

000 001 Building Rooftop Extraction in Aerial Image 002 with Dual-Task Fully Convolutional Network

003
004 Anonymous ACCV 2016 submission
005

006 Paper ID ***
007
008

009 **Abstract.** Automatic building extraction from remote sensing image
010 plays a critical role in a diverse range of applications. However, it is a huge
011 challenge to extract any-size buildings with large variational appear-
012 ances. In this article, we propose a novel multi-scale fully convolutional
013 network(MSFCN), as the name suggests, which can handle multi-size
014 objects only in a single convolutional neural network. Fully convolutional
015 networks(FCNs) are proved to be a appropriate method to complete
016 semantic segmentation task. It can take whole images as inputs without
017 warping or cropping and output segmentation results directly. However,
018 it has less ability to process a large number of small buildings in aerial
019 image. For the sake of handling this problem, we integrate all predictions
020 from layers with skipped receptive fields stage by stage. The experiment
021 results test on a publicly available aerial imagery dataset proved that our
022 proposed methodology significantly shorten time-consuming and surpass
023 the performance of state-of-the-art.
024

025 1 Introduction

026 With the rapid development of remote sensing technologies and popularization
027 of geospatial related commercial software, very high resolution satellite images
028 are very easily accessible, these valuable data provides a huge fuel for interpret-
029 ing real terrestrial scenes. Building rooftops is one the most important part of
030 terrestrial objects because it is essential for a wide range of technologies, such as,
031 urban planning, automated map making, 3D city modelling, disaster assessment,
032 military reconnaissance, etc. However, it is very costly and time-consuming for
033 human experts to complete this task. Therefore, many researchers have made
034 massive attempts to extract buildings automatically.
035

036 In previous literatures, large amounts of efforts achieve good performance
037 in extracting buildings with special color, shapes, textures, or surroundings.
038 According to the information they concerned, these methods may fall into three
039 classes. Though these techniques are different, the common point is there are a
040 series of manually set rules or features in their systems.

041 One popular way of extracting buildings is employing their shape informa-
042 tion. It is observed that rooftops have more regular shapes, which usually are
043 rectangular or combinations of several rectangles. A dozen years ago, Noronha
044 and Nevatia [1] designed a system that detects and constructs 3D models for rec-
tilinear buildings from multiple aerial images. Firstly, hypotheses for rectangular

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

2 ACCV-16 submission ID ***

045 rooftop were generated by grouping lines, then they were verified by searching
 046 for presence of predicted walls and shadows. Nosrati and Saeedi [2] pointed out
 047 that polygonal rooftops correspond to closed loops in a graph which represents
 048 the relationship between intersections of a pair of edge in an efficient way. In [3],
 049 Izadi and Saeedi exploited a graph-based search to establish a set of rooftop hy-
 050 potheses through examining the relationships of lines and line intersections. Cui
 051 *et al.* [4] extracted buildings through the Hough transform (HT), but HT has
 052 notable drawbacks in parameter tuning and time complexity. Cote and Saeedi
 053 [5] generated rooftop outline from selected corners in multiple color and color-
 054 invariance spaces, further refined to fit the best possible boundaries through
 055 level-set curve evolution. In [6], Wang *et al.* presented a graph search-based
 056 perceptual grouping approach to hierarchically group line segments detected by
 057 EDLines [7] into candidate rectangular buildings, computation complexity of the
 058 approach was reduced dramatically compared to [1] [3] [5] [8]. However, geom-
 059 etric primitives based methods suffer from three serious shortcomings. Firstly, they
 060 lack the ability of detecting arbitrarily shaped building rooftop. Secondly, they
 061 fail to extract credible geometric features in buildings with inhomogeneous color
 062 distribution or low contrast with surroundings. Thirdly, it is hardly possible to
 063 process large-scale scenes because of their high computational-complexity.

064 Several studies reported that buildings are often composed of homogeneous
 065 regions with similar color or texture nearby shadows in remote sensing images.
 066 Spectral features is a distinctive feature for object detection, for instance, shad-
 067ows are commonly dark grey or black, vegetations are usually green or yellow
 068 with particular textures, and main roads are dim gray with different road marks
 069 in most case. According to these prior knowledge mentioned above, Ghaffarian
 070 *et al.* [9] proposed an purposive fast independent component analysis (PFastI-
 071 CA) technique to separate building area from remote sensing image. However,
 072 Ghaffarian's approach fails to detect the buildings with significantly different
 073 coloured rooftops. In [10], illumination direction and shadow area information
 074 of training samples were collected firstly, and then a improved parallelepiped
 075 classification method was applied to classify the image pixels into building and
 076 non-building areas. In [11], Chen *et al.* proposed a supervised building detection
 077 framework. At first, source image was divided into super-pixels using the SLIC
 078 [12] algorithm, then shadow patches are recognized using LDA color feature and
 079 the SVM classifier. The rough segmentation of buildings is employed by an adap-
 080 tive regional growth algorithm that considers the spatial relationship between
 081 shadows and buildings. Finally, buildings are segmented accurately using a level
 082 set model. Dornaika *et al.* [13] proposed a similar framework, Firstly, remote
 083 sensing image is segmentation by statistical region merging (SRM) algorithm
 084 , hybrid descriptor composed by color histograms and local binary patterns is
 085 used to represent each segmented region. Finally, each region was classified using
 086 machine learning tools and a gallery of training descriptors. However, a major
 087 problem of these methods is that they have less capability to detect building
 088 rooftop with significantly varying illumination in different parts or no surround-
 089 ing shadows.

045
 046
 047
 048
 049
 050
 051
 052
 053
 054
 055
 056
 057
 058
 059
 060
 061
 062
 063
 064
 065
 066
 067
 068
 069
 070
 071
 072
 073
 074
 075
 076
 077
 078
 079
 080
 081
 082
 083
 084
 085
 086
 087
 088
 089

090 Other studies presented synthetic approaches with fusion of multiple features
 091 to extract building footprints from aerial images. Baluyan *et al.* [14] proposed
 092 a method combining spectral and spatial features. Firstly, image is segmented
 093 into a set of rooftop candidates. Secondly, a SVM classifier is trained to dis-
 094 tinguish rooftop regions or nonrooftop regions using extracted multiple features
 095 in dataset. Finally, "histogram method" is devised to detect missed rooftops in
 096 previous step. In [15], Ngo *et al.* presented a novel approach for automated detec-
 097 tion of rectangular buildings. At the first step, image is decomposed into small
 098 homogeneous regions as candidates. At the second step, a merging process is
 099 then performed over regions having similar spectral traits to produce rectilinear
 100 building region in accordance with position of shadows. Li *et al.* [16] proposed a
 101 higher order conditional random field (HCRF) based method, which incorporates
 102 both pixel-and segment-level information for the segmentation of rooftops. They
 103 claimed that the proposed model outperforms the best unsupervised methods at
 104 rooftops with complex structures and sizes. Nonetheless, feature selection and
 105 parameter tuning are considerable troubles.

106

107 1.1 Related Works

108 In recent years, deep neural networks have been widely deployed in general image
 109 segmentation or scene labelling tasks. Mnih [17] presented a patch-based frame-
 110 work for learning to label aerial images. A carefully designed neural network
 111 architecture is designed for predicting buildings in aerial imagery, and then the
 112 output of this network is processed by conditional random fields(CRFs). Satito
 113 [18] improved Mnih's networks for extracting multiple kinds of objects simul-
 114 taneously, two techniques consisting of model averaging with spatial displace-
 115 ment(MA) and channel-wise inhibited softmax (CIS) are introduced to enhance
 116 the performance. However, these two methods need to crop test image into a
 117 fixed size, which not only increase the time loss, but also break the integrity
 118 of building. For example, these method achieve bad performance for jumbo-size
 119 building (see Fig. 5). Meanwhile, it takes about 72s to solve a 1500×1500 image
 120 with model averaging in single computer.

121

122 Mapping buildings from aerial image is essentially a problem of semantic
 123 segmentation. Recent work suggests a number of methods in processing natu-
 124 ral images. Chen et al. [19] present a system which combine the responses at
 125 the final convolutional layer with a fully connected Conditional Random Field
 126 (CRF). The system is able to localize segment boundaries at a quite high level
 127 of accuracy. An end-to-end network [20] which integrates CRF modelling with
 128 CNNs avoids off-line post-processing methods for object delineation. Noh et al.
 129 [21] apply a deconvolution network to each proposal in an input image, and con-
 130 struct the final semantic segmentation map by combining the results from all
 131 proposals in a simple manner.

132

133 Although these methods show promise in segmenting natural images, they
 134 have components not suited for building extraction. First, buildings are fre-
 135 quently occluded by shadows or trees (see Fig.2 (a)). It is impossible to directly
 136 delineate building boundaries even for human experts. Nonetheless, human can

090
 091
 092
 093
 094
 095
 096
 097
 098
 099
 100
 101
 102
 103
 104
 105
 106
 107
 108
 109
 110
 111
 112
 113
 114
 115
 116
 117
 118
 119
 120
 121
 122
 123
 124
 125
 126
 127
 128
 129
 130
 131
 132
 133
 134

135 estimate the actual footprints according to a strong prior knowledge that almost
 136 rooftops are regular polygons. Though [19] [20] achieve excellent performance
 137 in natural image, both of them can't achieve such a high intelligence. **Second,**
 138 **buildings have significantly variational appearance even in a single**
 139 **one.** (see Fig.2 (b)). Moreover, some buildings' appearance are very close to the
 140 plot on the ground or road. Based on our observation, there are a handful of
 141 samples emerged in PASACAL VOC 2012 dataset. Last but not the least, the
 142 number of objects in a remote sensing image is in a wide range. For example,
 143 some imagery include a large number of tiny buildings Fig.2 (c) and some ones
 144 are composed of moderate quantity of small-scale rooftops and a few of large-
 145 scale rooftops Fig.2 (d). [21] claimed that it handles objects in multiple scales,
 146 but it only suitable to multi-class object segmentation. Therefore, a system which
 147 can handle any-size rooftops is in high demand.

148 The rest of this article is organized as follows. The section 2.1 outlines some
 149 key ideas of fully convolutional network(FCN), the section 2.2 provides details
 150 of our neural network architecture. In section 2.3 the formulation of our system
 151 is proposed. The section 4.1 presents the datasets which is used for training
 152 and testing in our experiments. The section 3.2 introduces training settings and
 153 strategies of our proposed network. The section 3.3, we compare our results with
 154 two patch based methods using same dataset with us. In section 4, we discuss
 155 the experimental results and summarize whole article.

157 2 Network Architecture

159 In this section, we firstly introduce a well-known semantic segmentation network,
 160 called fully convolutional network(FCN)[22]. We attempt to extract footprints of
 161 building using architecture proposed by author and two modified versions. Our
 162 experiments show that these networks are not enough to accomplish building
 163 extraction task. **Therefore, we introduce a multi-scale fully convolutional**
 164 **network(MSFCN) for extracting rooftops, it turns out to be achieve**
 165 **state-of-the-art performance in a public aerial image dataset.**



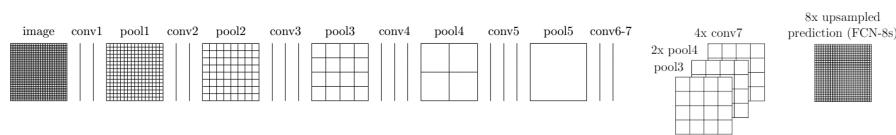
177 **Fig. 1.** Examples of aerial image

180 **2.1 Fully Convolutional Network** 180
181

182 As mentioned in introduction, patch-based building extraction methods [17] [18]
 183 suffer from shortcoming in processing big buildings because of the fixed-size inputs.
 184 As far as we know, fully convolutional network(FCN) takes whole image
 185 as inputs and directly outputs building rooftops by one pass of forward prop-
 186 agation. Because it can process input images of any sizes without warping or
 187 cropping, integrality of object is protected much better.

188 Long *et al.* put forward the concept of fully convolutional network(FCN) in
 189 [22] for the first time. **It said:** “While a general deep net computes a gen-
 190 eral nonlinear function, a net with only layers of this form computes
 191 a nonlinear filter, which we call a deep filter or fully convolutional
 192 network.” In practice, the fully connected layers in traditional CNNs are
 193 transformed to convolutional layers with 1×1 kernels. Due to the mechanism of
 194 pooling, the output of the network is a coarse heat map. For pixelwise prediction,
 195 skip-net technique is used to connect various coarse outputs back to pixels. For
 196 example, the framework of FCN-8s (see Fig. 2.1) is described as following. We
 197 firstly obtain the 16 stride predictions by fusing the predictions of *pool4* with
 198 the $2 \times$ upsampling of predictions from *conv7*(convolutionalized *fc7*). Then we
 199 continue in this fashion by fusing predictions from *pool3* with a $2 \times$ upsampling
 200 of the 16 stride predictions, denoted as 8 stride predictions. Finally, the 8 stride
 201 predictions are upsampled back to image with original size.

202 In our experiments, we first directly apply the FCN-8s to extract building
 203 rooftop by replacing the loss function with sigmoid cross-entropy loss. In order
 204 to achieving better performance, the network is extended to FCN-4s, FCN-2s.
 205 Our experiments show that FCN-4s and FCN-2s get the best performance with
 206 overall recall of 70.19 % in break-even point. According to our results, the first
 207 row of Fig. 3 indicate that FCN-4s has less ability to handle objects with small
 208 size. The second row of Fig. 3 shows that it is hard for FCN to discriminate
 209 buildings with ground if their color is similar. To overcome these drawbacks, a
 210 improved FCN is introduced in next section.

217 **Fig. 2.** The framework of FCN-8s 217
218219 **2.2 Our Architecture** 219
220

221 Recently, VGGNet [?] has been seen to achieve state-of-the-art performance
 222 in the ImageNet challenge, with great depth(16 convolutional layers), great
 223

224

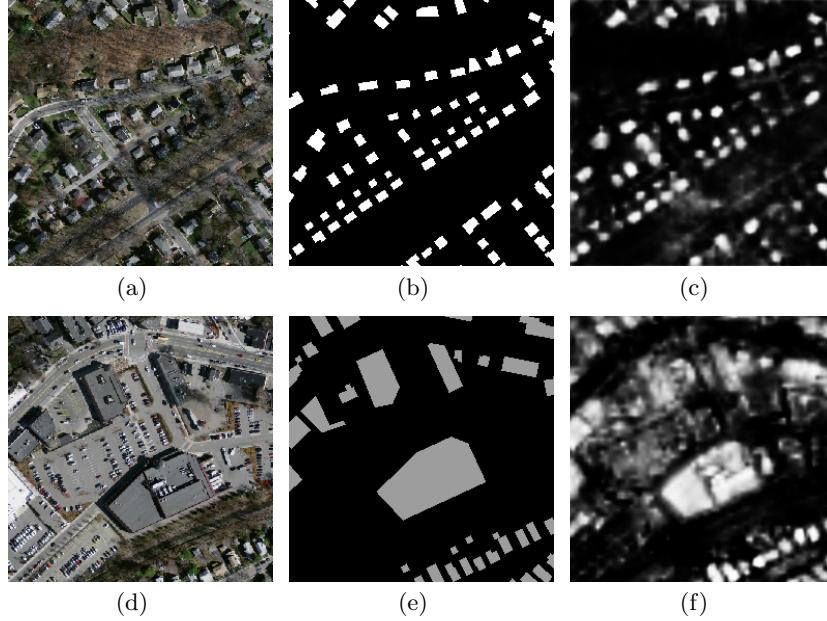


Fig. 3. (a),(d) Input image. (b),(e) Ground truth. (c),(f) FCN-4s prediction.

density(stride-1 convolutional kernels), and multiple stages(five 2-stride down-sampling layers). Many semantic segmentation architecture [22] [21] [19] [20] are based on this net. We therefore adopt the VGGNet architecture but make the following modifications:(a) two full connected layers and fifth pooling layer are cut. Because the interpolated prediction map of output of a 32 stride layer is too fuzzy to utilize. (b)

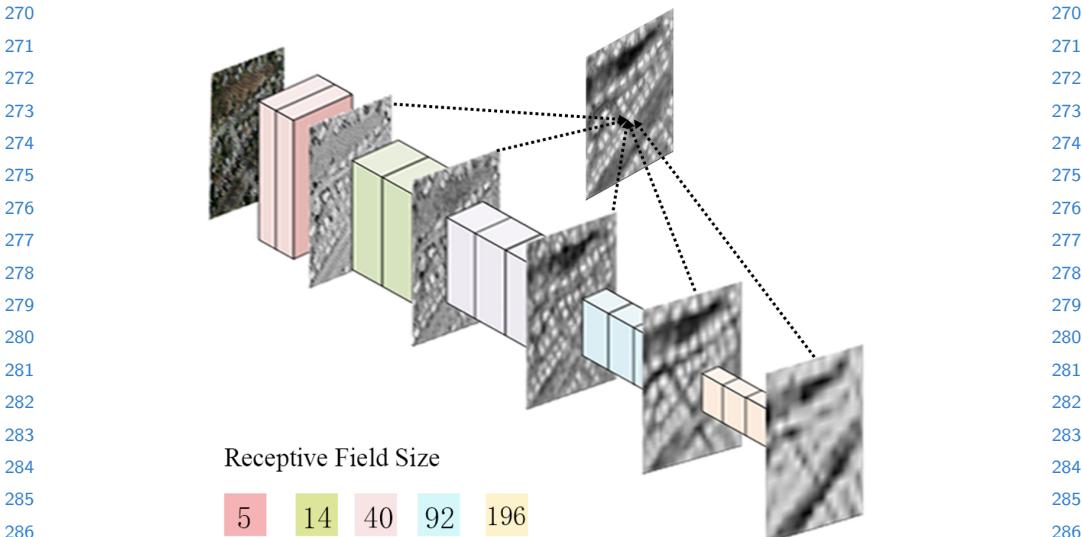
(1) Last stage of VGGNet is cut, including the 5th pooling layer and all the fully connected layers. (2) In order to integrating multi-scale predictions, we connect some output layers to the last convolutional layer in each stage, respectively conv1_2, conv2_2, conv3_3, conv4_3, conv5_3. Xie and Tu [23] have proved this network archive good performance for edge detection. Edge detection work on a very fine level, and our system need to extract many tiny objects, so, it is reasonable to utilize the similar architecture to extract rooftops.

layer	c1_1	c1_2	c2_1	c2_2	c3_1	c3_2	c3_3	c4_1	c4_2	c4_3	c5_1	c5_2	c5_3
rf size	3	5	10	14	24	32	40	48	76	92	108	164	196

2.3 Formulation

Our goal is to predict probabilistic label image \mathbf{M} from an input aerial image \mathbf{S} . Fig ?? shows an example of \mathbf{S} and \mathbf{M} . We directly learn a mapping from raw pixels in \mathbf{S} to a true label image \mathbf{M} by training the whole network.

Here we formulate our approach for building extraction. We denote our input training data set by $\mathbf{I} = \{(\mathbf{S}_n, \mathbf{M}_n), n = 1, \dots, |\mathbf{S}_n|\}$, where sample $\mathbf{S}_n =$

**Fig. 4.** Our Network Architecture

$\{s_j^{(n)}, j = 1, \dots, |\mathbf{S}_n|\}$ denotes the raw input image and $\mathbf{M}_n = \{m_j^{(n)}, j = 1, \dots, |\mathbf{S}_n|\}, m_j^{(n)} \in \{0, 1\}$ denotes the corresponding ground truth binary labelling map for satellite image \mathbf{S}_n . Taking account of each image holistically and independently, thus, we adopt the subscript n for notational simplicity. Our goal is to have a network that learns features from which it is possible to produce building maps approaching the ground truth. In our image-to-image training, the loss function is computed over all pixels in a training image $\mathbf{S} = \{s_j, j = 1, \dots, |\mathbf{S}|\}$ and building map $\mathbf{M} = \{m_j, j = 1, \dots, |\mathbf{S}|\}, m_j \in \{0, 1\}$. For simplicity, we denote the collection of all standard network layer parameters as \mathbf{W} . For each pixel j in a training image, the possibility that assigns it to building is denoted as $\hat{m}_j = Pr(m_j = 1 | \mathbf{S}; \mathbf{W})$. Specifically, the definition of sigmoid cross-entropy loss function is shown in Eq (1).

$$\mathcal{L}_{single} = - \sum_{j \in \mathbf{S}} [m_j \log \hat{m}_j + (1 - m_j) \log (1 - \hat{m}_j)] \quad (1)$$

3 Experiments

In this section, we discuss our detailed implementation and report the performance of our proposed algorithm.

3.1 Dataset

In our experiments, we use *Massachusetts Buildings Dataset* (Mass. Buildings) proposed by Mnih [17] and publicly available on website [http://www.cs.toronto.e](http://www.cs.toronto.edu/~mnih/)

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

8 ACCV-16 submission ID ***

315 du/vmnih/data/. The dataset consists of 151 aerial images of the Boston area,
 316 with each of the images being 1500×1500 pixels for an area of 2.25 square
 317 kilometers. Hence, the entire dataset covers roughly 340 square kilometers. The
 318 data is split into a training set of 137 images, a test set of 10 images and a
 319 validation set of 4 images. To train the network, we create image tiles of size
 320 256×256 by means of cropping entire image using a sliding window with size of
 321 256×256 and stride of 64 pixels. When scanning the whole dataset, image tiles
 322 with too many white pixels are removed. After scanning, train and validation
 323 dataset consists of 75938 and 2500 tiles and corresponding building masks. For
 324 testing, we use ten 1500×1500 entire images covering area excluded from the
 325 training data. In our experiments, input image should be scaled into range [0,1].
 326

327 3.2 Implementation

328 The implementation of our framework is based on the publicly available *Caffe*
 329 [24] Library. In our system, the whole network is fine-tuned from an initialization
 330 with the pre-trained VGG-16 Net model. Thus, it cost about four hour for
 331 training the network after 4000 iterations on a single NVIDIA Titan 12GB GPU.
 332

333 The network is trained in an end-to-end manner. No pre or post-processing
 334 is used. We train the networking using stochastic gradient descent with 20 im-
 335 ages as a mini-batch. The weight update rule is used with fixed learning rate
 336 10^{-5} , momentum 0.9, and weight decay 5×10^{-3} . To prevent gradient decreasing
 337 significantly, clip_gradients is set to 10000.

338 3.3 Results

339 4 Conclusions

340 References

- 344 Noronha, S., Nevatia, R.: Detection and modeling of buildings from multiple aerial
 345 images. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **23**
 346 (2001) 501–518
- 347 Nosrati, M.S., Saeedi, P.: A novel approach for polygonal rooftop detection in satel-
 348 lite/aerial imageries. In: *Image Processing (ICIP), 2009 16th IEEE International
 Conference on, IEEE* (2009) 1709–1712
- 349 Izadi, M., Saeedi, P.: Three-dimensional polygonal building model estimation from
 350 single satellite images. *Geoscience and Remote Sensing, IEEE Transactions on* **50**
 351 (2012) 2254–2272
- 352 Cui, S., Yan, Q., Reinartz, P.: Complex building description and extraction based
 353 on hough transformation and cycle detection. *Remote Sensing Letters* **3** (2012)
 354 151–159
- 355 Cote, M., Saeedi, P.: Automatic rooftop extraction in nadir aerial imagery of
 356 suburban regions using corners and variational level set evolution. *Geoscience and
 Remote Sensing, IEEE Transactions on* **51** (2013) 313–328
- 357 Wang, J., Yang, X., Qin, X., Ye, X., Qin, Q.: An efficient approach for automatic
 358 rectangular building extraction from very high resolution optical satellite imagery.
 359 *Geoscience and Remote Sensing Letters, IEEE* **12** (2015) 487–491

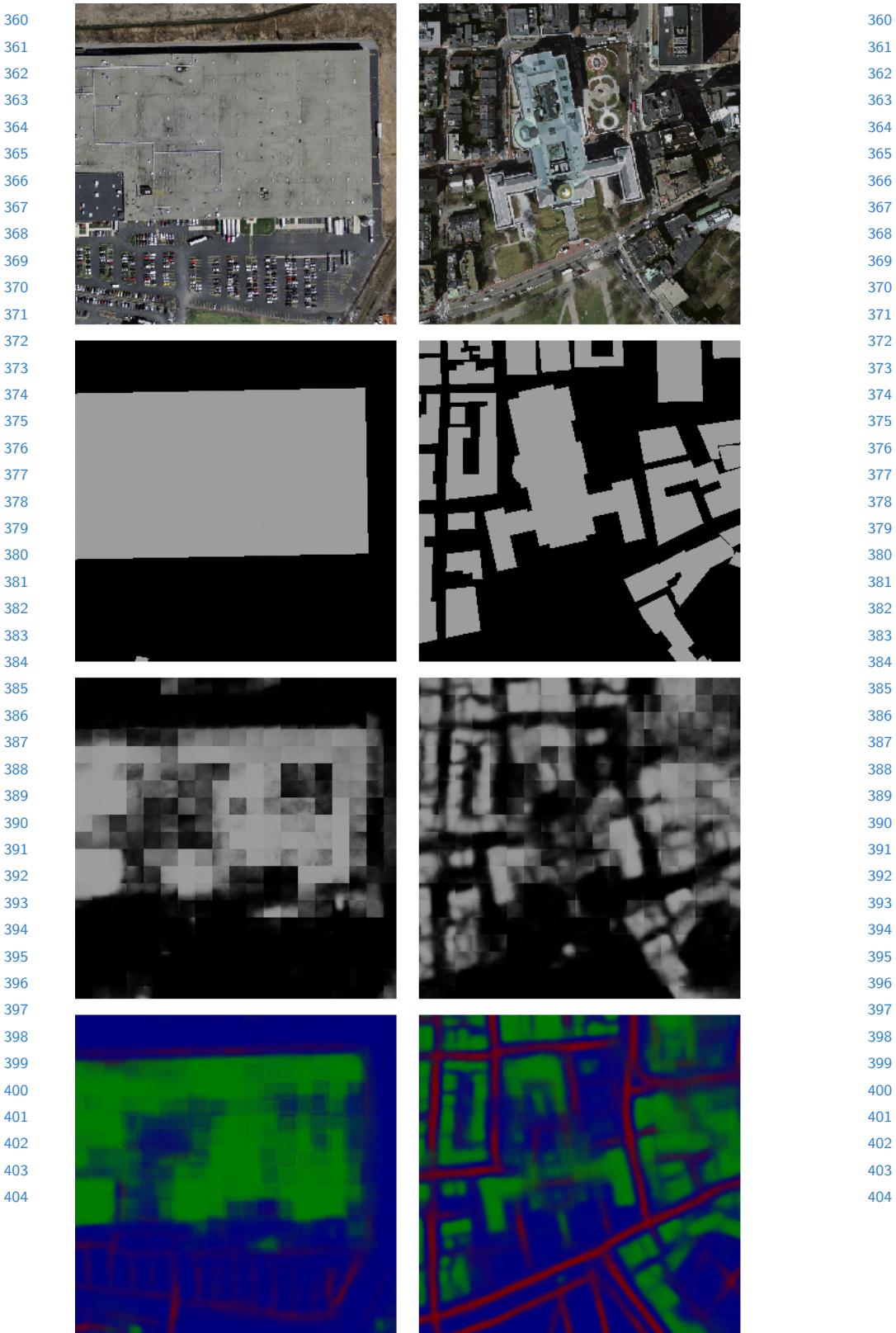


Fig. 5. (a) Input image. (b) Ground Truth. (c) Mnih's result. (d) Saito's result

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

10 ACCV-16 submission ID ***

- 405 7. Akinlar, C., Topal, C.: Edlines: A real-time line segment detector with a false
406 detection control. Pattern Recognition Letters **32** (2011) 1633–1642
407 8. Mayunga, S., Coleman, D., Zhang, Y.: A semi-automated approach for extracting
408 buildings from quickbird imagery applied to informal settlement mapping. International Journal of Remote Sensing **28** (2007) 2343–2357
409 9. Ghaffarian, S., Ghaffarian, S.: Automatic building detection based on purposive
410 fastica (pfica) algorithm using monocular high resolution google earth images. IS-
411 PRS Journal of Photogrammetry and Remote Sensing **97** (2014) 152–159
412 10. Ghaffarian, S., Ghaffarian, S.: Automatic building detection based on supervised
413 classification using high resolution google earth images. J. Photogr. Remote Sens.,
414 XL-3 (2014) 101–106
415 11. Chen, D., Shang, S., Wu, C.: Shadow-based building detection and segmentation
416 in high-resolution remote sensing image. Journal of Multimedia **9** (2014) 181–188
417 12. Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., Susstrunk, S.: Slic superpixels
418 compared to state-of-the-art superpixel methods. Pattern Analysis and
Machine Intelligence, IEEE Transactions on **34** (2012) 2274–2282
419 13. Dornaika, F., Moujahid, A., Bosaghzadeh, A., El Merabet, Y., Ruichek, Y.: Object
420 classification using hybrid holistic descriptors: Application to building detection in
421 aerial orthophotos. Polibits (2015) 11–17
422 14. Baluyan, H., Joshi, B., Al Hinai, A., Woon, W.L.: Novel approach for rooftop
423 detection using support vector machine. ISRN Machine Vision **2013** (2013)
424 15. Ngo, T.T., Collet, C., Mazet, V.: (Automatic rectangular building detection from
425 vhr aerial imagery using shadow and image segmentation)
426 16. Li, E., Femiani, J., Xu, S., Zhang, X., Wonka, P.: Robust rooftop extraction from
427 visible band images using higher order crf. Geoscience and Remote Sensing, IEEE
428 Transactions on **53** (2015) 4483–4495
429 17. Mnih, V.: Machine learning for aerial image labeling. Doctoral (2013)
430 18. Saito, S., Yamashita, Y., Aoki, Y.: Multiple object extraction from aerial imagery
431 with convolutional neural networks. Journal of Imaging Science & Technology **60**
432 (2016)
433 19. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Semantic
434 image segmentation with deep convolutional nets and fully connected crfs. In:
435 ICLR. (2015)
436 20. Zheng, S., Jayasumana, S., Romera-Paredes, B., Vineet, V.: Conditional random
fields as recurrent neural networks. (2015) 1529–1537
437 21. Noh, H., Hong, S., Han, B.: Learning deconvolution network for semantic segmen-
438 tation. (2015)
439 22. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic
440 segmentation. Computer Science **79** (2014) 1337–1342
441 23. Xie, S., Tu, Z.: Holistically-nested edge detection. Computer Science (2015) 1395–
1403
442 24. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadar-
443 rama, S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding.
444 Eprint Arxiv (2014) 675–678
445
446
447
448
449