about $2\sqrt{N}$ weights rather than N in non-local module.

To achieve the goal of capturing the full-image dependencies, we innovatively and simply take a recurrent operation for the criss-cross attention module. In particular, the local features are firstly passed through one criss-cross attention module to collect the contextual information in horizontal and vertical directions. Then, by feeding the feature map produced by the first criss-cross attention module into the second one, the additional contextual information obtained from the criss-cross path finally enables the full-image dependencies for all positions. As demonstrated in Fig. 1 (b), each position (e.g. red) in the second feature map can collect information from all others to augment the position-wise representations. We share parameters of the criss-cross modules to keep our model slim. Since the input and output are both convolutional feature maps, criss-cross attention module can be easily plugged into any fully convolutional neural network, named as CCNet, for learning full-image contextual information in an end-to-end manner. Thanks to the good usability of criss-cross attention module, CCNet is straight forward to extend to 3D networks for capturing long-range temporal context information.

In addition, to drive the proposed recurrent criss-cross attention method to learn more discriminative features, we introduce a category consistent loss to augment CCNet. Particularly, the category consistent loss enforces the network to map each pixel in the image to an n-dimensional vector in the feature space, such that feature vectors of pixels that belong to the same category lie close together while feature vectors of pixels that belong to different categories lie far apart.

We have carried out extensive experiments on multiple large-scale datasets. Our proposed CCNet achieves top performance on four most competitive semantic segmentation datasets, i.e., Cityscapes [15], ADE20K [16], LIP [17] and CamVid [18]. In addition, the proposed criss-cross attention even improves the state-of-the-art instance segmentation method, i.e., Mask R-CNN with ResNet-101 [19]. These results well demonstrate that our criss-cross attention module is generally beneficial to the dense prediction tasks. In summary, our main contributions are three-fold:

- We propose a novel criss-cross attention module in this work, which can be leveraged to capture contextual information from full-image dependencies in a more efficient and effective way.
- We propose category consistent loss which can enforce criss-cross attention module to produce more discriminative features.
- We propose CCNet by taking advantages of recurrent criss-cross attention module, achieving leading performance on segmentation-based benchmarks, including Cityscapes, ADE20K, LIP, CamVid and COCO.

Compare with our original conference version [20], the following improvements are conducted: 1) We further enhance the segmentation ability of CCNet by augmenting a simple yet effective category consistent loss; 2) we propose a more generic CCNet by extending the criss-cross attention module from 2D to 3D; 3) we include more extensive exper-

iments on the LIP, CamVid and COCO datasets to verify the effectiveness and generalization ability of our CCNet.

The rest of this paper is organized as follows. We first review related work in Section 2 and describe the architecture of our network in Section 3. In Section 4, ablation studies are given and experimental results are analyzed. Section 5 presents our conclusion and future work.

2 RELATED WORK

2.1 Semantic segmentation

The last years have seen a renewal of interest on semantic segmentation. FCN [4] is the first approach to adopt fully convolutional network for semantic segmentation. Later, FCN-based methods have made remarkable progress in image semantic segmentation. Chen *et al.* [21] and Yu *et al.* [22] removed the last two downsample layers to obtain dense prediction and utilized dilated convolutions to enlarge the receptive field. Unet [23], DeepLabv3+ [24], MSCI [25], SPGNet [26], RefineNet [27] and DFN [28] adopted encoder-decoder structures that fuse the information in low-level and high-level layers to make dense predictions. The scale-adaptive convolutions (SAC) [29] and deformable convolutional networks (DCN) [30] methods improved the standard convolutional operator to handle the deformation and various scales of objects. CRF-RNN [22] and DPN [31] used Graph model, *i.e.*, CRF, MRF, for semantic segmentation. AAF [32] used adversarial learning to capture and match the semantic relations between neighboring pixels in the label space. BiSeNet [33] was designed for real-time semantic segmentation. DenseDecoder [34] built feature-level long-range skip connections on cascaded architecture. VideoGCRF [35] used a densely-connected spatio-temporal graph for video semantic segmentation. RTA [36] proposed the region-based temporal aggregation for leveraging the temporal information in videos. In addition, some works focus on human parsing task. JPPNet [17] embed pose estimation into human parsing task. CE2P [37] proposed a simple yet effective framework for computing context embedding while preserving edges. SANet [38] used parallel branches with scale attention to handle large scale variance in human parsing.

2.2 Contextual information aggregation

It is a common practice to aggregate contextual information to augment the feature representation in semantic segmentation networks. Deeplabv2 [6] proposed atrous spatial pyramid pooling (ASPP) to use different dilation convolutions to capture contextual information. DenseASPP [39] brought dense connections into ASPP to generate features with various scale. DPC [40] utilized architecture search techniques to build multi-scale architectures for semantic segmentation. Chen *et al.* [41] made use of several attention masks to fuse feature maps or prediction maps from different branches. PSPNet [7] utilized pyramid spatial pooling to aggregate contextual information. Recently, Zhao *et al.* [12] proposed the point-wise spatial attention network which uses predicted attention map to guide contextual information collection. Auto-Deeplab [42] utilized neural architecture search to search an effective context modeling. He *et al.* [43] proposed an adaptive pyramid context module

for semantic segmentation. Liu *et al.* [44] utilized recurrent neural networks (RNNs) to capture long-range dependencies.

There are some works use graph models to model the contextual information. Conditional random field (CRF) [21], [36], [45], Markov random field (MRF) [31] were also utilized to capture long-range dependencies for semantic segmentation. Vaswani *et al.* [14] applied a self-attention model on machine translation. Wang *et al.* [5] proposed the non-local module to generate the huge attention map by calculating the correlation matrix between each spatial point on the feature maps, then the attention map guided dense contextual information aggregation. OCNet [46] and DANet [47] utilized Non-local module [5] to harvest the contextual information. PSA [12] learned an attention map to aggregate contextual information for each individual point adaptively and specifically. Chen *et al.* [48] proposed graph-based global reasoning networks which implements relation reasoning via graph convolution on a small graph.

CCNet vs. Non-Local vs. GCN. Here, we specifically discuss the differences among GCN [49], Non-local Network [5] and CCNet. In term of contextual information aggregation, only the center point can perceive the contextual information from all pixels by the global convolution filters in GCN [49]. In contrast, Non-local Network [5] and CCNet guarantee that a pixel at any position perceives contextual information from all pixels. Though GCN [49] alternatively decomposes the square-shape convolutional operation to horizontal and vertical linear convolutional operations which is related to CCNet, CCNet takes the criss-cross way to harvest contextual information which is more effective than the horizontal-vertical separate way. Moreover, CCNet is proposed to mimic Non-local Network [5] for obtaining dense contextual information through a more effective and efficient recurrent criss-cross attention module, in which dissimilar features get low attention weights and features with high attention weights are similar ones. GCN [49] is a conventional convolution neural network, while CCNet is a graph neural network in which each pixel in the convolutional feature map is considered as a node and the relation/context among nodes can be utilized to generate better node features.

2.3 Graph neural networks

Our work is related to deep graph neural network (GNN). Prior to graph neural networks, graphical models, such as the conditional random field (CRF) [21], [36], [45], markov random field (MRF) [31], were widely used to model the long-range dependencies for image understanding. GNNs were early studied in [11], [50], [51]. Inspired by the success of CNNs, a large number of methods adapt graph structure into CNNs. These methods could be divided into two main streams, the spectral-based approaches [52], [53], [54], [55] and the spatial-based approaches [5], [56], [57], [58]. The proposed CCNet belongs to the latter.

3 APPROACH

In this section, we give the details of the proposed Criss-Cross Network (CCNet) for semantic segmentation. We

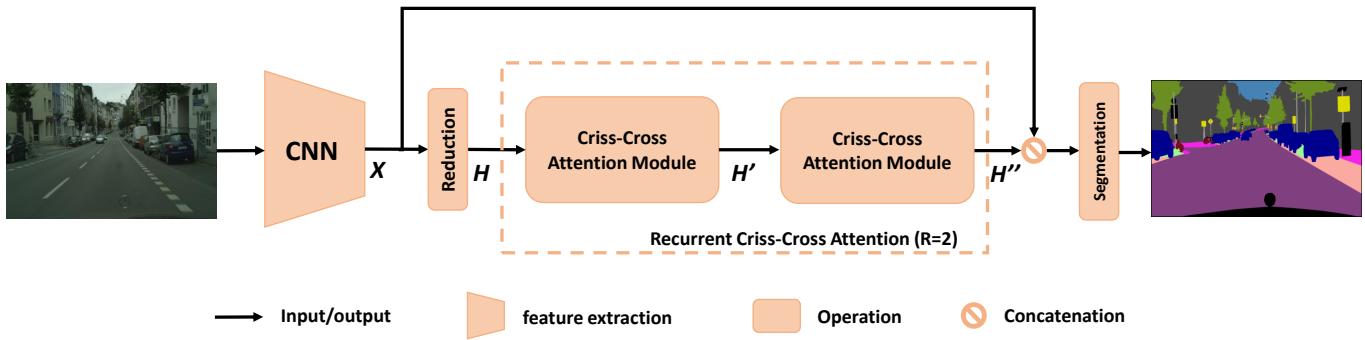


Fig. 2. Overview of the proposed CCNet for semantic segmentation.

first present a general framework of our CCNet. Then, the 2D criss-cross attention module which captures contextual information in horizontal and vertical directions will be introduced. To capture the dense and global contextual information, we propose to adopt a recurrent operation for the criss-cross attention module. To further improve RCCA, we introduce a discriminative loss function to drive RCCA to learn category consistent features. Finally we propose the 3D criss-cross attention module for leveraging temporal and spatial contextual information simultaneously.

3.1 Network Architecture

The network architecture is given in Fig. 2. An input image is passed through a deep convolutional neural network (DCNN), which is designed in a fully convolutional fashion [6], to produce feature map X with the spatial size of $H \times W$. In order to retain more details and efficiently produce dense feature maps, we remove the last two down-sampling operations and employ dilation convolutions in the subsequent convolutional layers, leading to enlarging the width/height of the output feature map X to 1/8 of the input image.

Given X , we first apply a convolutional layer to obtain the feature map H of dimension reduction. Then, H is fed into the criss-cross attention module to generate a new feature map H' which aggregate contextual information together for each pixel in its criss-cross path. The feature map H' only contains the contextual information in horizontal and vertical directions which are not powerful enough for accurate semantic segmentation. To obtain richer and denser context information, we feed the feature map H' into the criss-cross attention module again and output the feature map H'' . Thus, each position in H'' actually gathers the information from all pixels. Two criss-cross attention modules before and after share the same parameters to avoid adding too many extra parameters. We name this recurrent structure as recurrent criss-cross attention (RCCA) module.

Then, we concatenate the dense contextual feature H'' with the local representation feature X . It is followed by one or several convolutional layers with batch normalization and activation for feature fusion. Finally, the fused features are fed into the segmentation layer to predict the final segmentation result.

3.2 Criss-Cross Attention

To model full-image dependencies over local feature representations using light-weight computation and memory, we introduce a criss-cross attention module. The criss-cross attention module collects contextual information in horizontal and vertical directions to enhance pixel-wise representative capability. As shown in Fig. 3, given a local feature map $\mathbf{H} \in \mathbb{R}^{C \times W \times H}$, the module first applies two convolutional layers with 1×1 filters on \mathbf{H} to generate two feature maps \mathbf{Q} and \mathbf{K} , respectively, where $\{\mathbf{Q}, \mathbf{K}\} \in \mathbb{R}^{C' \times W \times H}$. C' is the number of channel, which is less than C for dimension reduction.

After obtaining \mathbf{Q} and \mathbf{K} , we further generate an attention map $\mathbf{A} \in \mathbb{R}^{(H+W-1) \times (W \times H)}$ via **Affinity** operation. At each position \mathbf{u} in the spatial dimension of \mathbf{Q} , we can obtain a vector $\mathbf{Q}_{\mathbf{u}} \in \mathbb{R}^{C'}$. Meanwhile, we can also obtain the set $\Omega_{\mathbf{u}} \in \mathbb{R}^{(H+W-1) \times C'}$ by extracting feature vectors from \mathbf{K} which are in the same row or column with position \mathbf{u} . $\Omega_{i,\mathbf{u}} \in \mathbb{R}^{C'}$ is the i -th element of $\Omega_{\mathbf{u}}$. The **Affinity** operation is then defined as follows.

$$d_{i,\mathbf{u}} = \mathbf{Q}_{\mathbf{u}} \Omega_{i,\mathbf{u}}^T, \quad (1)$$

where $d_{i,\mathbf{u}} \in \mathbf{D}$ is the degree of correlation between features $\mathbf{Q}_{\mathbf{u}}$ and $\Omega_{i,\mathbf{u}}$, $i = [1, \dots, H + W - 1]$, and $\mathbf{D} \in \mathbb{R}^{(H+W-1) \times (W \times H)}$. Then, we apply a softmax layer on \mathbf{D} over the channel dimension to calculate the attention map \mathbf{A} .

Another convolutional layer with 1×1 filters is applied on \mathbf{H} to generate $\mathbf{V} \in \mathbb{R}^{C \times W \times H}$ for feature adaptation. At each position \mathbf{u} in the spatial dimension of \mathbf{V} , we can obtain a vector $\mathbf{V}_{\mathbf{u}} \in \mathbb{R}^C$ and a set $\Phi_{\mathbf{u}} \in \mathbb{R}^{(H+W-1) \times C}$. The set $\Phi_{\mathbf{u}}$ is a collection of feature vectors in \mathbf{V} which are in the same row or column with position \mathbf{u} . The contextual information is collected by an **Aggregation** operation defined as follows.

$$\mathbf{H}'_{\mathbf{u}} = \sum_{i=0}^{H+W-1} \mathbf{A}_{i,\mathbf{u}} \Phi_{i,\mathbf{u}} + \mathbf{H}_{\mathbf{u}}, \quad (2)$$

where $\mathbf{H}'_{\mathbf{u}}$ is a feature vector in $\mathbf{H}' \in \mathbb{R}^{C \times W \times H}$ at position \mathbf{u} and $\mathbf{A}_{i,\mathbf{u}}$ is a scalar value at channel i and position \mathbf{u} in \mathbf{A} . The contextual information is added to local feature \mathbf{H} to augment the pixel-wise representation. Therefore, it has a wide contextual view and selectively aggregates contexts

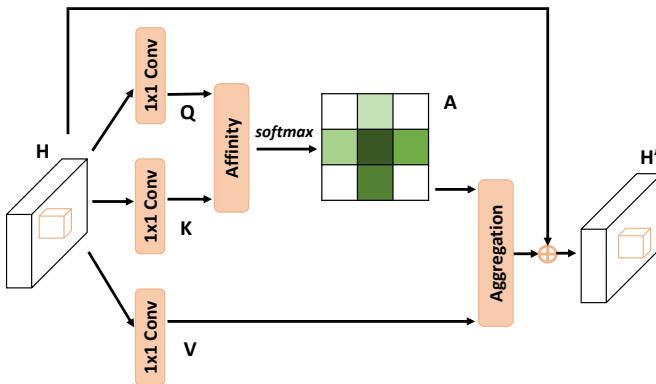


Fig. 3. The details of criss-cross attention module.

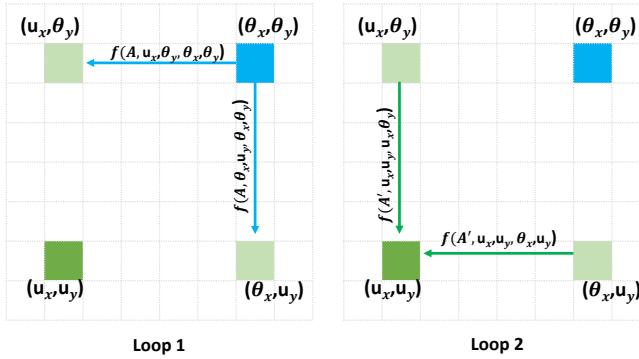


Fig. 4. An example of information propagation when the loop number is 2.

according to the spatial attention map. These feature representations achieve mutual gains and are more robust for semantic segmentation.

3.3 Recurrent Criss-Cross Attention (RCCA)

Despite the criss-cross attention module can capture contextual information in horizontal and vertical directions, the connections between one pixel and its around ones that are not in the criss-cross path are still absent. To tackle this problem, we innovatively and simply introduce a RCCA operation based on the criss-cross attention. The RCCA module can be unrolled into R loops. In the first loop, the criss-cross attention takes the feature map \mathbf{H} extracted from a CNN model as the input and output the feature map \mathbf{H}' , where \mathbf{H} and \mathbf{H}' are with the same shape. In the second loop, the criss-cross attention takes the feature map \mathbf{H}' as the input and output the feature map \mathbf{H}'' . As shown in Fig. 2, the RCCA module is equipped with two loops ($R = 2$) which is able to harvest full-image contextual information from all pixels to generate new features with dense and rich contextual information.

We denote \mathbf{A} and \mathbf{A}' as the attention maps in loop 1 and loop 2, respectively. Since we are interested only in contextual information spreads in spatial dimension rather than in channel dimension, the convolutional layer with 1×1 filters can be viewed as the identical connection. In the case of $R = 2$, the connections between any two spatial positions in the feature map built up by the RCCA module

can be clearly and quantitatively described by introducing function f defined as follows.

$$\exists i \in \mathbb{R}^{H+W-1}, s.t. \mathbf{A}_{i,\mathbf{u}} = f(\mathbf{A}, u_x^{CC}, u_y^{CC}, u_x, u_y),$$

where $\mathbf{u}(u_x, u_y) \in \mathbb{R}^{H \times W}$ is any spatial position in \mathbf{H} and $\mathbf{u}^{CC}(u_x^{CC}, u_y^{CC}) \in \mathbb{R}^{H+W-1}$ is a position in the criss-cross structure centered at \mathbf{u} . The function f is actually an **one-to-one mapping** from the position pair $(\mathbf{u}^{CC}, \mathbf{u}) \in \mathbb{R}^{(H+W-1) \times (H \times W)}$ in the feature map to a particular element $\mathbf{A}_{i,\mathbf{u}} \in \mathbb{R}^{(H+W-1) \times (H \times W)}$ in the attention map $\mathbf{A} \subset \mathbb{R}^{(H+W-1) \times (H \times W)}$, where \mathbf{u}^{CC} maps to a particular row i in \mathbf{A} and \mathbf{u} maps to a particular column in \mathbf{A} .

With the help of function f , we can easily describe the information propagation between any position \mathbf{u} in \mathbf{H}'' and any position θ in \mathbf{H} . It is obvious that information could flow from θ to \mathbf{u} when θ is in the criss-cross path of \mathbf{u} .

Then, we focus on another situation in which $\theta(\theta_x, \theta_y)$ is NOT in the criss-cross path of $\mathbf{u}(u_x, u_y)$. To make it easier to understand, we visualize the information propagation in Fig. 4. The position (θ_x, θ_y) , which is blue, firstly passes the information into the (u_x, θ_y) and (θ_x, u_y) (light green) in the loop 1. The propagation could be quantified by function f . It should be noted that these two points (u_x, θ_y) and (θ_x, u_y) are in the criss-cross path of $\mathbf{u}(u_x, u_y)$. Then, the positions (u_x, θ_y) and (θ_x, u_y) pass the information into the (u_x, u_y) (dark green) in the loop 2. Thus, the information in $\theta(\theta_x, \theta_y)$ could eventually flow into $\mathbf{u}(u_x, u_y)$ even if $\theta(\theta_x, \theta_y)$ is NOT in the criss-cross path of $\mathbf{u}(u_x, u_y)$.

In general, our RCCA module makes up for the deficiency of criss-cross attention that cannot obtain the dense contextual information from all pixels. Compared with criss-cross attention, the RCCA module ($R = 2$) does not bring extra parameters and can achieve better performance with the cost of a minor computation increment.

3.4 Learning Category Consistent Features

For semantic segmentation tasks, the pixels belonging to the same category should have the similar features, while the pixels from different categories should have far apart features. We name such a characteristic as category consistency. The deep features produced by RCCA have full-image context; however, the aggregated feature may have the problem of over-smoothing, which is a common issue in graph neural networks. To address this potential issue, beside the cross-entropy loss ℓ_{seg} to penalize the mismatch between the final predicted segmentation maps and ground truth, we further introduce the category consistent loss to drive RCCA module to learn category consistent features directly.

In [59], a discriminative loss function with three competing terms is proposed for instance segmentation. In particular, the three terms, denoted as $\ell_{var}, \ell_{dis}, \ell_{reg}$, are adopted to 1) penalize large distances between features with the same label for each instance, 2) penalize small distances between the mean features of different labels, and 3) draw mean features of all categories towards the origin, respectively.

Motivated by [59], we first adapt a discriminative loss for semantic segmentation rather than instance segmentation, then replace the first term with more robust one: instead

of using quadratic function as the distance function to penalize mismatch all along, we design a piece-wise distance function to make the optimization more robust.

Let C be the set of classes that are present in the mini-batch images. N_c is the number of valid elements belonging to category $c \in C$. $h_i \in \mathbf{H}$ is the feature vector at spatial position i . μ_c is the mean feature of category $c \in C$ (the cluster center). φ is a piece-wise distance function. δ_v and δ_d are respectively the margins. In particular, Eq. 6 is a piece-wise distance function and the function φ_{var} will be zero, quadratic, and linear function when the distance from the center μ_c is within d_v , in range of $(\delta_v, \delta_d]$, and exceeds δ_d , respectively.

$$\ell_{var} = \frac{1}{|C|} \sum_{c \in C} \frac{1}{N_c} \sum_{i=1}^{N_c} \varphi_{var}(h_i, \mu_c), \quad (3)$$

$$\ell_{dis} = \frac{1}{|C|(|C|-1)} \sum_{c_a \in C} \sum_{c_b \in C, c_a \neq c_b} \varphi_{dis}(\mu_{c_a}, \mu_{c_b}), \quad (4)$$

$$\ell_{reg} = \frac{1}{|C|} \sum_{c \in C} \|\mu_c\|, \quad (5)$$

$$\varphi_{var} = \begin{cases} \|\mu_c - h_i\| - \delta_d + (\delta_d - \delta_v)^2, & \|\mu_c - h_i\| > \delta_d \\ ((\|\mu_c - h_i\| - \delta_v)^2, & \delta_v < \|\mu_c - h_i\| \leq \delta_d \\ 0, & \|\mu_c - h_i\| \leq \delta_v \end{cases} \quad (6)$$

$$\varphi_{dis} = \begin{cases} (2\delta_d - \|\mu_{c_a} - \mu_{c_b}\|)^2, & \|\mu_{c_a} - \mu_{c_b}\| \leq 2\delta_d \\ 0, & \|\mu_{c_a} - \mu_{c_b}\| > 2\delta_d \end{cases} \quad (7)$$

To reduce the computation load, we first apply a convolutional layer with 1×1 filters on the output of RCCA module for dimension reduction and then apply these three loss on the feature map with fewer channels. The final loss ℓ is weighted sum of all losses.

$$\ell = \ell_{seg} + \alpha \ell_{var} + \beta \ell_{dis} + \gamma \ell_{reg}, \quad (8)$$

where α , β and ℓ are the weight parameters. In our experiments we set $\delta_v = 0.5$, $\delta_d = 1.5$, $\alpha = \beta = 1$, $\gamma = 0.001$ and 16 as the number of channels for dimension reduction.

3.5 3D Criss-Cross Attention

To adapt our method from 2D applications to 3D dense prediction tasks, we introduce 3D Criss-Cross Attention. In general, the architecture of 3D Criss-Cross Attention is an extension the 2D version by additional collecting more contextual information from the temporal dimension. As shown in Fig. 5, given a local feature map $\mathbf{H} \in \mathbb{R}^{C \times T \times W \times H}$, where T is axial dimension (*i.e.*, temporal dimension in video data). The module firstly applies two convolutional layers with $1 \times 1 \times 1$ filters on \mathbf{H} to generate two feature maps \mathbf{Q} and \mathbf{K} , respectively, where $\{\mathbf{Q}, \mathbf{K}\} \in \mathbb{R}^{C' \times T \times W \times H}$.

After obtaining the feature maps \mathbf{Q} and \mathbf{K} , we further generate an attention map $\mathbf{A} \in \mathbb{R}^{(T+H+W-2) \times T \times W \times H}$ via the **Affinity** operation. At each position u of \mathbf{Q} , we can obtain a vector $\mathbf{Q}_u \in \mathbb{R}^{C'}$. u contains three coordinate values (t, x, y) . We can also obtain the set $\Omega_u \in \mathbb{R}^{(T+H+W-2) \times C'}$

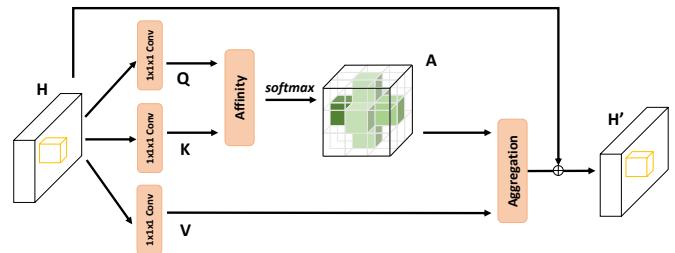


Fig. 5. The details of 3D criss-cross attention module.

by extracting feature vectors from \mathbf{K} with at least two coordinate values equal to u . $\Omega_{i,u} \in \mathbb{R}^{C'}$ is the i -th element of Ω_u . The **Affinity** operation is then defined as follows.

$$d_{i,u} = \mathbf{Q}_u \Omega_{i,u}^\top, \quad (9)$$

where $d_{i,u} \in \mathbf{D}$ is the degree of correlation between feature \mathbf{Q}_u and $\Omega_{i,u}$, $i = [1, \dots, (T+H+W-2)]$, $\mathbf{D} \in \mathbb{R}^{(T+H+W-2) \times T \times W \times H}$. Then, we apply a softmax layer on \mathbf{D} over the first dimension to calculate the attention map \mathbf{A} .

Another convolutional layer with $1 \times 1 \times 1$ filters is applied on \mathbf{H} to generate $\mathbf{V} \in \mathbb{R}^{C \times T \times W \times H}$ for feature adaptation. At each position u in the spatial dimension of \mathbf{V} , we can obtain a vector $\mathbf{V}_u \in \mathbb{R}^C$ and a set $\Phi_u \in \mathbb{R}^{(T+H+W-2) \times C}$. The set Φ_u is a collection of feature vectors in \mathbf{V} which are in the criss-cross structure centered at u . The contextual information is collected by the **Aggregation** operation:

$$\mathbf{H}'_u = \sum_{i=0}^{T+H+W-2} \mathbf{A}_{i,u} \Phi_{i,u} + \mathbf{H}_u, \quad (10)$$

where \mathbf{H}'_u is a feature vector in the output feature map $\mathbf{H}' \in \mathbb{R}^{C \times T \times W \times H}$ at position u . $\mathbf{A}_{i,u}$ is a scalar value at channel i and position u in \mathbf{A} .

4 EXPERIMENTS

To evaluate the effectiveness of the CCNet, we carry out comprehensive experiments on the Cityscapes dataset [15], the ADE20K dataset [16], the COCO dataset [60], the LIP dataset [17] and the CamVid dataset [61]. Experimental results demonstrate that CCNet achieves state-of-the-art performance on Cityscapes, ADE20K and LIP. Meanwhile, CCNet can bring constant performance gain on COCO for instance segmentation. In the following subsections, we first introduce the datasets and implementation details, then we perform a series of ablation experiments on Cityscapes dataset. Finally, we report our results on ADE20K, LIP, COCO and CamVid datasets.

4.1 Datasets and Evaluation Metrics

We adopt Mean IoU (mIoU, mean of class-wise intersection over union) for Cityscapes, ADE20K, LIP and CamVid and the standard COCO metrics Average Precision (AP) for COCO.

- **Cityscapes** is tasked for urban segmentation. Only the 5,000 finely annotated images are used in our experiments and are divided into 2,975/500/1,525 images for training, validation, and testing, respectively.
- **ADE20K** is a recent scene parsing benchmark containing dense labels of 150 stuff/object categories. The dataset includes 20k/2k/3k images for training, validation and testing, respectively.
- **LIP** is a large-scale single human parsing dataset. There are 50,462 images with fine-grained annotations at pixel-level with 19 semantic human part labels and one background label. Those images are further divided into 30k/10k/10k for training, validation and testing, respectively.
- **COCO** is a very challenging dataset for instance segmentation that contains 115k images over 80 categories for training, 5k images for validation and 20k images for testing.
- **CamVid** is one of the datasets focusing on semantic segmentation for autonomous driving scenarios. It is composed of 701 densely annotated images with size 720×960 from five video sequences.

4.2 Implementation Details

Network Structure For semantic segmentation, we choose the ImageNet pre-trained ResNet-101 [19] as our backbone network, remove its last two down-sampling operations, and employ dilated convolutions in the subsequent convolutional layers following the previous work [21], resulting in the output stride as 8. For human parsing, we choose CE2P [37] as our baseline and replace the Context Embedding module with RCCA. For instance segmentation, we choose Mask-RCNN [62] as our baseline. For video semantic segmentation, we also choose Cityscapes pre-trained ResNet-101 [19] as our backbone network with 3D RCCA.

Training settings SGD with mini-batch is used for training. For semantic segmentation, the initial learning rate is $1e-2$ for Cityscapes and ADE20K. Following the prior works [6], [10], we employ a poly learning rate policy where the initial learning rate is multiplied by $1 - (\frac{\text{iter}}{\max_{\text{iter}}})^{\text{power}}$ with $\text{power} = 0.9$. We use the momentum of 0.9 and a weight decay of 0.0001. For Cityscapes, the training images are augmented by randomly scaling (from 0.75 to 2.0), then randomly cropping out high-resolution patches (769×769) from the resulting images. Since the images from ADE20K are with various sizes, we adopt an augmentation strategy of resizing the short side of input image to a length randomly chosen from the set {300, 375, 450, 525, 600}. For human parsing, the model are trained and tested with the input size of 473×473 . For instance segmentation, we take the same training settings as that of Mask-RCNN [62]. For video semantic segmentation, we sample 5 temporally ordered frames from a training video as training data and the input size is 504×504 .

4.3 Experiments on Cityscapes

4.3.1 Comparisons with state-of-the-arts

Results of other state-of-the-art semantic segmentation solutions on Cityscapes are summarized in Tab. 1. For val

TABLE 1
Comparison with state-of-the-arts on Cityscapes (test).

| Method | Backbone | mIOU(%) |
|--------------------------------|----------------|-------------|
| <i>Performance on val set</i> | | |
| DeepLabv3 [8] | ResNet-101 | 79.3 |
| DeepLabv3+ [24] | Xception-65 | 79.1 |
| DPC [40] † | Xception-71 | 80.8 |
| CCNet | ResNet-101 | 80.5 |
| <i>Performance on test set</i> | | |
| DeepLab-v2 [6] | ResNet-101 | 70.4 |
| RefineNet [27] ‡ | ResNet-101 | 73.6 |
| SAC [29] ‡ | ResNet-101 | 78.1 |
| GCN [49] ‡ | ResNet-101 | 76.9 |
| DUC [63] ‡ | ResNet-101 | 77.6 |
| ResNet-38 [64] | WiderResnet-38 | 78.4 |
| PSPNet [7] | ResNet-101 | 78.4 |
| BiSeNet [33] ‡ | ResNet-101 | 78.9 |
| AAF [32] | ResNet-101 | 79.1 |
| PSANet [12] ‡ | ResNet-101 | 80.1 |
| DFN [28] ‡ | ResNet-101 | 79.3 |
| DenseASPP [39] ‡ | DenseNet-161 | 80.6 |
| CCNet ‡ | ResNet-101 | 81.9 |

† use extra COCO dataset for training.

‡ train with both the train-fine and val-fine datasets.

set, we provide these results for reference and emphasize that these results should not be simply compared with our method, since these methods are trained on different (even larger) training sets or different basic network. Among these approaches, DeepLabv3 [8] adopts multi-scale testing strategy. DeepLabv3+ [24] and DPC [40] both use a more stronger backbone (*i.e.*, Xception-65 & 71 vs. ResNet-101). In addition, DPC [40] makes use of additional dataset, *i.e.*, COCO, for pre-training beyond the training set of Cityscapes. The results show that the proposed CCNet with single-scale testing still achieve comparable performance without bells and whistles.

Additionally, we also train the best learned CCNet with ResNet-101 as the backbone using both training and validation sets and make the evaluation on the test set by submitting our test results to the official evaluation server. Most of methods [6], [7], [12], [27], [28], [29], [32], [33], [49], [63] adopt the same backbone as ours and the others [39], [64] utilize stronger backbones. From Tab. 1, it can be observed that our CCNet substantially outperforms all the previous state-of-the-arts on test set. Among the approaches, PSANet [12] is the most related to our method which generates sub attention map for each pixel. One of the differences is that the sub attention map has $2 \times H \times W$ weights in PSANet and $H + W - 1$ weights in CCNet. Even with lower computation cost and memory usage, our method still achieves better performance.

4.3.2 Ablation studies

To verify the rationality of the CCNet, we conduct extensive ablation experiments on the validation set of Cityscapes with different settings for CCNet.

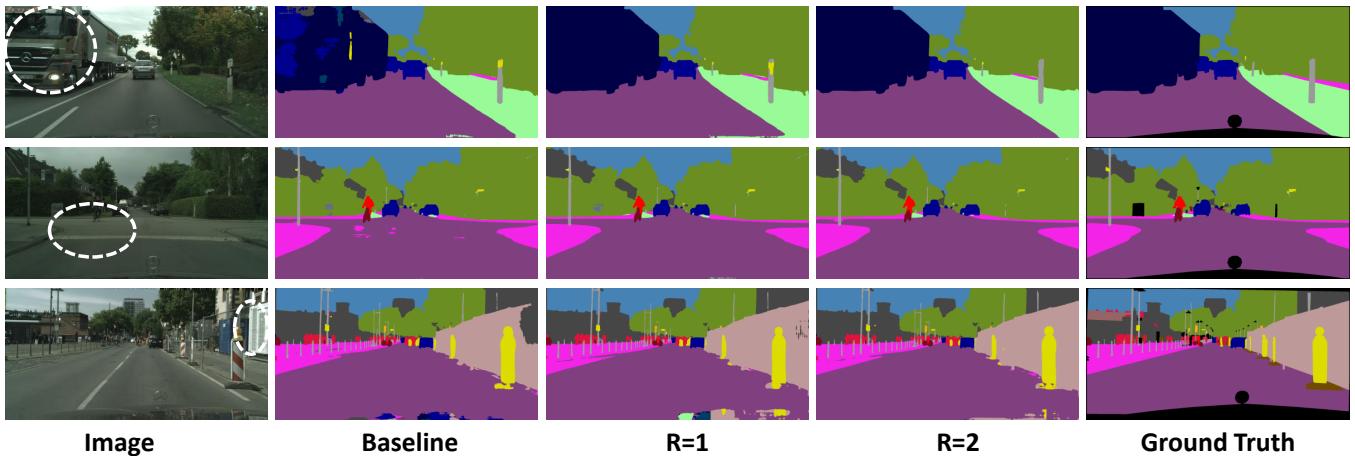


Fig. 6. Visualization results of RCCA with different loops on Cityscapes validation set.

TABLE 2

Performance on Cityscapes (val) for different number of loops in RCCA. FLOPs and memory increment are estimated for an input of $1 \times 3 \times 769 \times 769$.

| Loops | GFLOPs(Δ) | Memory($M\Delta$) | mIOU(%) |
|----------|--------------------|---------------------|---------|
| baseline | 0 | 0 | 75.1 |
| $R = 1$ | 8.3 | 53 | 78.0 |
| $R = 2$ | 16.5 | 127 | 79.8 |
| $R = 3$ | 24.7 | 208 | 80.2 |

The effect of the RCCA module Tab. 2 shows the performance on the Cityscapes validation set by adopting different number of loop in RCCA. All experiments are conducted using ResNet-101 as the backbone. Besides, the input size of training images is 769×769 and the size of the input feature map H of RCCA is 97×97 . Our baseline network is the ResNet-based FCN with dilated convolutional module incorporated at stage 4 and 5, *i.e.*, dilation rates are set to 2 and 4 for these two stages respectively. The increment of FLOPs and memory usage are estimated when $R = 1, 2, 3$, respectively.

We observe that adding a criss-cross attention module into the baseline, denoted as $R = 1$, improves the performance by 2.9%, which can effectively demonstrates the significance of criss-cross attention. Furthermore, increasing the number of loops from 1 to 2 can further improve the performance by 1.8%, demonstrating the effectiveness of dense contextual information. Finally, increasing loops from 2 to 3 slightly improves the performance by 0.4%. Meanwhile, with the increasing the number of loops, the FLOPs and usage of GPU memory keep increasing. These results prove that the proposed criss-cross attention can significantly improve the performance by capturing contextual information in horizontal and vertical direction. In addition, the proposed RCCA is effective in capturing the dense and global contextual information, which can finally benefit the performance of semantic segmentation. To balance the performance and resource usage, we choose $R = 2$ as default settings in all the following experiments.

TABLE 3

Performance on Cityscapes (val) for different kinds of category consistent loss.

| Function Type | Successes | Mean mIOU(%) |
|---------------------|-----------|--------------|
| Quadratic function | 6/10 | 79.2 |
| Piece-wise function | 9/10 | 79.3 |

To further validate the effectiveness of the criss-cross module, we provide the qualitative comparisons in Fig. 6. We leverage the *white circles* to indicate those challenging regions that are easily to be misclassified. It can be seen that these challenging regions are progressively corrected with the increasing the number of loops, which can well prove the effectiveness of dense contextual information aggregation for semantic segmentation.

The effect of the category consistent loss Tab. 4 also shows the performance on the Cityscapes validation set by adopting the proposed category consistent loss. The category consistent loss is denoted as “CCL” in the table. As we can see, adopting the category consistent loss could stably bring 0.7% mIoU gain with both Resnet-101 and Resnet-50, which prove the effectiveness of the proposed category consistent loss for semantic segmentation. To prove that the proposed piece-wise function is more robust than the original one, we conduct 10 times of the training processes using ResNet-50 for each kind of loss function. The training is deemed to fail when the loss value is NaN, thus we can calculate the success rate (number of successful training / total number of training). The experimental results in Table 3 demonstrate that using the piece-wise function has higher training success rate than using the original one. Besides, using the piece-wise function could achieve slightly better performance than a single quadratic function. Because we relax the punishment in the Eq. 6 to reduce the numerical values and gradients especially when the distance from the center exceeds δ_d . This relaxation makes the optimization much more stable.

Comparison of other context aggregation approaches We

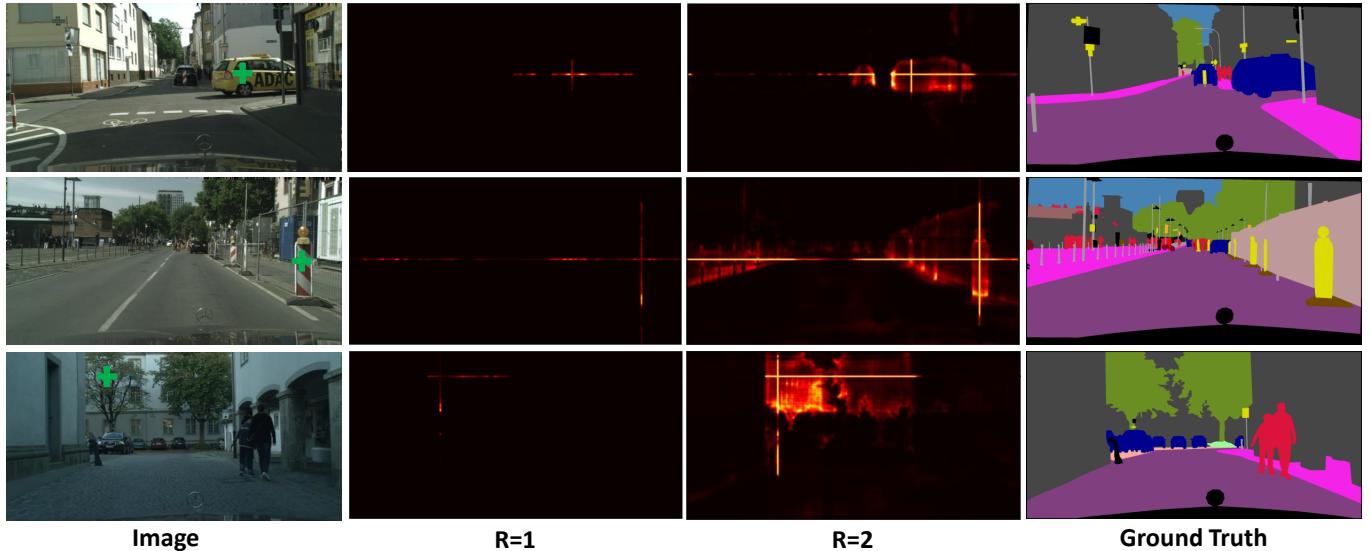


Fig. 7. Visualization of attention module on Cityscapes validation set. The left column is the input images, the 2 and 3 columns are pixel-wise attention maps when $R = 1$ and $R = 2$ in RCCA.

TABLE 4
Comparison of context aggregation approaches on Cityscapes (val).

| Method | mIOU(%) |
|-----------------------------|-------------|
| ResNet101-Baseline | 75.1 |
| ResNet101+GCN | 78.1 |
| ResNet101+PSP | 78.5 |
| ResNet101+ASPP | 78.9 |
| ResNet101+NL | 79.1 |
| ResNet101+RCCA($R=2$) | 79.8 |
| ResNet101+RCCA($R=2$)+CCL | 80.5 |
| ResNet50-Baseline | 73.3 |
| ResNet50+GCN | 76.2 |
| ResNet50+PSP | 76.4 |
| ResNet50+ASPP | 77.1 |
| ResNet50+NL | 77.3 |
| ResNet50+HV | 77.3 |
| ResNet50+HV&VH | 77.8 |
| ResNet50+RCCA($R=2$) | 78.5 |
| ResNet50+RCCA($R=2$)+CCL | 79.3 |

compare the performance of several different context aggregation approaches on the Cityscapes validation set with ResNet-50 and ResNet-101 as backbone networks.

Specifically, the baselines of context aggregation mainly include: 1) Peng *et al.* [49] utilized global convolution filters for contextual information aggregation, denoted as “+GCN”. 2) Zhao *et al.* [7] proposed Pyramid pooling which is the simple and effective way to capture global contextual information, denoted as “+PSP”; 3) Chen *et al.* [8] used different dilation convolutions to harvest pixel-wise contextual information at the different range, denoted as “+ASPP”; 4) Wang *et al.* [5] introduced non-local network for context aggregation, denoted as “+NL”.

In Tab. 4, both “+NL” and “+RCCA” achieve better

performance compared with the other context aggregation approaches, which demonstrates the importance of capturing full-image contextual information. More interestingly, our method achieves better performance than “+NL”. This reason may be attributed to the sequentially recurrent operation of criss-cross attention. Concretely, “+NL” generates an attention map directly from the feature which has limit receptive field and short-range dependencies. In contrast, our “+RCCA” takes two steps to form dense contextual information, leading to that the latter step can learn a better attention map benefiting from the feature map produced by the first step in which some long-range dependencies has already been embedded.

To prove the effectiveness of attention with criss-cross shape, we compare criss-cross shape with other shapes in Tab. 4. “+HV” means stacking horizontal attention and vertical attention. “+HV&VH” means summing up features of two parallel branches, i.e. “HV” and “VH”.

We further explore the amount of computation and memory footprint of RCCA. As shown in Tab. 5, compared with “+NL” method, the proposed “+RCCA” requires 11× less GPU memory usage and significantly reduces FLOPs by about 85% of non-local block in computing full-image dependencies, which shows that CCNet is an efficient way to capture full-image contextual information in the least amount of computation and memory footprint. To further prove the effectiveness of the recurrent operation, we also run non-local module in the recurrent way, denoted as “+NL($R=2$)”. As we can seen, the recurrent operation can bring more than 1 point gain. Because the recurrent operation leads to that the latter step can learn a better attention map benefiting from the feature map produced by the first step in which some long-range dependencies has already been embedded. However, compared with “+RCCA”, “+NL($R=2$)” needs huge GPU memory usage, which limits the use of self-attention.

Visualization of Attention Map To get a deeper under-

TABLE 5

Comparison of Non-local module and RCCA. FLOPs and memory increment are estimated for an input of $1 \times 3 \times 769 \times 769$.

| Method | GFLOPs(Δ) | Memory(M Δ) | mIOU(%) |
|------------|--------------------|---------------------|---------|
| baseline | 0 | 0 | 73.3 |
| +NL | 108 | 1411 | 77.3 |
| +NL(R=2) | 216 | 2820 | 78.7 |
| +RCCA(R=2) | 16.5 | 127 | 78.5 |

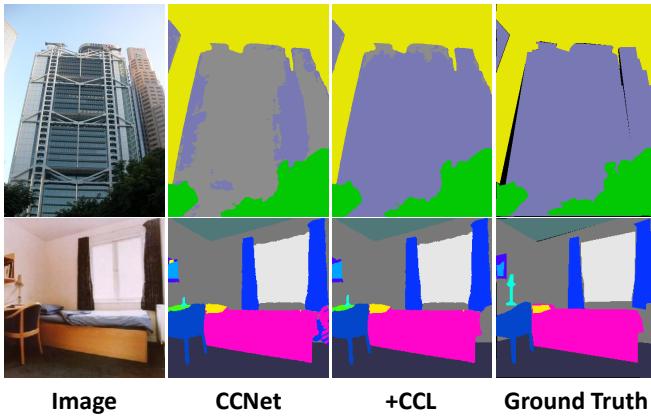


Fig. 8. Visualized examples on ADE20K val set with/without category consistent loss (CCL).

standing of our RCCA, we visualize the learned attention masks as shown in Fig. 7. For each input image, we select one point (cross in green) and show its corresponding attention maps when $R = 1$ and $R = 2$ in columns 2 and 3, respectively. It can be observed that only contextual information from the criss-cross path of the target point is captured when $R = 1$. By adopting one more criss-cross module, *i.e.*, $R = 2$, RCCA can finally aggregate denser and richer contextual information compared with that of $R = 1$. Besides, we observe that the attention module could capture semantic similarity and full-image dependencies.

4.4 Experiments on ADE20K

In this subsection, we conduct experiments on the AED20K dataset, which is a very challenging scene parsing dataset. As shown in Tab. 6, CCNet with CCL achieves the state-of-the-art performance of 45.76%, outperforms the previous state-of-the-art methods by more than 1.1% and also outperforms the conference version CCNet by 0.5%. Some successful segmentation results are given in Fig 8. Among the approaches, most of methods [7], [10], [12], [29], [65], [66] adopt the ResNet-101 as backbone and RefineNet [27] adopts a more powerful network, *i.e.*, ResNet-152, as the backbone. EncNet [10] achieves previous best performance among the methods and utilizes global pooling with image-level supervision to collect image-level context information. In contrast, our CCNet adopts an alternative way to integrate contextual information by capture full-image dependencies and achieve better performance.

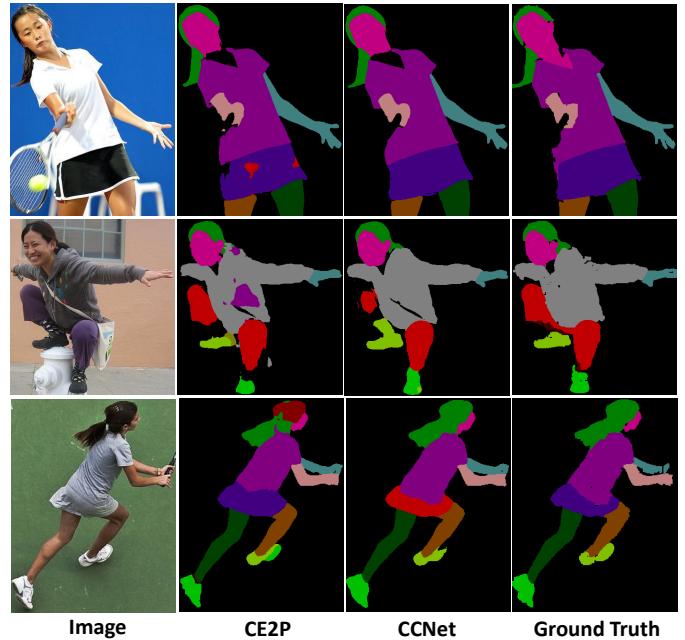


Fig. 9. Visualized examples for human parsing result on LIP val set.

4.5 Experiments on LIP

In this subsection, we conduct experiments on the LIP dataset, which is a very challenging human parsing dataset. The framework of CE2P [37] is utilized, with ImageNet pre-trained ResNet-101 as backbone and using RCCA ($R=2$) rather than PSP [7] as context embedding module. The category consistent loss is used to boost the performance. The hyper-parameter setting strictly follows that in the CE2P [37]. Among the approaches, Deeplab (VGG-16) [21], Attention [41] and SAN [38] adopt the VGG-16 as backbone and Deeplab (ResNet-101) [6], JPPNet [17], CE2P [37] and CCNet adopt ResNet-101 as the backbone. As shown in Tab. 7, CCNet achieves the state-of-the-art performance of 55.47%, outperforms the previous state-of-the-art methods by more than 2.3%. This significant improvement demonstrates the effectiveness of proposed method on human parsing task. Fig. 9 shows some visualized segmentation results. The top two rows show some successful segmentation results. It shows our method can produce accurate segmentation even for complicated poses. The third row shows a failure segmentation result where the “skirt” is misclassified as “pants”. But it’s difficult to recognize even for humans.

4.6 Experiments on COCO

To further demonstrate the generality of CCNet, we conduct the instance segmentation task on COCO [60] using the competitive Mask R-CNN model [62] as the baseline. Following [5], we modify the Mask R-CNN backbone by adding the RCCA module right before the last convolutional residual block of res4. We evaluate a standard baseline of ResNet-50/101. All models are fine-tuned from ImageNet pre-training. We use the official implementation¹ with end-to-end joint training whose performance is almost the same as the baseline reported in [5]. For fair comparison, we do

1. <https://github.com/facebookresearch/maskrcnn-benchmark>

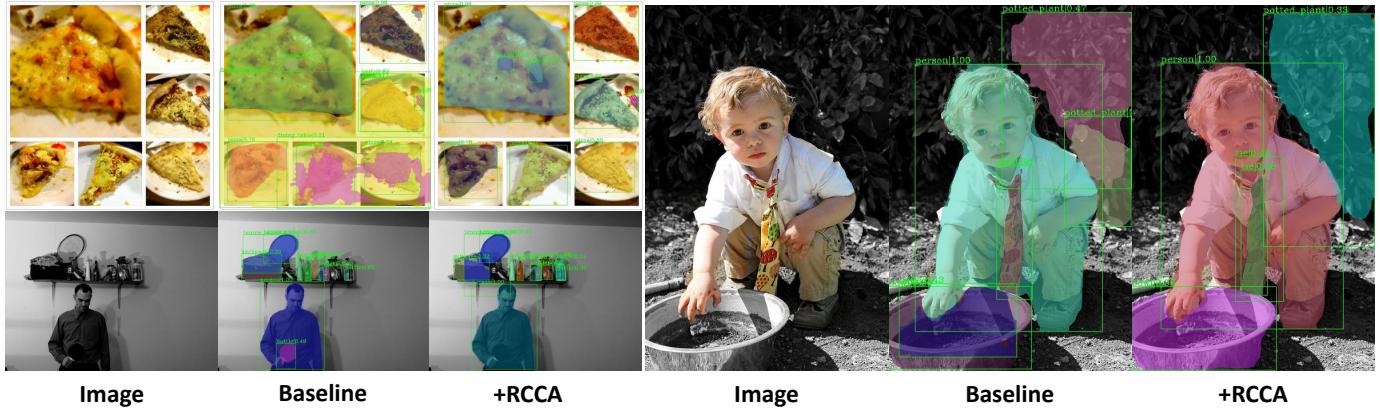


Fig. 10. Visualized examples for instance segmentation result on COCO val set.

not use the category consistent loss in our method. We report the results in terms of box AP and mask AP in Tab. 8 on COCO. The results demonstrate that our method substantially outperforms the baseline in all metrics. Some segmentation results for comparing baseline with “+RCCA” are given in Fig 10. Meanwhile, the network with “+RCCA” also achieves the better performance than the network with one non-local block “+NL”.

4.7 Experiments on CamVid

To further demonstrate the effectiveness of 3D-RCCA, we carry out the experiments on CamVid [61], which is one of the first datasets focusing on video semantic segmentation for driving scenarios. We follow the standard protocol proposed in [67] to split the dataset into 367 training, 101 validation and 233 test images. For fair comparison, we only report single-scale evaluation scores. As can be seen in Tab. 9, we achieve an mIoU of 79.1%, outperforming all other methods by a large margin.

To demonstrate the effectiveness of our proposed techniques, we perform training under the same settings with the different length of input frames. We apply the CNNs on each frame for extracting features and then concatenate and reshape them to satisfy the required shape of 3D Criss-Cross Attention module. We use the $R = 3$ for collecting dense spatial and temporal contextual information. Here, to make a training sample, we try two kinds of length (T) of input frames. For $T = 1$, we randomly sample 1 frame from a training video, denoted as “CCNet3D ($T = 1$)”. For $T = 5$, we sample 5 temporally ordered frames from a training video, denoted as “CCNet3D ($T = 5$)”. As can be seen in Tab. 9, “CCNet3D ($T = 5$)” outperforms “CCNet3D ($T = 1$)” by 1.2%.

5 CONCLUSION AND FUTURE WORK

In this paper, we have presented a Criss-Cross Network (CCNet) for deep learning based dense prediction tasks, which adaptively captures contextual information on the criss-cross path. To obtain dense contextual information, we introduce RCCA which aggregates contextual information from all pixels. The experiments demonstrate that RCCA

TABLE 6
Comparison with state-of-the-arts on ADE20K (val).

| Method | Backbone | mIoU(%) |
|----------------|------------|--------------|
| RefineNet [27] | ResNet-152 | 40.70 |
| SAC [29] | ResNet-101 | 44.30 |
| PSPNet [7] | ResNet-101 | 43.29 |
| PSANet [12] | ResNet-101 | 43.77 |
| DSSPN [65] | ResNet-101 | 43.68 |
| UperNet [66] | ResNet-101 | 42.66 |
| EncNet [10] | ResNet-101 | 44.65 |
| CCNet | ResNet-101 | 45.76 |

TABLE 7
Comparison with state-of-the-arts on LIP (val).

| Method | pixel acc | mean acc | mIoU |
|--------------------------|--------------|--------------|--------------|
| DeepLab (VGG-16) [6] | 82.66 | 51.64 | 41.64 |
| Attention [41] | 83.43 | 54.39 | 42.92 |
| SAN [38] | 84.22 | 55.09 | 44.81 |
| DeepLab (ResNet-101) [6] | 84.09 | 55.63 | 44.80 |
| JPPNet [17] | 86.39 | 62.32 | 51.37 |
| CE2P [37] | 87.37 | 63.20 | 53.10 |
| CCNet | 88.01 | 63.91 | 55.47 |

TABLE 8
Comparison on COCO (val).

| Method | AP^{box} | AP^{mask} |
|--------|------------|-------------|
| R50 | baseline | 38.2 |
| | +NL | 39.0 |
| | +RCCA | 39.3 |
| R101 | baseline | 40.1 |
| | +NL | 40.8 |
| | +RCCA | 41.0 |

captures full-image contextual information in less computation cost and less memory cost. Besides, to learn discriminative features, we introduce the category consistent loss.

TABLE 9
Results on the CamVid test set.

| Method | Backbone | mIoU (%) |
|-------------------|------------|-------------|
| SegNet [67] | VGG16 | 60.1 |
| RTA [36] | VGG16 | 62.5 |
| Dilate8 [22] | Dilate | 65.3 |
| BiSeNet [33] | ResNet18 | 68.7 |
| PSPNet [7] | ResNet50 | 69.1 |
| DenseDecoder [34] | ResNeXt101 | 70.9 |
| VideoGCRF‡ [35] | ResNet101 | 75.2 |
| CCNet3D (T=1) ‡ | ResNet101 | 77.9 |
| CCNet3D (T=5) ‡ | ResNet101 | 79.1 |

‡ the initialized model is pre-trained on Cityscapes.

Our CCNet achieves outstanding performance consistently on several semantic segmentation datasets, *i.e.*, Cityscapes, ADE20K, LIP, CamVid and instance segmentation dataset, *i.e.*, COCO. The source codes of CCNet are released to facilitate related research and applications.

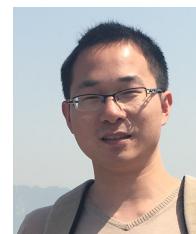
ACKNOWLEDGEMENTS

This work was in part supported by NSFC (No. 61733007 and No. 61876212), ARC DECRA DE190101315, ARC DP200100938, HUST-Horizon Computer Vision Research Center, and IBM-ILLINOIS Center for Cognitive Computing Systems Research (C3SR) - a research collaboration as part of the IBM AI Horizons Network.

REFERENCES

- [1] J. Fritsch, T. Kuehnl, and A. Geiger, "A new performance measure and evaluation benchmark for road detection algorithms," in *ITSC*, 2013, pp. 1693–1700. [1](#)
- [2] R. T. Azuma, "A survey of augmented reality," *Presence: Teleoperators & Virtual Environments*, vol. 6, no. 4, pp. 355–385, 1997. [1](#)
- [3] M. Evening, *Adobe Photoshop CS3 for photographers: a professional image editor's guide to the creative use of Photoshop for the Macintosh and PC*. Focal press, 2012. [1](#)
- [4] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *CVPR*, 2015, pp. 3431–3440. [1, 3](#)
- [5] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *CVPR*, 2018, pp. 7794–7803. [1, 2, 3, 9, 10](#)
- [6] L.-C. 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 TPAMI*, vol. 40, no. 4, pp. 834–848, 2018. [1, 2, 3, 4, 7, 10, 11, 11](#)
- [7] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *CVPR*, 2017, pp. 2881–2890. [1, 2, 3, 7, 9, 10, 11, 12](#)
- [8] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking atrous convolution for semantic image segmentation," *arXiv preprint arXiv:1706.05587*, 2017. [2, 7, 9](#)
- [9] H. Ding, X. Jiang, B. Shuai, A. Qun Liu, and G. Wang, "Context contrasted feature and gated multi-scale aggregation for scene segmentation," in *CVPR*, June 2018. [2](#)
- [10] H. Zhang, K. Dana, J. Shi, Z. Zhang, X. Wang, A. Tyagi, and A. Agrawal, "Context encoding for semantic segmentation," in *CVPR*, 2018. [2, 7, 10, 11](#)
- [11] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE TNN*, vol. 20, no. 1, pp. 61–80, 2008. [2, 3](#)
- [12] H. Zhao, Y. Zhang, S. Liu, J. Shi, C. C. Loy, D. Lin, and J. Jia, "Psanet: Point-wise spatial attention network for scene parsing," in *ECCV*, 2018, pp. 270–286. [2, 3, 7, 10, 11](#)
- [13] J. Cheng, L. Dong, and M. Lapata, "Long short-term memory-networks for machine reading," 2016. [2](#)
- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *NeurIPS*, 2017, pp. 5998–6008. [2, 3](#)
- [15] 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 *CVPR*, 2016, pp. 3213–3223. [2, 6](#)
- [16] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, "Scene parsing through ade20k dataset," in *CVPR*, 2017. [2, 6](#)
- [17] X. Liang, K. Gong, X. Shen, and L. Lin, "Look into person: Joint body parsing & pose estimation network and a new benchmark," *IEEE TPAMI*, vol. 41, no. 4, pp. 871–885, 2018. [2, 3, 6, 10, 11](#)
- [18] G. J. Brostow, J. Fauqueur, and R. Cipolla, "Semantic object classes in video: A high-definition ground truth database," *Pattern Recognition Letters*, vol. 30, no. 2, pp. 88–97, 2009. [2](#)
- [19] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778. [2, 7](#)
- [20] Z. Huang, X. Wang, L. Huang, C. Huang, Y. Wei, and W. Liu, "Ccnet: Criss-cross attention for semantic segmentation," in *ICCV*, 2019. [2](#)
- [21] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Semantic image segmentation with deep convolutional nets and fully connected crfs," *ICLR*, 2015. [3, 7, 10](#)
- [22] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *ICLR*, 2016. [3, 12](#)
- [23] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*, 2015, pp. 234–241. [3](#)
- [24] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," *ECCV*, 2018. [3, 7](#)
- [25] D. Lin, Y. Ji, D. Lischinski, D. Cohen-Or, and H. Huang, "Multi-scale context intertwining for semantic segmentation," in *ECCV*, 2018, pp. 603–619. [3](#)
- [26] B. Cheng, L.-C. Chen, Y. Wei, Y. Zhu, Z. Huang, J. Xiong, T. Huang, W.-M. Hwu, and H. Shi, "Spynet: Semantic prediction guidance for scene parsing," in *ICCV*, 2019. [3](#)
- [27] G. Lin, A. Milan, C. Shen, and I. D. Reid, "Refinenet: Multi-path refinement networks for high-resolution semantic segmentation," in *CVPR*, 2017. [3, 7, 10, 11](#)
- [28] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang, "Learning a discriminative feature network for semantic segmentation," *CVPR*, 2018. [3, 7](#)
- [29] R. Zhang, S. Tang, Y. Zhang, J. Li, and S. Yan, "Scale-adaptive convolutions for scene parsing," in *ICCV*, 2017, pp. 2031–2039. [3, 7, 10, 11](#)
- [30] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei, "Deformable convolutional networks," *ICCV*, 2017. [3](#)
- [31] Z. Liu, X. Li, P. Luo, C.-C. Loy, and X. Tang, "Semantic image segmentation via deep parsing network," in *ICCV*, 2015, pp. 1377–1385. [3](#)
- [32] T.-W. Ke, J.-J. Hwang, Z. Liu, and S. X. Yu, "Adaptive affinity field for semantic segmentation," *ECCV*, 2018. [3, 7](#)
- [33] C. Yu, J. Wang, C. Peng, C. Gao, G. Yu, and N. Sang, "Bisenet: Bilateral segmentation network for real-time semantic segmentation," *ECCV*, 2018. [3, 7, 12](#)
- [34] P. Bilinski and V. Prisacariu, "Dense decoder shortcut connections for single-pass semantic segmentation," in *CVPR*, 2018, pp. 6596–6605. [3, 12](#)
- [35] S. Chandra, C. Couprie, and I. Kokkinos, "Deep spatio-temporal random fields for efficient video segmentation," in *CVPR*, 2018, pp. 8915–8924. [3, 12](#)
- [36] P.-Y. Huang, W.-T. Hsu, C.-Y. Chiu, T.-F. Wu, and M. Sun, "Efficient uncertainty estimation for semantic segmentation in videos," in *ECCV*, 2018, pp. 520–535. [3, 12](#)
- [37] T. Ruan, T. Liu, Z. Huang, Y. Wei, S. Wei, and Y. Zhao, "Devil in the details: Towards accurate single and multiple human parsing," in *AAAI*, vol. 33, 2019, pp. 4814–4821. [3, 7, 10, 11](#)
- [38] Z. Huang, C. Wang, X. Wang, W. Liu, and J. Wang, "Semantic image segmentation by scale-adaptive networks," *IEEE TIP*, 2019. [3, 10, 11](#)
- [39] M. Yang, K. Yu, C. Zhang, Z. Li, and K. Yang, "Denseaspp for semantic segmentation in street scenes," in *CVPR*, 2018, pp. 3684–3692. [3, 7](#)

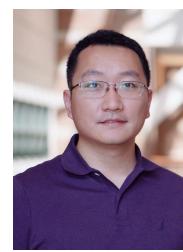
- [40] L.-C. Chen, M. D. Collins, Y. Zhu, G. Papandreou, B. Zoph, F. Schroff, H. Adam, and J. Shlens, "Searching for efficient multi-scale architectures for dense image prediction," *NeurIPS*, 2018. 3, 7
- [41] L.-C. Chen, Y. Yang, J. Wang, W. Xu, and A. L. Yuille, "Attention to scale: Scale-aware semantic image segmentation," in *CVPR*, 2016, pp. 3640–3649. 3, 10, 11
- [42] C. Liu, L.-C. Chen, F. Schroff, H. Adam, W. Hua, A. L. Yuille, and L. Fei-Fei, "Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation," in *CVPR*, 2019, pp. 82–92. 3
- [43] J. He, Z. Deng, L. Zhou, Y. Wang, and Y. Qiao, "Adaptive pyramid context network for semantic segmentation," in *CVPR*, 2019, pp. 7519–7528. 3
- [44] S. Liu, S. De Mello, J. Gu, G. Zhong, M.-H. Yang, and J. Kautz, "Learning affinity via spatial propagation networks," in *NeurIPS*, 2017, pp. 1520–1530. 3
- [45] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr, "Conditional random fields as recurrent neural networks," in *ICCV*, 2015, pp. 1529–1537. 3
- [46] Y. Yuan and J. Wang, "Ocnet: Object context network for scene parsing," *arXiv preprint arXiv:1809.00916*, 2018. 3
- [47] J. Fu, J. Liu, H. Tian, Z. Fang, and H. Lu, "Dual attention network for scene segmentation," *CVPR*, 2019. 3
- [48] Y. Chen, M. Rohrbach, Z. Yan, Y. Shuicheng, J. Feng, and Y. Kalantidis, "Graph-based global reasoning networks," in *CVPR*, 2019, pp. 433–442. 3
- [49] C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun, "Large kernel matters—improve semantic segmentation by global convolutional network," in *CVPR*, 2017, pp. 1743–1751. 3, 7, 9
- [50] A. Sperduti and A. Starita, "Supervised neural networks for the classification of structures," *IEEE TNN*, vol. 8, no. 3, pp. 714–735, 1997. 3
- [51] M. Gori, G. Monfardini, and F. Scarselli, "A new model for learning in graph domains," in *IJCNN*, vol. 2, 2005, pp. 729–734. 3
- [52] M. Henaff, J. Bruna, and Y. LeCun, "Deep convolutional networks on graph-structured data," *arXiv preprint arXiv:1506.05163*, 2015. 3
- [53] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *NeurIPS*, 2016, pp. 3844–3852. 3
- [54] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016. 3
- [55] R. Levie, F. Monti, X. Bresson, and M. M. Bronstein, "Cayleynets: Graph convolutional neural networks with complex rational spectral filters," *IEEE TSP*, vol. 67, no. 1, pp. 97–109, 2018. 3
- [56] J. Atwood and D. Towsley, "Diffusion-convolutional neural networks," in *NeurIPS*, 2016, pp. 1993–2001. 3
- [57] M. Niepert, M. Ahmed, and K. Kutzkov, "Learning convolutional neural networks for graphs," in *ICML*, 2016, pp. 2014–2023. 3
- [58] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," in *ICML*, 2017, pp. 1263–1272. 3
- [59] B. De Brabandere, D. Neven, and L. Van Gool, "Semantic instance segmentation with a discriminative loss function," *arXiv preprint arXiv:1708.02551*, 2017. 5
- [60] 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 *ECCV*, 2014, pp. 740–755. 6, 10
- [61] G. J. Brostow, J. Shotton, J. Fauqueur, and R. Cipolla, "Segmentation and recognition using structure from motion point clouds," in *ECCV*, 2008, pp. 44–57. 6, 11
- [62] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *ICCV*, 2017, pp. 2980–2988. 7, 10
- [63] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, and G. Cotterell, "Understanding convolution for semantic segmentation," in *WACV*, 2018, pp. 1451–1460. 7
- [64] Z. Wu, C. Shen, and A. v. d. Hengel, "Wider or deeper: Revisiting the resnet model for visual recognition," *arXiv preprint arXiv:1611.10080*, 2016. 7
- [65] X. Liang, H. Zhou, and E. Xing, "Dynamic-structured semantic propagation network," in *CVPR*, 2018, pp. 752–761. 10, 11
- [66] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun, "Unified perceptual parsing for scene understanding," *ECCV*, 2018. 10, 11
- [67] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE TPAMI*, vol. 39, no. 12, pp. 2481–2495, 2017. 11, 12



Zilong Huang is a Ph.D. student in the School of Electronics Information and Communications, Huazhong University of Science and Technology (HUST). He received his B.S. degree from HUST in 2015. His research interests include computer vision and machine learning. In particular, he focuses on semantic segmentation and object parsing.



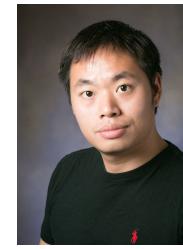
Xinggang Wang received the B.S. and Ph.D. degrees in Electronics and Information Engineering from Huazhong University of Science and Technology (HUST), Wuhan, China, in 2009 and 2014, respectively. He is currently an Associate Professor with the School of Electronic Information and Communications, HUST. His research interests include computer vision and machine learning.



Yunchao Wei received his Ph.D. degree from Beijing Jiaotong University, Beijing, China. He was a Postdoctoral Researcher at Beckman Institute, UIUC, from 2017 to 2019. He is currently an Assistant Professor with the Centre for Artificial Intelligence, University of Technology Sydney. He is ARC Discovery Early Career Researcher Award (DECRA) Fellow from 2019 to 2021. His current research interests include computer vision and machine learning.



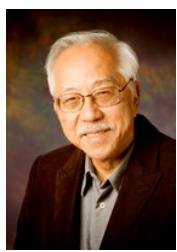
Lichao Huang received his B.S. degree in Software Engineering from Sun Yat-Sen University, Guangzhou, China in 2013, and the MSc degree in Computer Vision from Imperial College London in 2014. He was a senior researcher and developer at Institution of Deep Learning, Baidu, and now he is an algorithm engineer at Horizon Robotics. His research interests includes deep learning and computer vision.



Humphrey Shi received his B.Eng from Tsinghua University, and MS/PhD from UIUC. He is an Assistant Professor at the University of Oregon and a Graduate Faculty Member at UIUC. Before that, he was a Research Staff Member at IBM T.J. Watson Research Center. He is the winner of a dozen of major international AI competitions in visual recognition. He received highest recognitions and multiple corporate-level awards during his two-year tenure at IBM Research. He has been a PI/Co-PI of multiple industrial, academic and government projects with focuses on not only pushing the envelope of core AI research such as accurate and efficient visual understanding, but also connecting cutting edge AI research with real world data, problems and tasks at scale.



Wenyu Liu (SM'15) received the B.S. degree in Computer Science from Tsinghua University, Beijing, China, in 1986, and the M.S. and Ph.D. degrees, both in Electronics and Information Engineering, from Huazhong University of Science and Technology (HUST), Wuhan, China, in 1991 and 2001, respectively. He is now a professor and associate dean of the School of Electronic Information and Communications, HUST. His current research areas include computer vision, multimedia, and machine learning.



Thomas S. Huang (1936-2020) received his Sc.D. from MIT in 1963. He had been a faculty member at MIT, Purdue and UIUC over the past five decades. He was the Swanlund Endowed Chair Emeritus with UIUC. He was a leading pioneer in computer vision, image processing and multi-modal signal processing. He had received numerous awards, including the IEEE Jack Kilby Signal Processing Medal, the King-Sun Fu Prize of the IAPR, the Azriel Rosenfeld Life Time Achievement Award at ICCV, and ICIP Pioneer Award. He was a Life Fellow of IEEE, IAPR, SPIE and OSA. He was a member of the US National Academy of Engineering.