

Abstract

Errors in semantic segmentation could be classified into two types: the large area misclassification and inaccurate local boundaries. Previously attention-based methods typically capture rich global contextual information, which benefits the large area classification but cannot address the local errors of boundaries. In this paper, we propose a Global-Local Attention Network (GLANet) which can simultaneously consider the global context and local details. Specifically, our GLANet consists of two branches: (1) the global attention branch and (2) local attention branch. Furthermore, three different modules are embedded in GLANet for respectively modelling the semantic interdependencies in spatial, channel and boundary dimension. Lastly, we merge the outputs of different branches to enhance the feature representation further, resulting in more precise segmentation. Overall, the proposed method achieves the competitive segmentation accuracy on two public aerial image datasets, bringing significant improvements over the existing baselines.

Contributions

- We analysis error composition of semantic segmentation task and propose Global-Local Attention Network (GLANet) to diminish both the two type of errors.
- Our GLANet give a simultaneous consideration of both global context and local details with the help of global attention branch and local attention branch, in which three different attention modules are embedded.
- We achieve new state-of-the-art results on two popular aerial image datasets including Vaihingen dataset and Postdam dataset.

Method

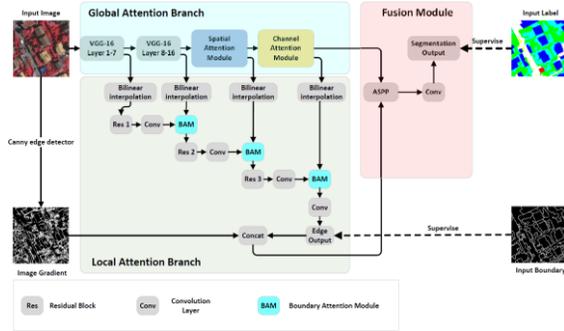


Fig.1. An overview of our Global-Local Attention Network (GLANet).

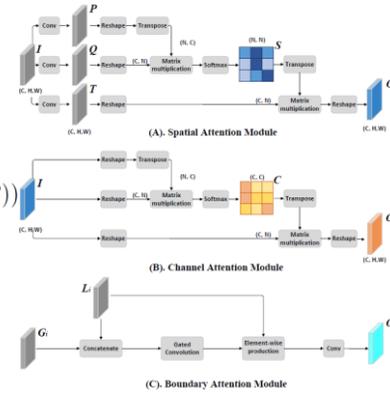
The total loss function:

$$L = \lambda_1 L_s + \lambda_2 L_b$$

$$L_s = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C y_{ij}^s \log p_{ij}^s$$

$$L_b = -\frac{1}{N} \sum_{i=1}^N (\beta y_i^b \log p_i^b + (1 - \beta)(1 - y_i^b) \log(1 - p_i^b))$$

Fig.2. The details of three attention module embedded in Global Attention Branch and Local Attention Branch.



Important References

- [1]. J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu, "Dual attention network for scene segmentation," in The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2019.
- [2]. T. Takikawa, D. Acuna, V. Jampani, and S. Fidler, "Gated-scnn: Gated shape cnns for semantic segmentation," in The IEEE International Conference on Computer Vision (ICCV), October 2019.
- [3]. L. Mou, Y. Hua, and X. X. Zhu, "A relation-augmented fully convolutional network for semantic segmentation in aerial scenes," in The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2019.

Acknowledgements

This work is supported by National Key R&D Program of China under contract No. 2017YFB1002203, NSFC projects under Grant 61976201, NSFC Key Projects of International (Regional) Cooperation and Exchanges under Grant 61860206004, and Ningbo 2025 Key Project of Science and Technology Innovation with No. 2018B10071.

Experiments

Method	SAM	CAM	BAM	Mean IoU%
Dilated FCN [7]				72.69
GLANet	✓			78.63
GLANet		✓		77.42
GLANet			✓	78.31
GLANet	✓	✓		79.75
GLANet	✓		✓	80.27
GLANet		✓	✓	80.09
GLANet	✓	✓	✓	80.67

Table 1. Ablation study about attention modules on the Vaihingen dataset.

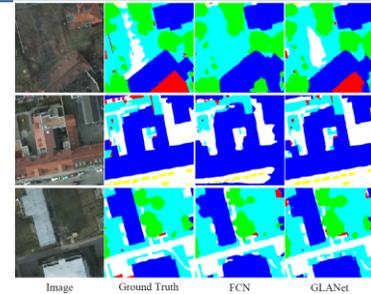


Fig.3. Examples of segmentation results on the Potsdam dataset.

Model Name	Imp. surf.	Build.	Low veg.	Tree	Car	mean F_1	mIoU	OA
Dilated FCN* [16]	86.52	90.78	83.01	78.41	90.42	85.83	-	84.14
SCNN [32]	88.37	92.32	83.68	80.94	91.17	87.30	77.72	85.57
FCN [7]	88.61	93.29	83.29	79.83	93.02	87.61	78.34	85.59
FCN-dCRF [25]	88.62	93.29	83.29	79.83	93.03	87.61	78.35	85.60
FCN-FR * [33]	89.31	94.37	84.83	81.10	93.56	88.63	-	87.02
S-RA-FCN [13]	91.33	94.70	86.81	83.47	94.52	90.17	82.38	88.59
GLANet	91.07	97.13	87.96	84.88	93.71	90.95	83.04	91.88

Table 2. Experimental Results on the Postdam Dataset.