

# BISVP: BUILDING FOOTPRINT EXTRACTION VIA BIDIRECTIONAL SERIALIZED VERTEX PREDICTION

Mingming Zhang<sup>1</sup>, Ye Du<sup>1</sup>, Zhenghui Hu<sup>2</sup>, Qingjie Liu<sup>\*1</sup>, Yunhong Wang<sup>1</sup>

<sup>1</sup> State Key Laboratory of Virtual Reality Technology and Systems, Beihang University, Beijing, China

<sup>2</sup> Hangzhou Innovation Institute, Beihang University

## ABSTRACT

Extracting building footprints from remote sensing images has been attracting extensive attention recently. Dominant approaches address this challenging problem by generating vectorized building masks with cumbersome refinement stages, which limits the application of such methods. In this paper, we introduce a new refinement-free and end-to-end building footprint extraction method, which is conceptually intuitive, simple, and effective. Our method, termed as BiSVP, represents a building instance with ordered vertices and formulates the building footprint extraction as predicting the serialized vertices directly in a bidirectional fashion. Moreover, we propose a cross-scale feature fusion (CSFF) module to facilitate high resolution and rich semantic feature learning, which is essential for the dense building vertex prediction task. Without bells and whistles, our BiSVP outperforms state-of-the-art methods by considerable margins on three building instance segmentation benchmarks, clearly demonstrating its superiority. The code and datasets will be made public available.

**Index Terms**— building footprint extraction, cross-scale feature fusion, bidirectional prediction, attention mechanism

## 1. INTRODUCTION

Extracting building footprints from remote sensing images has been receiving increasing attention due to its great potential value in many applications, such as urban change detection, city modeling, and cartography, which require precise geometric contours. Most prevalent methods address this task by vectorizing building segmentation masks, heavily relying on the performance of segmentation methods. Another line of works directly predict the order-agnostic vertex set, however, they require additional information to determine the vertex's order. In a nutshell, mainstream approaches tend to have complex model structures and tedious inference processes.

To tackle the above issues, we propose a Bidirectional Serialized Vertex Prediction (BiSVP) framework to predict the *ordered sequence* of building vertices. The method is refinement-free and can be trained end-to-endly. Considering that a polygon can be represented with sequential vertices in clockwise or counterclockwise direction, our BiSVP represents the building contour with ordered vertices and predicts building vertices sequentially in a bidirectional fashion. The

predictions from two directions are then combined to produce the final results. In this way, we manage to leverage the *bidirectional* information of building polygons to generate accurate building footprints. Besides, an attention mechanism is integrated into our BiSVP to enhance the ability of predicting long vertex sequences of complex buildings. Furthermore, we propose a cross-scale feature fusion (CSFF) module to obtain building features with high resolution and rich semantic information.

Our BiSVP can be seamlessly incorporated into existing object detectors (e.g., Faster RCNN [1]). Experimental results show that our method significantly outperforms state-of-the-art approaches on three building instance segmentation benchmarks. In short, the contributions are summarized as follows: First, we propose Bidirectional Serialized Vertex Prediction, a simple yet effective end-to-end framework to extract building footprints. Second, our method predicts the serialized vertices directly in a bidirectional fashion and proposes the cross-scale feature fusion module to enhance the building feature learning. Third, our method outperforms state-of-the-art approaches by considerable margins on three building instance segmentation benchmarks.

## 2. RELATED WORK

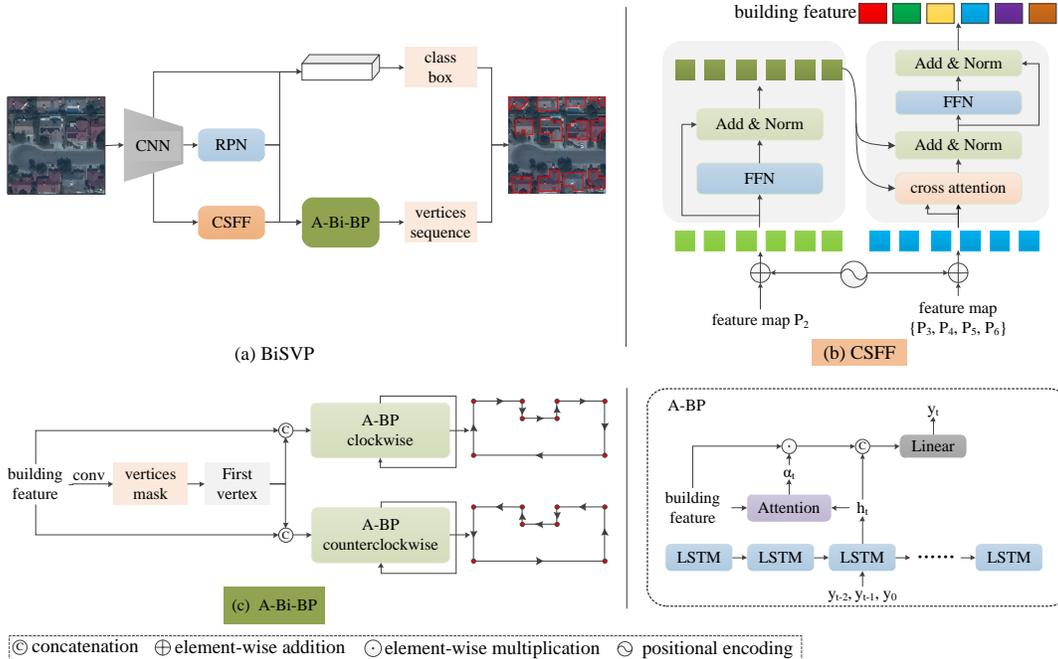
In this section, we review literature closely related to our research.

### 2.1. Building instance segmentation

Early works address building footprint extraction as a pixel-wise classification problem. To improve the segmentation quality, multi-source information, such as digital surface data and LIDAR data, are incorporated to obtain rich features [2, 3]. In the era of deep learning, numerous instance-level building segmentation methodologies have emerged [4, 5, 6, 7], which benefit a lot from instance segmentation networks. However, these methods provide building masks in raster format, which can not meet application needs.

### 2.2. Polygonal building segmentation

Polygonal building segmentation approaches extract building footprints in a vector format. Most current methods [8, 9,



**Fig. 1.** Overview of our proposed BiSVP. (a) BiSVP is a refinement-free and end-to-end framework. (b) Cross-scale feature fusion (CSFF) module is introduced to facilitate high resolution and rich semantic feature learning. (c) Attentional bidirectional building polygon (A-Bi-BP) prediction module is proposed to predict the serialized vertices in a bidirectional fashion directly.

10, 11, 12] vectorize building segmentation masks by post-processing. Girard et al. [13] predicts raster segmentation and frame fields to generate building polygons. These methods heavily rely on building segmentation masks and consequently generate irregular building polygons.

Another mainstream algorithms focus on directly predicting building vertices from the feature maps. Li et al. [14] and Zhu et al. [15] determine the building vertex’s order in a geometric way. PolyWorld [16] assigns the vertex connections for building vertices by solving a differentiable optimal transport problem. These methods decompose the polygonal building segmentation problem into multiple tasks and require complex polygon constraints, normally bringing intensive computation burden and leading to poor generalization. To determine the vertex’s order, PolyRNN [17] and PolyRNN++ [18] apply a CNN-RNN architecture to directly extract vertex sequences. Following this simple idea, many attempts [19, 20] have been made and achieved promising results. However, these methods are sensitive to inevitable occlusions and shadows, which may be challenging to extract complex building vertices.

### 3. METHOD

The overall architecture of BiSVP is shown in Figure 1, including a feature extraction network, a cross-scale feature fusion (CSFF) module, and an attentional bidirectional building

polygon (A-Bi-BP) prediction module.

#### 3.1. Feature Extraction Network

Our method can be seamlessly embedded into any polygonal models. To evaluate the effectiveness of it, we build a model on top of Mask RCNN because of its popularity and excellence performance. To be specific, BiSVP adopts a deep CNN backbone to extract multi-scale features  $F$ . To improve multi-scale building segmentation, we apply a feature pyramid network (FPN) [21] to fuse feature maps of different resolutions  $\{P_2, P_3, P_4, P_5, P_6\}$ . Then, a region proposal network (RPN) [1] proposes candidate building bounding boxes from the multi-scale features.

#### 3.2. Cross-scale Feature Fusion

As depicted in Figure 1 (b), CSFF adopts the transformer-based architecture to obtain the high resolution and rich semantic building representations, which are crucial for the subsequent vertex sequence prediction. Especially, CSFF primarily employs the cross attention module to aggregate feature maps of different resolutions in a coarse to fine manner, which can automatically focus on building boundaries.

Given an image  $I \in R^{3 \times H \times W}$ , CSFF firstly obtains the building queries  $B_q$  from the feature map  $P_2$  and the positional encoding. Then, CSFF takes in feature maps  $P_3, P_4, P_5, P_6$ , which embeds positional encoding with the corresponding feature map to localize building instances and

boundaries in the following cross attention. The cross attention module aggregates information between building queries  $B_q$  and feature maps  $P_i$  ( $i \in [3, 4, 5, 6]$ ). Besides, the residual connection and layer normalization are also applied after the attention and FFN modules. The output of the CSFF module is defined as:

$$B = \text{FFN}(\text{softmax}(\frac{B_q P_i^k}{\sqrt{d}}) \cdot P_i^v) \quad (1)$$

where  $B_q$  is the building queries from the feature map  $P_2$ , and FFN is the fully connected feed-forward network.  $P_i^k$  and  $P_i^v$  represent key and value from feature map  $P_i$  ( $i \in [3, 4, 5, 6]$ ) from FPN.

### 3.3. Attentional bidirectional building polygon

As illustrated in Figure 1 (c), A-Bi-BP module firstly takes in the building feature  $B$  to get the first vertex. Then, it outputs building vertices in two directions (i.e., clockwise and counterclockwise). In the following, we take step  $t$  as an example to introduce the attentional building polygon (A-BP) prediction, which is a branch of A-Bi-BP module.

Firstly, A-BP module outputs the hidden state  $h_t$  by taking in the building feature  $B$ , the previous predicted vertices  $y_{t-2}$  and  $y_{t-1}$ , and the predicted first vertex  $y_0$ . Then, a gaussian constrained attention [22] module integrating the hidden state  $h_t$  and the building feature  $B$  is applied to calculate the attention weight  $\alpha_t$ . Subsequently, the coefficient is calculated by the element-wise product of the attention weight  $\alpha_t$  and the building feature  $B$ . Finally, the next vertex  $y_t$  or the end signal (EOS) is captured from the concatenation of the coefficient and the hidden state  $h_t$ . The step  $t$  of the A-BP is defined as:

$$\begin{aligned} (h_t, c_t) &= \text{LSTM}(B, y_{t-2}, y_{t-1}, y_0), \\ y_t &= \text{softmax}(W[h_t; \alpha_t \odot B] + b) \end{aligned} \quad (2)$$

where  $\alpha_t$  is  $\text{attn}_{gc}(B, h_t)$ ,  $\alpha_t \odot B$  is the element-wise product of the attention weights and  $B$ , and  $[\cdot]$  is the concatenation operation.  $W$  and  $b$  are the trainable parameters and  $y_t$  is obtained by softmax operation.

### 3.4. Training objective

BiSVP loss includes binary classification loss  $L_{cls}$ , bounding box regression loss  $L_{reg}$ , and building vertex sequence prediction loss  $L_{ver}$ . This paper adopts the binary cross entropy loss and L1 loss to calculate  $L_{cls}$  and  $L_{reg}$ , respectively. As for  $L_{ver}$ , we calculate the average loss of the cross entropy loss between predicted polygons of two directions with the corresponding ground truth respectively. The total loss is defined as follows:

$$\begin{aligned} L_{ver} &= (\text{L}_{ce}(pred_p, gt) + \text{L}_{ce}(pred_p^c, gt^c))/2.0, \\ L &= L_{cls} + L_{reg} + L_{ver}. \end{aligned} \quad (3)$$

where  $pred_p$  and  $gt$  ( $pred_p^c$  and  $gt^c$ ) represent the predicted building polygon and ground truth in clockwise (counterclockwise).

## 4. EXPERIMENTS

### 4.1. Datasets

The proposed method is evaluated on three building datasets: (1) SpaceNet (LasVegas) [23] consists of over 3,800 images of size  $650 \times 650$  pixels. (2) 5M-Building [24] contains 109 images with a resolution ranging from  $2000 \times 2000$  to  $5000 \times 5000$ . We crop them into  $512 \times 512$  sub-images with an overlap of 64 pixels and then split it by 7:3 for training and testing. (3) CNData is a very challenging dataset, including 4200 images with a size of  $512 \times 512$ . The images are captured over different provinces of China and consist of residential, rural and industrial areas, where buildings vary greatly in size, structure and appearance. Especially in rural and urban villages, buildings are small and dense. The dataset has 101430 building instances with polygonal annotations, which is split by 8:1:1 for training, validation and testing.

### 4.2. Implementation Details

Our proposed model is trained in an end-to-end manner with SGD [25] optimizer. The backbone is a ResNet50 pre-trained on ImageNet, and we fine-tune it with a smaller initial learning rate  $1e-5$ ; for the other part of the model, the initial learning rate is set to  $1e-4$ . The weight decay is set to  $1e-4$ . The model is trained for 24 epochs, and we decrease the learning rate by 10 at the 16-th and 22-th epoch, respectively. We use MS COCO metrics [26] to evaluate the segmentation results. Furthermore,  $F1_{75}$  calculated from  $AP_{75}$  and  $AR_{75}$  is also employed to comprehensively evaluate different methods, since high accuracy delineation is vital for practical applications.

### 4.3. Comparison with State-of-the-arts

Since building footprint extraction is also a instance segmentation task, we compare it with the baseline model Mask R-CNN [27] and PANet [28]. Furthermore, we also compare the proposed BiSVP with PolyMapper [19] and the SOTA method Framefield [13] for polygonal building segmentation.

**Quantitative Evaluation.** Table 1 reports the building segmentation results on three building datasets. As illustrated in Table 1,  $F1_{75}$  of our method on the three datasets are significantly better than the baseline method by 13.34%, 2.93%, and 5.15%;  $AP_{75}$  are improved by 15.7%, 2.1%, and 4%, which indicates that our method can extract building footprint precisely. Besides, the recall metrics are all enhanced by our approach on three building test datasets, especially, +9.6% in terms of  $AR_{75}$  on SpaceNet. Moreover, our model outperforms the polygonal building segmentation methods by large margins.

**Qualitative Comparison.** Figure 2 shows some example results obtained by our approach. It can be seen that our method can generate high-quality polygonal building footprints.

Dataset	Method	AP	$AP_{50}$	$AP_{75}$	AR	$AR_{50}$	$AR_{75}$	$F1_{75}$
SpaceNet (LasVegas)	PANet [28]	46.9	85.1	45.9	54.6	87.8	60.9	52.35
	PolyMapper [19]	51.6	<b>87.4</b>	59.6	58.3	<b>89.5</b>	68.9	63.91
	FrameField [13]	<b>53.6</b>	84.5	<b>63.1</b>	58.5	87.9	68.8	65.83
	Baseline [27]	47.0	85.9	46.4	54.6	88.0	60.5	52.52
	BiSVP (ours)	53.2 <sub>+6.2</sub>	87.2 <sub>+1.3</sub>	62.1 <sub>+15.7</sub>	<b>59.3</b> <sub>+4.7</sub>	<b>89.5</b> <sub>+1.5</sub>	<b>70.1</b> <sub>+9.6</sub>	<b>65.86</b> <sub>+13.34</sub>
5M-Building	PANet [28]	31.3	59.6	28.9	45.8	74.7	48.5	36.22
	PolyMapper [19]	32.0	62.9	30.1	47.5	82.1	51.1	37.88
	FrameField [13]	18.4	36.4	16.4	31.0	54.7	31.0	21.45
	Baseline [27]	31.5	60.5	29.0	45.9	75.4	48.3	36.24
	BiSVP (ours)	<b>32.7</b> <sub>+1.2</sub>	<b>63.8</b> <sub>+3.3</sub>	<b>31.1</b> <sub>+2.1</sub>	<b>48.7</b> <sub>+2.8</sub>	<b>83.9</b> <sub>+8.5</sub>	<b>52.9</b> <sub>+4.6</sub>	<b>39.17</b> <sub>+2.93</sub>
CNData	PANet [28]	35.1	68.8	34.0	47.5	81.3	50.3	40.57
	PolyMapper [19]	36.4	70.6	35.7	50.6	86.1	54.6	43.17
	FrameField [13]	21.7	40.7	21.2	32.9	54.9	34.4	26.23
	Baseline [27]	35.1	68.4	33.7	47.7	81.6	50.2	40.33
	BiSVP (ours)	<b>37.9</b> <sub>+2.8</sub>	<b>71.5</b> <sub>+3.1</sub>	<b>37.7</b> <sub>+4.0</sub>	<b>52.7</b> <sub>+5.0</sub>	<b>88.3</b> <sub>+6.7</sub>	<b>57.3</b> <sub>+7.1</sub>	<b>45.48</b> <sub>+5.15</sub>

**Table 1.** Results on three building test datasets: SpaceNet (LasVegas), 5M-Building, and CNData. The best results in each dataset group are marked in bold.



**Fig. 2.** Qualitative results. Our method can generate geometric contours of buildings accurately.

#### 4.4. Ablation Study

We analyze the effectiveness of CSFF module, Bi-BP module, and the attention mechanism ( $Attn_{gc}$ ) in Bi-BP. In ablation studies, we add CSFF, Bi-BP, and  $Attn_{gc}$  respectively to the baseline method [27]. The experimental results are present in Table 2.

**Cross-scale feature fusion (CSFF).** The performance of the baseline model decreases on the three building datasets by 12.27%, 2.26%, and 4.73% in the indicator  $F1_{75}$ . The results shown in Table 2 indicate that CSFF module plays a vital role in aggregating features of different levels.

**Bidirectional building polygon (Bi-BP).** Table 2 shows that Bi-BP module can significantly improve the performance on the three building datasets, which demonstrates that Bi-BP module can leverage the *bidirectional* information of the building polygon. Especially, it improves the baseline by 12.8% in  $F1_{75}$  of SpaceNet (Las Vega), proving the effectiveness of the Bi-BP module.

**Attention mechanism.** We can observe from Table 2 that the model with the attention mechanism can achieve surprisingly good performance. Finally, the proposed BiSVP can improve the baseline from 52.52%, 36.24%, and 40.33% to 65.86%

Dataset	CSFF	Bi-BP	Attn	$F1_{75}$
SpaceNet (LasVegas)	✓			52.52
		✓		64.79
			✓	65.32
	✓	✓	✓	65.44
			<b>65.86</b>	
5M-Building	✓			36.24
		✓		38.50
			✓	38.58
	✓	✓	✓	39.03
			<b>39.17</b>	
CNData	✓			40.33
		✓		45.06
			✓	43.35
	✓	✓	✓	44.80
			<b>45.48</b>	

**Table 2.** Ablation study. "✓" means adding the corresponding module to the baseline. The last row in each dataset group is the value of BiSVP. The best result in each dataset group is marked in bold.

(+13.34%), 39.17% (+2.93%), and 45.48% (+5.15%) in  $F1_{75}$  over three building test datasets, respectively.

## 5. CONCLUSION

In this paper, we have presented Bidirectional Serialized Vertex Prediction (BiSVP), a new refinement-free and end-to-end framework to extract building footprints from remote sensing images. The proposed BiSVP represents the building contour with ordered vertices and predicts the serialized vertices directly in a bidirectional fashion. Furthermore, BiSVP proposes a cross-scale feature fusion (CSFF) module to fuse feature maps of different levels, obtaining the building feature with rich spatial and context information that is essential for the dense building vertex prediction task. The extensive experiments on three building instance segmentation datasets demonstrate the superiority of our method in building footprint extraction.

## 6. REFERENCES

- [1] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *NIPS*, 2015.
- [2] Mohammad Awrangjeb, Mehdi Ravanbakhsh, and Clive S. Fraser, “Automatic detection of residential buildings using lidar data and multispectral imagery,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 65, pp. 457–467, Sep. 2010.
- [3] Weijia Li, Conghui He, Jiarui Fang, and Haohuan Fu, “Semantic segmentation based building extraction method using multi-source gis map datasets and satellite imagery,” in *CVPRW*, 2018.
- [4] Qingpeng Li, Yunhong Wang, Qingjie Liu, and Wei Wang, “Hough transform guided deep feature extraction for dense building detection in remote sensing images,” in *ICASSP*, 2018.
- [5] Kang Zhao, Jungwon Kang, Jaewook Jung, and Gunho Sohn, “Building extraction from satellite images using mask r-cnn with building boundary regularization,” in *CVPRW*, 2018.
- [6] Stefano Zorzi and Friedrich Fraundorfer, “Regularization of building boundaries in satellite images using adversarial and regularized losses,” in *IGARSS*, 2019.
- [7] Lele Xu, Ye Li, Jinzhong Xu, and Lili Guo, “Gated spatial memory and centroid-aware network for building instance extraction,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–14, Apr 2021.
- [8] Qi Chen, Lei Wang, Steven L. Waslander, and Xiuguo Liu, “An end-to-end shape modeling framework for vectorized building outline generation from aerial images,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 170, pp. 114–126, Aug. 2020.
- [9] Justin Liang, Namdar Homayounfar, Wei-Chiu Ma, Yuwen Xiong, Rui Hu, and Raquel Urtasun, “Polytransform: Deep polygon transformer for instance segmentation,” in *CVPR*, 2020.
- [10] Muxingzi Li, Florent Lafarge, and Renaud Marlet, “Approximating shapes in images with low-complexity polygons,” in *CVPR*, 2020.
- [11] Stefano Zorzi, Ksenia Bittner, and Friedrich Fraundorfer, “Machine-learned regularization and polygonization of building segmentation masks,” in *ICPR*, 2021.
- [12] Yuhao Chen, Yifan Wu, Linlin Xu, and Alexander Wong, “Quantization in relative gradient angle domain for building polygon estimation,” in *ICPR*, 2021.
- [13] Nicolas Girard, Dmitriy Smirnov, Justin Solomon, and Yuliya Tarabalka, “Polygonal building extraction by frame field learning,” in *CVPR*, 2021.
- [14] Qingyu Li, Lichao Mou, Yuansheng Hua, Yao Sun, Pu Jin, Yilei Shi, and XiaoXiang Zhu, “Instance segmentation of buildings using keypoints,” in *IGARSS*, 2020.
- [15] Y. Zhu, B. Huang, J. Gao, E. Huang, and H. Chen, “Adaptive polygon generation algorithm for automatic building extraction,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–14, May 2022.
- [16] Stefano Zorzi, Shabab Bazrafkan, Stefan Habenschuss, and Friedrich Fraundorfer, “Polyworld: Polygonal building extraction with graph neural networks in satellite images,” in *CVPR*, 2022.
- [17] Lluís Castrejón, Kaustav Kundu, Raquel Urtasun, and Sanja Fidler, “Annotating object instances with a polygon-rnn,” in *CVPR*, 2017.
- [18] David Acuna, Huan Ling, Amlan Kar, and Sanja Fidler, “Efficient interactive annotation of segmentation datasets with polygon-rnn++,” in *CVPR*, 2018.
- [19] Zuoyue Li, Jan Dirk Wegner, and Aurélien Lucchi, “Topological map extraction from overhead images,” in *ICCV*, 2019.
- [20] Wufan Zhao, Claudio Persello, and Alfred Stein, “Building outline delineation: From aerial images to polygons with an improved end-to-end learning framework,” *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 175, pp. 119–131, May 2021.
- [21] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie, “Feature pyramid networks for object detection,” in *CVPR*, 2017.
- [22] Zhi Qiao, Xugong Qin, Yu Zhou, Fei Yang, and Weiping Wang, “Gaussian constrained attention network for scene text recognition,” in *ICPR*, 2021.
- [23] Adam Van Etten, Dave Lindenbaum, and Todd M. Bacastow, “Spacenet: A remote sensing dataset and challenge series,” *arXiv preprint arXiv:1807.01232*, 2018.
- [24] Zeshan Lu, Tao Xu, Kun Liu, Zhen Liu, Feipeng Zhou, and Qingjie Liu, “5m-building: A large-scale high-resolution building dataset with cnn based detection analysis,” in *ICTAI*, 2019.
- [25] Herbert Robbins and Sutton Monro, “A stochastic approximation method,” *The annals of mathematical statistics*, vol. 22, pp. 400–407, Sep. 1951.
- [26] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár, “Microsoft coco: Common objects in context,” in *ECCV*, 2014.
- [27] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick, “Mask r-cnn,” in *ICCV*, 2017.
- [28] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia, “Path aggregation network for instance segmentation,” in *CVPR*, 2018.