# Dense Dilated Convolutions Merging Network for Semantic Mapping of Remote Sensing Images

1st Qinghui liu
*SAMBA and Machine Learning Group*
*Norwegian Computing Center and UiT*
Oslo, Norway
Brian.Liu@nr.no

2nd Michael Kampffmeyer
*Machine Learning Group*
*UiT, the Arctic University of Norway*
Tromsø, Norway
Michael.C.Kampffmeyer@uit.no

3nd Robert Jenssen
*Machine Learning Group*
*UiT, the Arctic University of Norway*
Tromsø, Norway
Robert.Jenssen@uit.no

4th Arnt-Børre Salberg
*dept. SAMBA of NR*
*Norwegian Computing Center*
Oslo, Norway
Arnt-Borre.Salberg@nr.no

*Abstract*—We propose a network for semantic mapping called the Dense Dilated Convolutions Merging Network (DDCM-Net) to provide a deep learning approach that can recognize multi-scale and complex shaped objects with similar color and textures, such as buildings, surfaces/roads, and trees in very high resolution remote sensing images. The proposed DDCM-Net consists of dense dilated convolutions merged with varying dilation rates. This can effectively enlarge the kernels' receptive fields, and, more importantly, obtain fused local and global context information to promote surrounding discriminative capability. We demonstrate the effectiveness of the proposed DDCM-Net on the publicly available ISPRS Potsdam dataset and achieve a performance of 92.3% F1-score and 86.0% mean intersection over union accuracy by only using the RGB bands, without any post-processing. We also show results on the ISPRS Vaihingen dataset, where the DDCM-Net trained with IRRG bands, also obtained better mapping accuracy (89.8% F1-score) than previous state-of-the-art approaches.

*Index Terms*—Dense Dilated Convolutions Merging (DDCM), deep learning, semantic mapping, remote sensing

## I. INTRODUCTION

Automatic semantic interpretation of remote sensing images is important for a wide range of practical applications, such as urban land cover classification [1], traffic monitoring and vehicle detection. Large-scale semantic mapping is a challenging task which consists of the assignment of a semantic category to every pixel in very high resolution (VHR) aerial images. Due to the successes of deep learning methods, a large variety of modern approaches to pixel-to-pixel classification are based on deep convolutional neural networks (CNN), in particular end-to-end learning with fully convolutional neural networks (FCN) [2]. However, to achieve higher performance, CNN and FCN based methods [3]–[5] normally rely on deep multi-scale architectures which typically require a large number of trainable parameters and computation resources.
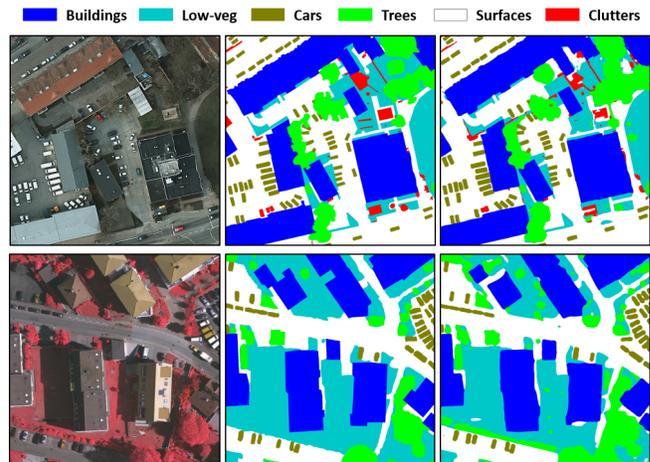
Fig. 1. Examples of semantic mapping of remote sensing images on RGB (top-left) and IRRG data (bottom-left) respectively with our DDCM network. From left to right, test images, ground truths, and mapping results.

In this work, we propose a novel network architecture, called the dense dilated convolutions merging network (DDCM-Net), which utilizes multiple dilated convolutions merged with various dilation rates. The proposed network learns with densely linked dilated convolutions and outputs a fusion of all intermediate features without losing resolutions during the extraction of multi-scale features. Our experiments demonstrate that the network achieves robust and accurate results with relatively few parameters and layers. Fig. 1 shows illustrative examples of semantic mapping results on RGB and IRRG data respectively with our DDCM-Net methods. These results will be further discussed in section III.

## II. METHODS

We first briefly revisit dilated convolutions which are used in DDCM networks. We then present our proposed DDCM

architecture and provide training details.

## A. Dilated Convolutions

Dilated convolutions [6] have been demonstrated to improve performance in many classification and segmentation tasks [7]–[10]. One key advantage is that they allow us to flexibly adjust the filter's receptive field to capture multi-scale information without resorting to down-scaling and up-scaling operations. A 2-D dilated convolution operator can be defined as

$$g_{i,j}(x_\ell) = \sum_{c=0}^{C_\ell} \theta_{k,r}^{i,j} * x_\ell^c \qquad (1)$$

where, $*$ denotes a convolution operator, $g_{i,j} : \mathbb{R}^{H_\ell \times W_\ell \times C_\ell} \rightarrow \mathbb{R}^{H_{\ell+1} \times W_{\ell+l}}$ convolves the input feature map $x_\ell \in \mathbb{R}^{H_\ell \times W_\ell \times C_\ell}$ within channel $c \in \{0, 1, \ldots, C_\ell\}$ at row $i$ and column $j$. A dilated convolution $\theta_{k,r}$ with a filter $k$ and dilation $r \in \mathbb{Z}^+$ is only nonzero for a multiple of $r$ pixels from the center. In dilated convolution, a kernel size $k$ is enlarged to $k + (k-1)(r-1)$ with the dilation factor $r$. As a special case, a dilated convolution with dilation rate $r = 1$ corresponds to a standard convolution.

## B. Dense Dilated Convolutions Merging Module

The proposed dense dilated convolutions merging (DDCM) module densely stacks multi-scale features and merges them to yield more accurate and robust representations with fewer parameters. Fig. 2 illustrates the basic structure of the DDCM module.

DDCM module consists of a number of Dilated CNN-stack (DCs) blocks with a merging module as output. A basic DCs block is composed of a dilated convolution followed by PReLU [11] non-linear activation and batch normalization (BN) [12]. It then stacks the output with its input together to feed the next layer, which can alleviate context information loss and problems with vanishing gradients when adding more layers. The final network output is computed by a merging layer composed of $1 \times 1$ filters with BN and PReLU in order to efficiently combine all stacked features generated by intermediate DCs blocks.

In a DDCM module, all feature maps are maximally utilized with high computational efficiency while preserving the input resolution throughout the network. In particular, densely connected DCs blocks, typically configured with linearly increasing dilation factors, enable DDCM networks to have very large receptive fields with just a few layers as well as to capture rich global representations by merging multi-scale features properly.

Fig. 3 shows the end-to-end pipeline of DDCM network (DDCM-Net) architecture combined with a pre-trained ResNet [13] for semantic mapping tasks. The proposed DDCM-Net is easy to implement, train and combine with existing architectures. In our work, we only utilize the first 3 bottleneck layers of ResNet50 and remove the last bottleneck layer and fully connected layers to reduce the number of parameters to train.
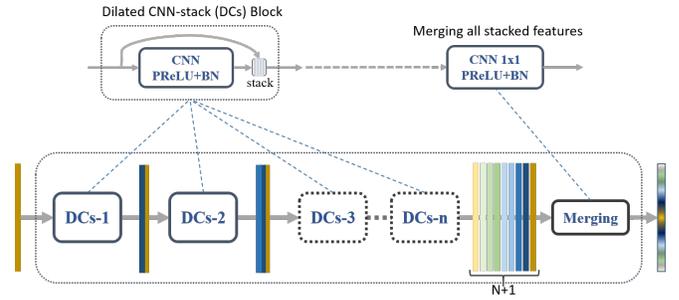


Fig. 2. Example of the DDCM architecture composed of $n$ DC blocks with various dilation rates $\{1, 2, 3, ..., n\}$.
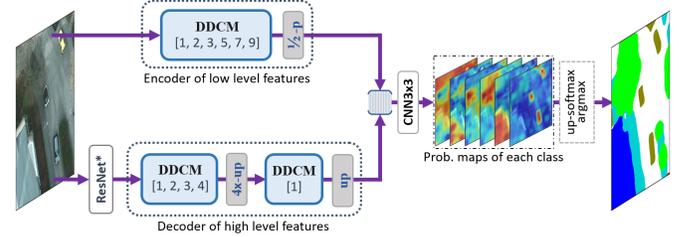


Fig. 3. End-to.end pipeline of DDCM-Net for semantic mapping of VHR images. The encoder of low level features encodes multi-scale contextual information from the initial input images by a DDCM module (output 3-channel) using $3 \times 3$ kernels with 6 different dilation rates $[1, 2, 3, 5, 7, 9]$. The decoder of high level features decodes highly abstract representations learned from ResNet50 (output 1024-channel) by 2 DDCM modules with rates $[1, 2, 3, 4]$ (output 36-channel) and $[1]$ (output 18-channel) separately. The transformed low-level and high-level feature maps by DDCMs are then fused together to infer pixel-wise class probabilities. Here, 'p' and 'up' denote pooling and up-sampling respectively.

## C. Data augmentation and normalization

We randomly sample 5000 image patches of size $256 \times 256$ in run time from VHR training images (of size $6000 \times 6000$) for each training epoch and flip and mirror images for data augmentation. These patches are normalized to [0.0, 1.0] by dividing by 255 for all bands (RGB, or IRRG). We used a pre-trained ResNet50 model that has been trained on ImageNet [14]. No mean and standard deviation normalization were used.

## D. Optimizer and weighted loss function

In our work, we choose Adam [15] with AMSGrad [16] as the optimizers with weight decay $5 \times 10^5$ and polynomial learning rate (LR) decay $(1 - \frac{cur\_iter}{max\_iter})^{0.9}$ with the maximum iterations of $10^8$ for the model. We also set $2 \times LR$ to all bias parameters in contrast to weights parameters. Guided by our empirical results, we use an initial LR of $\frac{8.5 \times 10^{-5}}{\sqrt{2}}$ and step-LR schedule method which drops the LR by 0.85 at every 15 epochs. As loss function, we apply a cross-entropy loss function with median frequency balancing as defined in [1].

## III. Experiments

We investigate the proposed network on the ISPRS 2D semantic labeling dataset [17] which is comprised of very high resolution aerial images over two cities: Potsdam and

Vaihingen in Germany. In this work, we only use RGB bands of Potsdam dataset and IRRG (Infrared-Red-Green) bands of Vaihingen dataset.

## A. Dataset

The ISPRS Potsdam dataset contains 38 RGB images ($6000 \times 6000$) annotated with six different labels including impervious surfaces, buildings, trees, low vegetation, cars and clutters. Originally, 24 of the images were public available and 14 were included in a hold-out test set. The Vaihingen dataset has 33 IRRG images, where 16 were from the original public dataset, and 17 were included the hold-out test set. To evaluate our models, the original public part (24 images) of the Potsdam dataset was divided into training, validation (areas: 4_10 and 7_10), and local test set (areas: 5_11, 6_9 and 7_11). The the original public part (16 images) of the Vaihingen dataset was similarly split into training, validation (tiles of 7 and 28), and local test set (tiles of 5, 15, 21, and 30). Please note that our trained models are evaluated both on the local test sets to compare with our previous work [1], [18], and on the hold-out test sets to compare with other related published work.

## B. Evaluation methods

We train and validate all networks with patches of size $256 \times 256$ as input and batch size of 5. All hyper-parameters settings, except the learning rates, were shared for the different models. At test time, we apply test time augmentation (TTA) in terms of flipping and mirroring. We use sliding windows (with $448 \times 448$ size at a 100px stride) on a test image and stitch the results together by averaging the predictions of the over-lapping TTA regions to form the output. The performance is measured by both the F1-score [1], and mean Intersection over Union (IoU) [18].

## C. Results

Table I shows our results on the hold-out test sets and our local test sets of ISPRS Potsdam and Vaihingen separately with a single trained model. The mean F1-score (mF1) and the mean IoU (mIou) are computed as the average measure of all classes except the clutter class. Fig. 4 shows a qualitative comparison of the semantic mapping results from our model and the ground truths.

TABLE I
RESULTS ON THE HOLD-OUT TEST IMAGES OF ISPRS POTSDAM AND VAIHINGEN DATASETS WITH A SINGLE TRAINED DDCM-R50 MODEL SEPARATELY.

| Potsdam | Avg. | Building | Tree | Low-veg | Surface | Car | OA |
|---|---|---|---|---|---|---|---|
| F1-score | 0.923 | 0.969 | 0.894 | 0.877 | 0.929 | 0.949 | 0.908 |
| | 0.925* | 0.983* | 0.892* | 0.865* | 0.946* | 0.939* | 0.931* |
| IoU | 0.860 | 0.940 | 0.809 | 0.781 | 0.867 | 0.902 | |
| | 0.863* | 0.966* | 0.805* | 0.762* | 0.898* | 0.885* | |
| **Vaihingen** | | | | | | | |
| F1-score | 0.898 | 0.953 | 0.894 | 0.833 | 0.927 | 0.883 | 0.904 |
| | 0.909* | 0.973* | 0.914* | 0.814* | 0.934* | 0.909* | 0.921* |
| IoU | 0.817 | 0.909 | 0.808 | 0.713 | 0.863 | 0.790 | |
| | 0.837* | 0.948* | 0.842* | 0.686* | 0.876* | 0.832* | |

\* Results were measured on our local test images.

We also compare our results to other related published work on the ISPRS Potsdam RGB dataset and Vaihingen IRRG



Fig. 4. Mapping results for test images of Vaihingen tile-27 (left) and Potsdam tile-3_14 (right). A. Test images (top), B. the Ground truths (center), C. DDCM-R50 (bottom).

dataset. These results are shown in Table II and III respectively. Our single model with overall F1-score (92.3%) on Potsdam RGB dataset, achieves around 0.5 percent higher than the secondary best model - FuseNet+OSM [19] which used Open-StreetMap (OSM) as an additional data source. In other words, our model achieves better performance with fewer labeled training data. Similarly, our model trained on Vaihingen IRRG images, also obtained the best overall performance with 89.8% F1-score which is around 1.1% higher than the second best model GSN [20]. It's worth noting that our model is the only one that works equally well on both Vaihingen IRRG dataset and Potsdam RGB dataset, which outperforms the DST_2 [21] model with 3.9% and 0.6% higher F1-score on Vaihingen and Potsdam dataset respectively.

TABLE II
COMPARISONS BETWEEN OUR METHOD WITH OTHER PUBLISHED METHODS ON THE HOLD-OUT RGB TEST IMAGES OF ISPRS POTSDAM DATASET.

| Models | OA | Building | Tree | Low-veg | Surface | Car | mF1 |
|---|---|---|---|---|---|---|---|
| HED+SEG.H-Sc1 [22] | 0.851 | 0.967 | 0.686 | 0.842 | 0.850 | 0.858 | 0.846 |
| RiFCN [23] | 0.883 | 0.930 | 0.819 | 0.837 | 0.917 | 0.937 | 0.861 |
| RGB+I-ensemble [24] | 0.900 | 0.936 | 0.845 | 0.822 | 0.870 | 0.892 | 0.873 |
| Hallucination [24] | 0.901 | 0.938 | 0.848 | 0.821 | 0.873 | 0.882 | 0.872 |
| SegNet RGB [19] | 0.897 | 0.929 | 0.851 | 0.850 | 0.930 | 0.951 | 0.902 |
| DST_2 [21] | 0.903 | 0.964 | 0.880 | 0.867 | 0.925 | 0.947 | 0.917 |
| FuseNet+OSM [19] | **0.923** | 0.959 | 0.851 | 0.863 | **0.953** | **0.968** | 0.918 |
| DDCM-R50 (ours) | 0.908 | **0.969** | **0.894** | **0.877** | 0.929 | 0.949 | **0.923** |
| | (-1.5%) | (+0.5%) | (+1.4%) | (+1.0%) | (-2.4%) | (-1.9%) | (+0.5%) |

In addition, we evaluate our method on the local Potsdam test set to compare with other popular architectures reviewed and re-implemented in [18]. Our DDCM-R50 model achieved the highest mIoU score (80.8%) compared to others while using much fewer parameters and computational cost (FLOPs) as shown in Table IV. Note that the performance on full reference ground truths is slightly lower than on eroded boundary ground truths as the boundary pixels are not ignored during evaluation.

TABLE III
COMPARISONS BETWEEN OUR METHOD WITH OTHER PUBLISHED
METHODS ON THE HOLD-OUT IRRG TEST IMAGES OF ISPRS VAIHINGEN
DATASET.

| Models | OA | Building | Tree | Low-veg | Surface | Car | mF1 |
|---|---|---|---|---|---|---|---|
| UOA [25] | 0.876 | 0.921 | 0.882 | 0.804 | 0.898 | 0.820 | 0.865 |
| ADL_3 [26] | 0.880 | 0.932 | 0.882 | 0.823 | 0.895 | 0.633 | 0.833 |
| DST_2 [21] | 0.891 | 0.937 | 0.892 | 0.834 | 0.905 | 0.726 | 0.859 |
| ONE_7 [27] | 0.898 | 0.945 | 0.899 | 0.844 | 0.910 | 0.778 | 0.875 |
| DLR_9 [22] | 0.903 | 0.952 | 0.899 | 0.839 | 0.924 | 0.812 | 0.885 |
| GSN [20] | 0.903 | 0.951 | 0.899 | 0.837 | 0.922 | 0.824 | 0.887 |
| DDCM-R50 (ours) | 0.904 | 0.953 | 0.894 | 0.833 | 0.927 | 0.883 | 0.898 |
| | (+0.1%) | (+0.1%) | (-0.5%) | (-1.1%) | (+0.3%) | (+5.9%) | (+1.1%) |

TABLE IV
QUANTITATIVE COMPARISON OF PARAMETERS SIZE, FLOPs (MEASURED
ON INPUT IMAGE SIZE OF $3 \times 256 \times 256$), AND mIoU ON ISPRS
POTSDAM RGB DATASET.

| Models | Backbones | Parameters (Million) | FLOPs (Giga) | mIoU* |
|---|---|---|---|---|
| UNet [28] | VGG16 | 31.04 | 15.25 | 0.715 |
| FCN8s [2] | VGG16 | 134.30 | 73.46 | 0.728 |
| SegNet [29] | VGG19 | 39.79 | 60.88 | 0.781 |
| GCN [5] | ResNet50 | 23.84 | 5.61 | 0.774 |
| PSPNet [3] | ResNet50 | 46.59 | 44.40 | 0.789 |
| DUC [4] | ResNet50 | 30.59 | 32.26 | 0.793 |
| DDCM-R50 (ours) | ResNet50 | 9.99 | 4.86 | 0.808 |

* mIoU was measured on full reference ground truths of our local test images 5_11, 6_9 and
7_11 in order to fairly compare with our previous work [18].

## IV. CONCLUSIONS

In this paper, we presented a dense dilated convolutions merging (DDCM) network architecture for semantic mapping in very high-resolution aerial images. The proposed architecture applies dilated convolutions to learn features at varying dilation rates, and merges the feature map of each layer with the feature maps from all previous layers. On both the Potsdam and Vahingen datasets, the DDCM-Net architecture achieves the best mean $F_1$ score compared to the other architectures, but with much fewer parameters and feature maps. DDCM-Net is easy to adapt to address a wide range of different problems by using various combinations of dilation rates, is fast to train, and achieves accurate results even on small datasets.

## REFERENCES

[1] M. Kampffmeyer, A.-B. Salberg, and R. Jenssen, "Semantic segmentation of small objects and modeling of uncertainty in urban remote sensing images using deep convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 1–9, 2016.

[2] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3431–3440, 2015.

[3] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," *arXiv preprint arXiv:1612.01105*, 2016.

[4] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, and G. Cottrell, "Understanding convolution for semantic segmentation," *arXiv preprint arXiv:1702.08502*, 2017.

[5] C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun, "Large kernel matters–improve semantic segmentation by global convolutional network," *arXiv preprint arXiv:1703.02719*, 2017.

[6] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.

[7] 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 transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2018.

[8] D. M. Pelt and J. A. Sethian, "A mixed-scale dense convolutional neural network for image analysis," *Proceedings of the National Academy of Sciences*, vol. 115, no. 2, pp. 254–259, 2018.

[9] Y. Li, X. Zhang, and D. Chen, "Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1091–1100, 2018.

[10] Y. Wei, H. Xiao, H. Shi, Z. Jie, J. Feng, and T. S. Huang, "Revisiting dilated convolution: A simple approach for weakly-and semi-supervised semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7268–7277, 2018.

[11] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.

[12] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.

[13] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

[14] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.

[15] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014.

[16] S. J. Reddi, S. Kale, and S. Kumar, "On the convergence of adam and beyond," *CoRR*, 2018.

[17] I. S. for Photogrammetry and R. S. (ISPRS), "2D Semantic Labeling Contest." online, 2018.

[18] Q. Liu, A. Salberg, and R. Jenssen, "A comparison of deep learning architectures for semantic mapping of very high resolution images," in *IGARSS 2018 - 2018 IEEE International Geoscience and Remote Sensing Symposium*, pp. 6943–6946, July 2018.

[19] N. Audebert, B. Le Saux, and S. Lefèvre, "Joint learning from earth observation and openstreetmap data to get faster better semantic maps," in *EARTHVISION 2017 IEEE/ISPRS CVPR Workshop. Large Scale Computer Vision for Remote Sensing Imagery*, 2017.

[20] H. Wang, Y. Wang, Q. Zhang, S. Xiang, and C. Pan, "Gated convolutional neural network for semantic segmentation in high-resolution images," *Remote Sensing*, vol. 9, no. 5, p. 446, 2017.

[21] J. Sherrah, "Fully convolutional networks for dense semantic labelling of high-resolution aerial imagery," *CoRR*, vol. abs/1606.02585, 2016.

[22] D. Marmanis, K. Schindler, J. D. Wegner, S. Galliani, M. Datcu, and U. Stilla, "Classification with an edge: Improving semantic image segmentation with boundary detection," *CoRR*, vol. abs/1612.01337, 2016.

[23] L. Mou and X. X. Zhu, "Rifcn: Recurrent network in fully convolutional network for semantic segmentation of high resolution remote sensing images," *CoRR*, vol. abs/1805.02091, 2018.

[24] M. Kampffmeyer, A.-B. Salberg, and R. Jenssen, "Urban land cover classification with missing data modalities using deep convolutional neural networks," *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 11, no. 6, pp. 1758–1768, 2018.

[25] G. Lin, C. Shen, A. Van Den Hengel, and I. Reid, "Efficient piecewise training of deep structured models for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3194–3203, 2016.

[26] S. Paisitkriangkrai, J. Sherrah, P. Janney, V.-D. Hengel, *et al.*, "Effective semantic pixel labelling with convolutional networks and conditional random fields," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 36–43, 2015.

[27] N. Audebert, B. Le Saux, and S. Lefèvre, "Semantic segmentation of earth observation data using multimodal and multi-scale deep networks," in *Asian Conference on Computer Vision*, pp. 180–196, Springer, 2016.

[28] 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*, pp. 234–241, Springer, 2015.

[29] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *arXiv preprint arXiv:1511.00561*, 2015.