

DEEP FEATURE SELECTION-AND-FUSION FOR RGB-D SEMANTIC SEGMENTATION

Yuejiao Su, Yuan Yuan, Zhiyu Jiang*

School of Computer Science and School of Artificial Intelligence, Optics and Electronics (iOPEN),
Northwestern Polytechnical University, Xi'an 710072, P.R. China
yuejiao@mail.nwpu.edu.cn; y.yuan1.ieee@gmail.com; jiangzhiyu@nwpu.edu.cn

ABSTRACT

Scene depth information can help visual information for more accurate semantic segmentation. However, how to effectively integrate multi-modality information into representative features is still an open problem. Most of the existing work uses DCNNs to implicitly fuse multi-modality information. But as the network deepens, some critical distinguishing features may be lost, which reduces the segmentation performance. This work proposes a unified and efficient feature selection-and-fusion network (FSFNet), which contains a symmetric cross-modality residual fusion module used for explicit fusion of multi-modality information. Besides, the network includes a detailed feature propagation module, which is used to maintain low-level detailed information during the forward process of the network. Compared with the state-of-the-art methods, experimental evaluations demonstrate that the proposed model achieves competitive performance on two public datasets.

Index Terms— RGB-D Semantic Segmentation, Multi-modality, Skip-connection, Attention Mechanism

1. INTRODUCTION

Semantic segmentation refers to the pixel-wise classification of images according to semantic information. Besides widely-used visual information, the depth information is regarded as another supplementary information to improve the scenario understanding performance due to the development of depth sensors. Depth modality contains 3D geometric information, which is insensitive to illumination changes and can distinguish objects with similar appearance. Therefore, depth cues can make up for some of the defects of semantic segmentation using only visual cues. RGB-D semantic segmentation is very important for many applications such as autonomous driving [1], robot vision and understanding [2], and land cover classification [3], *etc.*

With the development of deep learning, two-stream networks have achieved remarkable performance in RGB-D semantic segmentation [4, 5, 6]. As we all know, the information of RGB and depth modalities are complementary. However, how to effectively fuse RGB and depth information into a unified and distinguishing representation is still a basic yet difficult issue in RGB-D semantic segmentation. There are many methods to try to solve this problem. FuseNet [7] and RedNet [8] integrated different modalities by directly adding the feature maps of depth to RGB. RFBNet [9] proposed a residual fusion block to achieve bottom-up interaction and fusion between two modalities. ACNet [10] proposed an attention complementary module to assign different modality-weights for better integration. RDFNet [11] used single-modality residual learning to learn residual RGB and depth features and their combinations to exploit the complementary characteristics. Although these methods provide structured models to integrate the two kinds of information, it is still an unresolved problem to ensure that the network makes full use of the information from both modalities for fine semantic segmentation.

Moreover, the loss of detailed information during down-sampling is an inherent property of convolution operation. For the encoding part of the segmentation network, reducing the resolution of the feature map to a very small level through various pooling layers is not conducive to accurate mask generation, which can lead to inaccurate segmentation results. To further make up for the information lost in the encoding stage, U-Net [12] proposed skip-connection to reuse the feature to assist up-sampling learning and recover the fine segmentation results. Although it can realize the reuse of some lost features, it lacks pertinence and does not explicitly model the recovery of detailed information.

To solve the above problems, this work proposes a novel feature selection-and-fusion network to explicitly strengthen features in RGB-D semantic segmentation model from two aspects: multi-modality representations and decoder features. The key idea of our proposed network is to select discriminative information from one modality to supplement the other modality to obtain well-informed representation. In addition, this work focuses on the lost information in the encoder and finds ways to make it helpful in predicting the final result.

2021 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works..
*Corresponding author: Zhiyu Jiang (jiangzhiyu@nwpu.edu.cn)

These two aspects correspond to two modules respectively. For the former, the Symmetric Cross-modality Residual Fusion module (SCRf) is designed to effectively fuse the complementary information of two modalities, while maintaining the specificity of the specific modality during the information interaction process at the encoder stage. For the latter, a Detailed Feature Propagation module (DFP) is designed to encourage the network to spotlight the missing vital details in the encoder and reuse them in the decoder to improve the segmentation performance. Both modules are designed as two steps: feature selection and feature fusion.

The main contributions of this work are described as follows:

- In order to solve the multi-modality information fusion in RGB-D semantic segmentation, this work designs the SCRf module in FSFNet. The core of the module is the cross-modality residual connection, which can retain the advantages of the residual connection and can explicitly select and fuse complementary information into distinguishing and effective representations.
- In response to the loss of some important information during the down-sampling process, this work designs the DFP module in the network. The DFP module firstly selects vital information that may be lost in the encoder stage by attention mechanism. And then the module propagates and fuse the selected features with decoder features for further segmentation.
- With proposed modules, our FSFNet uses a relatively simple architecture to achieve excellent performance. We verify the effectiveness of FSFNet and its modules through a series of experiments and achieve competitive or superior performance on NYUDv2 and SUN RGB-D datasets.

2. RELATED WORK

The main difference between RGB-D semantic segmentation and RGB semantic segmentation is that the former not only use visual information but also leverages scene depth information to achieve better accuracy [13]. According to whether the deep learning methods are used or not, RGB-D semantic segmentation methods are roughly divided into traditional and deep learning-based methods.

In the early stage, researchers preferred to use depth information directly. Silberman *et al.* [14] recovered support relationships by parsing indoor scenes into floor, walls, supporting surfaces, and object regions using depth information. Ren *et al.* [15] achieved high labeling accuracy by a combination of color and depth features using kernel descriptors, and by combining MRF with segmentation tree. Gupta *et al.* [16] proposed algorithms for object boundary detection and hierarchical segmentation that generalize the gPb-ucm [17] approach by making effective use of depth information.

Although the traditional methods can mine the internal relationship between RGB and depth information by explicitly utilizing depth information, they need a lot of prior knowledge and specific descriptors. Compared with traditional methods, deep learning-based approaches implicitly utilizes deep information to assist RGB semantic segmentation through various networks. Couprie *et al.* [18, 19] regarded the depth information as another channel to concatenate with RGB image and then extract features with RGB semantic segmentation network. However, RGB channel and depth channel contain inconsistent features and cannot be processed by shared network feature extractors. Besides, Gupta *et al.* [5] encoded depth into HHA (Horizontal disparity, Height above ground, and Angle with gravity). Some studies [8, 9, 10, 20] used two-stream networks to process the RGB images and HHA information. These methods prove that depth data can improve the performance of semantic segmentation. Many of the later researchers focus on the effective fusion of multi-modality data. Proposed by Hazirbas *et al.* [7], FuseNet integrated depth characteristics into RGB feature mapping by dense fusion or sparse fusion strategies as the network deepens. Lin *et al.* [21] divided the image into several branches with different scene resolutions based on depth information aiming at the multi-scale problem. Compared with single-scale networks, multi-scale networks [22] have better segmentation performance, but they also require more computation. Different from the previous work, this work extends the idea of residual connection [23] to multi-modality and designs a SCRf module based on multi-modal residual connection to clearly promote multi-modality feature fusion.

3. PROPOSED METHOD

3.1. Overview

The RGB-D semantic segmentation model is usually based on the encoder-decoder architecture. However, a superior encoder can unify both the complementary characteristics of two modalities and their specific characteristics into effective representation. To this end, inspired by the residual connection [23], we put forward a unified encoder framework that includes the symmetric cross-modality residual fusion module as described in Section 3.2, which aims to encourage explicit fusion of cross-modality information and preserve single-modality specific characteristics as completely as possible. In addition, in the encoder phase, a lot of detailed information is lost due to cascading downsampling, which is fatal to semantic segmentation. Our goal is to automatically reuse the essential details lost in the encoder stage and make them helpful to form the final segmentation mask. For this purpose, the detailed feature propagation module is proposed between encoder and decoder, as described in Section 3.3.

The overall frame diagram of the proposed method is shown in Fig. 1(a). Take the three-channel RGB image

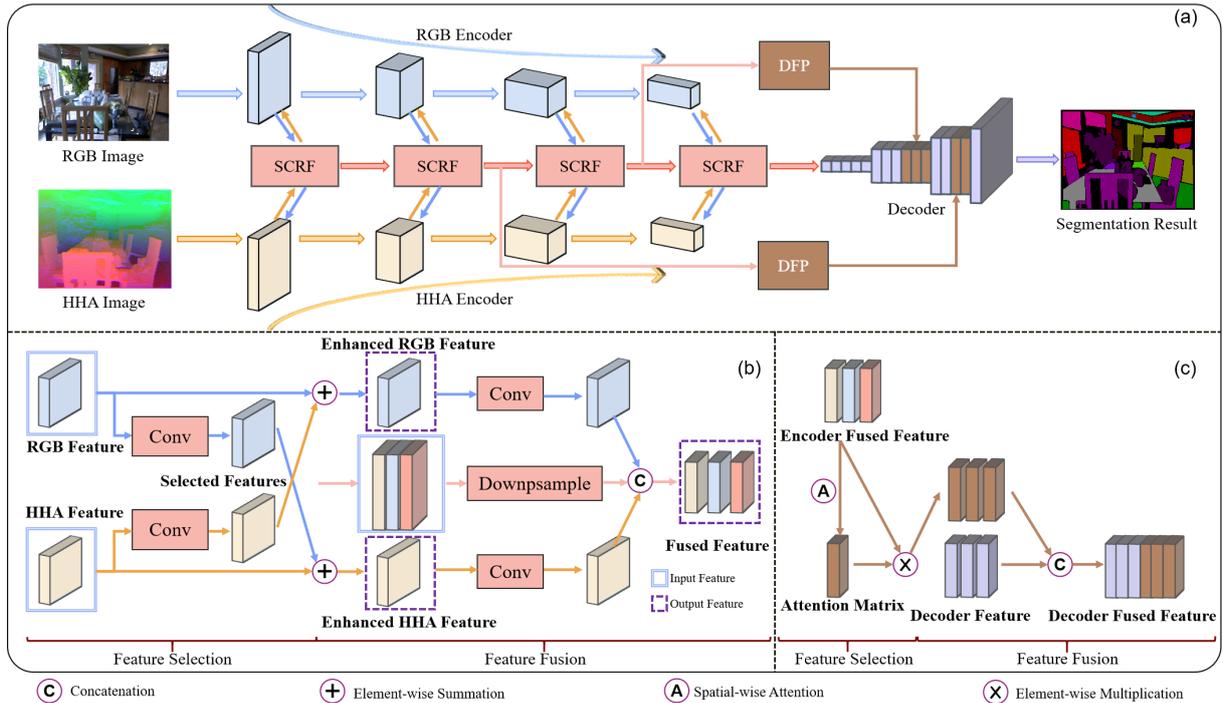


Fig. 1. (a) Overview of the proposed framework. Given RGB and HHA images as input, the upper and lower encoder branches extract the characteristics of specific modalities from input respectively. The middle fusion branch uses cascaded SCRF modules to fuse the two modal features. The fused features of the middle two layers of the fusion branch are selected by the DFP module and propagated to the corresponding decoder layer for joint prediction. (b) Details of SCRF module. It is based on the cross-modality residual connection. The SCRF firstly select features that complement another modality from one modality, and then perform feature fusion between modalities and levels. (c) Details of DFP module. The DFP firstly uses spatial-wise attention to select important but possibly lost detailed information from the fusion features in the middle two layers of the encoder stage, and then merge them with the corresponding features of the decoder stage for the final joint segmentation.

$\mathbf{I}_R \in \mathbb{R}^{H \times W \times 3}$ and three-channel HHA image $\mathbf{I}_H \in \mathbb{R}^{H \times W \times 3}$ as input, our network selects and enhances the information representation capability of two modalities through cascaded SCRF modules, and at the same time encourages to save the specific features of the specific modality as much as possible. In addition, partial details of the encoder are selected and transmitted by the DFP module to the corresponding stage of the decoder to make full use of the lost important details. Each component will be described in more details in the remaining parts of this section.

3.2. Symmetric Cross-modality Residual Fusion

To encourage the complementary fusion of information between two modalities, we consider this issue from two aspects. Firstly, we argue that if the complementary part can be modeled more explicitly, cross-modality information fusion can be accomplished better. For this motivation, the SCRF module in every layer is designed into two steps, feature selection and feature fusion, as shown in Fig. 1(b). Unlike other work using the simple single-modality residual con-

nection [11], we propose a cross-modality residual connection to select and fuse complementary features from another modality to make both steps effectively. Secondly, as we know, shallow features can identify edge information, while deep features can study the global context to locate salient objects [24]. Therefore, the deep features are very irregular, while the shallow features are very noisy and chaotic, and a lot of previous work lacked cross-level interaction in the fusion branch. So in this paper, each SCRF module of the fusion branch is cascaded from shallow to deep, so that the network can gradually select complementary features across levels for joint decision-making. Assume that the RGB feature map of j -th layer is $\mathbf{F}_{\text{rgb}}^j \in \mathbb{R}^{H \times W \times C}$, and the HHA feature map of j -th layer is $\mathbf{F}_{\text{hha}}^j \in \mathbb{R}^{H \times W \times C}$. For simplicity, H , W , and C are the height, width, and channel of the feature map respectively. So the selection process can be denoted as:

$$\mathbf{S}_{\text{hha}}^j = f_{s_1}(\mathbf{F}_{\text{hha}}^j), \mathbf{S}_{\text{rgb}}^j = f_{s_2}(\mathbf{F}_{\text{rgb}}^j), \quad (1)$$

where $\mathbf{S}_{\text{rgb}}^j$ and $\mathbf{S}_{\text{hha}}^j$ are selected RGB and HHA features of j -th layer. They are automatically selected and they can help

improve the distinguishing ability of another modality. Note that $f_{s_1}(\cdot)$ and $f_{s_2}(\cdot)$ are feature selection operations. In practice, we use 1×1 convolution with stride 1.

And then the selected features are sent to the feature fusion step. Suppose that the fusion feature of the previous layer of j -th layer is $\mathbf{F}_{\text{fuse}}^{j-1} \in \mathbb{R}^{H \times W \times C}$. Then the feature after fusion is:

$$\mathbf{F}_{\text{fuse}}^j = [f_{\text{down}}(\mathbf{F}_{\text{fuse}}^{j-1}), f_{\text{conv}}(\mathbf{S}_{\text{hha}}^j + \mathbf{F}_{\text{rgb}}^j), f_{\text{conv}}(\mathbf{S}_{\text{rgb}}^j + \mathbf{F}_{\text{hha}}^j)]. \quad (2)$$

Here $\mathbf{F}_{\text{fuse}}^j$ is the concatenated fused feature. $f_{\text{down}}(\cdot)$ and $f_{\text{conv}}(\cdot)$ mean the down-sample and convolution operation respectively. Note that when $j = 1$, only the last two items in equation (2) are concatenated.

Therefore cross-modality residual function we defined includes feature selection, feature fusion, and the final convolution part, that is:

$$\begin{aligned} f_{r_1} : \mathbf{F}_{\text{rgb}}^j &\rightarrow f_{\text{conv}}(\mathbf{S}_{\text{hha}}^j + \mathbf{F}_{\text{rgb}}^j), \\ f_{r_2} : \mathbf{F}_{\text{hha}}^j &\rightarrow f_{\text{conv}}(\mathbf{S}_{\text{rgb}}^j + \mathbf{F}_{\text{hha}}^j), \end{aligned} \quad (3)$$

where $f_{r_1}(\cdot)$ and $f_{r_2}(\cdot)$ are the cross-modality residual functions we defined.

The simple but effective SCRF module not only continues the advantages of residual connection, that is, it can accelerate the convergence of the network, but also it uses adaptive weight learning to intelligently and explicitly select and fuse supplementary features from another modality. And our work enables the network to automatically select deep or shallow features by cascading the module. These all allow us to effectively aggregate multi-modality information.

3.3. Detailed Feature Propagation

As mentioned above, the semantic information of some small-scale objects and outlines will be lost in the down-sampling stage, which is unfavorable for semantic segmentation. Therefore, inspired by the idea of skip-connection [12], we design the detailed feature propagation module to reuse the lost but essential information, as shown in Fig. 1(c).

The DFP module is still designed with the same two steps as the SCRF module, *i.e.* feature selection and feature fusion. Different from SCRF, DFP automatically selects some detailed features that may be lost in the encoder stage by spatial-wise attention, and then propagates the selected features to the corresponding decoder part for fusion with the features being up-sampled. In this way, the features of some small scale objects and outlines can be well used.

For feature selection, we make the model select detailed features only from the second and third layers of the encoder. This is because the first layer contains a lot of noise and is very complex, and the features of the fourth layer no longer contain enough detailed features after the cascading of the down-sampling operations.

Accordingly, our feature fusion operation is carried out in the corresponding decoder layer, to achieve the enhancement of features in the decoder, which can improve the segmentation ability of some small objects and outlines. Suppose fused feature maps in i -th layer of encoder denoted as $\mathbf{F}_{\text{encoder}}^i \in \mathbb{R}^{H \times W \times C}$, the feature maps of the corresponding decoder are $\mathbf{F}_{\text{decoder}}^{m-i} \in \mathbb{R}^{H \times W \times C}$, then the enhanced decoder features after feature selection-and-fusion can be calculated by:

$$\tilde{\mathbf{F}}_{\text{decoder}}^{m-i} = f_{\text{fuse}}(f_{\text{select}}(\mathbf{F}_{\text{encoder}}^i), \mathbf{F}_{\text{decoder}}^{m-i}), i \in \{2, 3\}, m = 4. \quad (4)$$

Here, $\tilde{\mathbf{F}}_{\text{decoder}}^{m-i}$ represents the feature maps after fusion, $i \in \{2, 3\}$ means that the DFP module only works in the second and third layers, while $f_{\text{select}}(\cdot)$ and $f_{\text{fuse}}(\cdot)$ represents feature selection operation and feature fusion operation, respectively. In practice, we use the spatial-wise attention as the feature selection block to select some useful but may be lost features, and we transfer the processed features to the corresponding decoder layer like skip-connection. Simple concatenation is used as the feature fusion operation.

Through explicitly selecting and fusing detailed features, the impact of the loss of detailed features in down-sampling on the final segmentation can be reduced. And the semantic information of some small-scale objects and outlines can be preserved to realize the reuse of features, which can play their role in the prediction of the final segmentation mask.

In addition, we use pyramid supervision [10] in our work. In other words, the output of the last three layers of the decoder is also used as supervision information to ensure rapid network convergence. The only difference between the supervision information of the intermediate output and the ground truth is the resolution. We obtain the supervision information of the intermediate output through the nearest neighbor interpolation down-sampling. The final loss function is as follows:

$$L_{\text{final}} = \sum_{i=1}^3 \lambda_i l_i, \quad (5)$$

where l_i and λ_i represents the loss function and its weight of the i -th layer, respectively. The weighted cross-entropy loss function is the loss function of each layer. Our hyperparameter settings refer to ACNet [10].

4. EXPERIMENTS

4.1. Datasets and Evaluation Metrics

NYUDv2 [14] and SUN RGB-D [25] datasets are used to evaluate the proposed network. NYUDv2 dataset contains 1449 densely labeled pairs of RGB-D images captured by Microsoft Kinect. There are 795 pairs of images for training and 654 pairs for testing. The SUN RGB-D dataset is the largest RGB-D semantic segmentation dataset currently, with

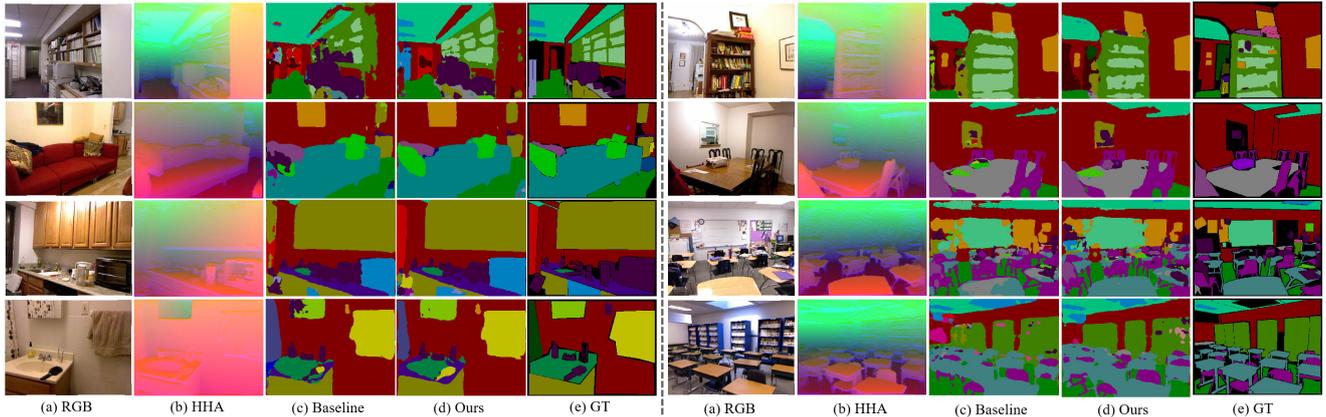


Fig. 2. Visualization results of this work in NYUDv2 (left) and SUN RGB-D (right) test sets. For each dataset, we show (a) RGB image, (b) HHA image, (c) result of baseline, (d) result of ours, and (e) ground truth. The baseline directly adds feature maps in HHA branch into RGB branch.

10,335 densely annotated RGB-D images taken from 20 different scenes. It is captured by four different sensors (Kinect V1, Kinect V2, Xtion, and RealSense). The officially divided training set consists of 5285 pairs of RGB-D images and labels, and the remaining 5050 pairs are used for testing. The number of classes in both datasets is 40.

As with most related work, mean Intersection-over-Union (mIoU) and Pixel Accuracy (Pixel Acc.) are used as our performance metrics.

4.2. Implementation Details

PyTorch framework is used in our work and we use Stochastic Gradient Descent (SGD) optimizer to train our model with a momentum of 0.9 and a weight decay of 0.0005. Initial learning rate in our work is set to 0.02 and decreased at a rate of 0.9. The input RGB-D images are cropped to 480×480 , and the batch size is set to 10.

4.3. Experimental Results

To prove the effectiveness of the proposed model, we compare it with several state-of-the-art methods, as shown in Table 1. The result shows that on NYUDv2 dataset, our FSFNet performs equivalently to the current best model using more simple architecture. And on SUN RGB-D dataset, our FSFNet outperforms other models in mIoU. That’s because our model uses cascaded SCRF modules to implicitly integrate multi-modality features and uses the DFP module to reuse the selected detailed information from the encoder, which improves the ability to segment small-scale objects and contours.

In order to verify the effectiveness of the SCRF and DFP modules, we perform ablation studies on the NYUDv2 dataset under the same parameter settings as shown in Table 1. We use ResNet-101 as our backbone and the first row in Table 2 is the baseline that fuses RGB and depth feature maps by

Table 1. Comparing with state-of-the-art methods on NYUDv2 and SUN RGB-D test sets.

Method	NYUDv2		SUN RGB-D	
	mIoU	Pixel Acc.	mIoU	Pixel Acc.
3DGNN [26]	43.1%	-	45.9%	-
Kong <i>et al.</i> [27]	44.5%	72.1%	45.1%	80.3%
RedNet [8]	-	-	47.8%	81.3%
CFN [21]	47.7%	-	48.1%	-
ACNet [10]	48.3%	-	48.1%	-
PAP [28]	50.4%	76.2%	50.5%	83.8%
SA-Gate [20]	52.4%	77.9%	49.4%	82.5%
Ours	52.0%	77.9%	50.6%	81.8%

Table 2. Ablation study for SCRF and DFP modules on NYUDv2 dataset.

Method	mIoU(%)
Res-101 + SUM	47.9
Res-101 + DFP	49.6(1.7%↑)
Res-101 + SCRF	50.8(2.9%↑)
Res-101 + SCRF + DFP	52.0(4.1%↑)

element-wise summation in each encoder layer. We can observe that SCRF and DFP can improve performance by 2.9% and 1.7% respectively. And when two modules exist at the same time, the performance improvement is more obvious than the baseline. This experiment proves the effectiveness and importance of the two proposed modules.

In order to display the results of the model more intuitively, we show a part of the visualization results in Fig. 2. It can be observed that compared to the baseline, our model has a better segmentation performance on different classes of objects, such as ceiling, table, *etc.*, which can prove that our model combines the characteristics of RGB and depth information well. In addition, our model can distinguish small-scale objects and can perform more accurate contour segmen-

tation, which can prove the effectiveness of DFP.

5. CONCLUSIONS

In this work, we propose a feature selection-and-fusion network to address two main challenges in RGB-D semantic segmentation. Firstly, for the effective fusion of multi-modality information, we put forward the cascaded SCRF module to obtain unified multi-modality representations. Secondly, aimed at the loss of detailed information in the down-sampling stage, we design the DFP module to make the important details helpful in predicting the results. Experimental results demonstrate our model achieves competitive performance on NYUDv2 and SUN RGB-D datasets.

6. REFERENCES

- [1] Kaihong Yang, Sheng Bi, and Min Dong, "Lightningnet: Fast and accurate semantic segmentation for autonomous driving based on 3D LIDAR point cloud," in *Proc. IEEE International Conference on Multimedia and Expo*, 2020, pp. 1–6.
- [2] Lingni Ma, Jörg Stückler, Christian Kerl, and Daniel Cremers, "Multi-view deep learning for consistent semantic mapping with RGB-D cameras," in *Proc. IEEE International Conference on Intelligent Robots and Systems*, 2017, pp. 598–605.
- [3] Shu Liu, Jia-Li He, and Sheng-Hui Liao, "Automatic detection of anatomical landmarks on geometric mesh data using deep semantic segmentation," in *Proc. IEEE International Conference on Multimedia and Expo*, 2020, pp. 1–6.
- [4] Yanhua Cheng, Rui Cai, Zhiwei Li, Xin Zhao, and Kaiqi Huang, "Locality-sensitive deconvolution networks with gated fusion for RGB-D indoor semantic segmentation," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1475–1483.
- [5] Saurabh Gupta, Ross B. Girshick, Pablo Andrés Arbeláez, and Jitendra Malik, "Learning rich features from RGB-D images for object detection and segmentation," in *Proc. European Conference on Computer Vision*, 2014, pp. 345–360.
- [6] Zhitong Xiong, Yuan Yuan, Nianhui Guo, and Qi Wang, "Variational context-deformable convnets for indoor scene parsing," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3991–4001.
- [7] Caner Hazirbas, Lingni Ma, Csaba Domokos, and Daniel Cremers, "FuseNet: Incorporating depth into semantic segmentation via fusion-based CNN architecture," in *Proc. Asian Conference on Computer Vision*, 2016, pp. 213–228.
- [8] Jindong Jiang, Lunan Zheng, Fei Luo, and Zhijun Zhang, "RedNet: Residual encoder-decoder network for indoor RGB-D semantic segmentation," *arXiv*, 2018.
- [9] Liuyuan Deng, Ming Yang, Tianyi Li, Yuesheng He, and Chunxiang Wang, "RFBNet: Deep multimodal networks with residual fusion blocks for RGB-D semantic segmentation," *arXiv*, 2019.
- [10] Xinxin Hu, Kailun Yang, Lei Fei, and Kaiwei Wang, "ACNET: Attention based network to exploit complementary features for RGBD semantic segmentation," in *Proc. IEEE International Conference on Image Processing*, 2019, pp. 1440–1444.
- [11] Seungyong Lee, Seong-Jin Park, and Ki-Sang Hong, "RDFNet: RGB-D multi-level residual feature fusion for indoor semantic segmentation," in *Proc. IEEE International Conference on Computer Vision*, 2017, pp. 4990–4999.
- [12] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Proc. Medical Image Computing and Computer-Assisted Intervention*, 2015, pp. 234–241.
- [13] Zhiyu Jiang, Yuan Yuan, and Qi Wang, "Contour-aware network for semantic segmentation via adaptive depth," *Neurocomputing*, vol. 284, pp. 27–35, 2018.
- [14] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus, "Indoor segmentation and support inference from RGBD images," in *Proc. European Conference on Computer Vision*, 2012, pp. 746–760.
- [15] Xiaofeng Ren, Liefeng Bo, and Dieter Fox, "RGB-(D) scene labeling: Features and algorithms," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 2759–2766.
- [16] Saurabh Gupta, Pablo Arbeláez, and Jitendra Malik, "Perceptual organization and recognition of indoor scenes from RGB-D images," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 564–571.
- [17] Pablo Arbeláez, Michael Maire, Charless C. Fowlkes, and Jitendra Malik, "Contour detection and hierarchical image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 5, pp. 898–916, 2011.
- [18] Camille Couprie, Clément Farabet, Laurent Najman, and Yann LeCun, "Indoor semantic segmentation using depth information," in *Proc. International Conference on Learning Representations*, 2013, pp. 2759–2766.
- [19] Jinghua Wang, Zhenhua Wang, Dacheng Tao, Simon See, and Gang Wang, "Learning common and specific features for RGB-D semantic segmentation with deconvolutional networks," in *Proc. European Conference on Computer Vision*, 2016, pp. 664–679.
- [20] Xiaokang Chen, Kwan-Yee Lin, Jingbo Wang, Wayne Wu, Chen Qian, Hongsheng Li, and Gang Zeng, "Bi-directional cross-modality feature propagation with separation-and-aggregation gate for RGB-D semantic segmentation," in *Proc. European Conference on Computer Vision*, 2020.
- [21] Di Lin, Guangyong Chen, Daniel Cohen-Or, Pheng-Ann Heng, and Hui Huang, "Cascaded feature network for semantic segmentation of RGB-D images," in *Proc. IEEE International Conference on Computer Vision*, 2017, pp. 1320–1328.
- [22] Yuan Yuan, Zhiyu Jiang, and Qi Wang, "HDP: Hierarchical deep probability analysis for scene parsing," in *Proc. IEEE International Conference on Multimedia and Expo*, 2017, pp. 313–318.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [24] Ting Zhao and Xiangqian Wu, "Pyramid feature attention network for saliency detection," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3085–3094.
- [25] Shuran Song, Samuel P. Lichtenberg, and Jianxiong Xiao, "SUN RGB-D: A RGB-D scene understanding benchmark suite," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 567–576.
- [26] Xiaojuan Qi, Renjie Liao, Jiaya Jia, Sanja Fidler, and Raquel Urtasun, "3D graph neural networks for RGBD semantic segmentation," in *Proc. IEEE International Conference on Computer Vision*, 2017, pp. 5209–5218.
- [27] Shu Kong and Charless C. Fowlkes, "Recurrent scene parsing with perspective understanding in the loop," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 956–965.
- [28] Zhenyu Zhang, Zhen Cui, Chunyan Xu, Yan Yan, Nicu Sebe, and Jian Yang, "Pattern-affinitive propagation across depth, surface normal and semantic segmentation," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4106–4115.