

HF-FCN: Hierarchical Fusion Fully Convolutional Network for Robust Building Extraction

Robust Building Extraction from Large-scale Aerial Scene with Hierarchical Fusion Fully Convolutional Network

Anonymous ACCV 2016 submission

Paper ID ***

Abstract. Currently, automatic building extraction from remote sensing image plays a critical role in a diverse range of applications. However, it is significantly challenging to extract arbitrary-size buildings with large variational appearances or occlusions. To tackle these problems, we propose a robust system employing a novel hierarchical fusion fully convolutional network (HF-FCN), which effectively integrates the information generated from a group of neurons with multi-scale receptive fields. Our architecture can take a whole aerial images as inputs without warping or cropping and output building map directly. The experiment results tested on a publicly available aerial imagery dataset proved that our proposed methodology significantly shorts time-consuming and surpass the performance of state-of-the-art.

1 Introduction

With the rapid development of remote sensing technologies and popularization of geospatial related commercial softwares, very high resolution satellite images are easily accessible. These valuable data provides a huge fuel for interpreting real terrestrial scenes. Building rooftops is one of the most important type of terrestrial objects because it is essential for a wide range of technologies, such as, urban planning, automated map making, 3D city modelling, disaster assessment, military reconnaissance, etc. However, it is very costly and time-consuming to manually delineate the footprint of buildings even for human experts.

In recent decades, many researchers have made massive attempts to extract buildings automatically. Much of the past work defines criteria according to the particular characteristics of rooftop, such as, **polygonal boundary**[1][2][3][4], **homogeneous color or texture** [5], **surrounding shadows** [6][7][8][9], and their combinations [10][11]. However, such approaches are weakly capable of handling real-world data because hand-coded rules or probability models learned from small samples have strong dependency with data.

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

2 ACCV-16 submission ID ***

045 For example, they assume that profile of buildings is polygon while the shape of
 046 stadiums always is circle or oval. For the sake of deploying a practical building
 047 extraction system, Mnih *et al.* [12] created a huge publicly dataset including
 048 large-scale aerial images and corresponding human-labeled maps, and proposed
 049 a patch-based convolutional neural network to extract object locations of objects
 050 automatically. Based on Mnih's work, [13] improved the performance further.
 051 Though these methods achieve high performance, they still have limited ability
 052 to deal with two often appearing cases: (1) Buildings are severely occluded by
 053 shadows or trees. (2) Buildings possess moderately variational appearances.

054 Here, we present a robust building extraction system by developing a hierarchical
 055 fusion fully convolutional network(HF-FCN) trained on a publicly
 056 available large aerial imagery dataset[12]. HF-FCN provides a strong ability to
 057 extracting building rooftops even with significantly variational appearances and
 058 severely occlusions without any post-processing. Meanwhile, it further improves
 059 the overall accuracy than [12][13]. Compared with patch based methods [12][13],
 060 HF-FCN takes whole images as inputs without cropping or warping them to a
 061 fixed size and directly outputs segmentation maps by one pass of forward propagation.
 062 Therefore, time consuming is tremendously decreased for predicting
 063 building map from a large aerial image.

064 The rest of this article is organized as follows. In section 1.1, we summarize
 065 main methods for building extraction. Section 2.1 outlines some key ideas of fully
 066 convolutional network(FCN), section 2.2 provides details of our neural network
 067 architecture. In section 2.3, the formulation of our system is proposed. Section
 068 3.1 presents the dataset used for training and testing in our experiments. Section
 069 3.2 introduces training settings and strategies of our proposed network. In sec-
 070 tion 3.3, we compare our results with two patch based methods using the same
 071 publicly dataset. In section 4, we discuss the experimental results and summarize
 072 whole article.

073

074 2 Related Works

075

076 In previous literatures, one popular way of extracting buildings is employing
 077 their shape information. It is observed that rooftops have more regular shapes,
 078 which usually are rectangular or combinations of several rectangles.[1][2][3][4]
 079 exploited a graph-based search to establish a set of rooftop hypotheses through
 080 examining the relationships of lines and line intersections, and then removed
 081 the fake hypotheses using a series of criteria. [5] generated rooftop outline from
 082 selected corners in multiple color and color-invariance spaces, further refined to
 083 fit the best possible boundaries through level-set curve evolution. Though these
 084 geometric primitives based methods achieve good performance in high contrast
 085 remote sensing imagery, they suffer from three shortcomings. Firstly, they lack
 086 the ability of detecting arbitrarily shaped building rooftop. Secondly, they fail
 087 to extract credible geometric features in buildings with inhomogeneous color
 088 distribution or low contrast with surroundings. Thirdly, it is time-consuming to
 089 process large-scale scenes because of their high computational complexity.

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 Apart from using shape information, spectral information is a distinctive fea-
 091 ture for terrestrial object detection. For instance, shadows are commonly dark
 092 grey or black, vegetations are usually green or yellow with particular textures,
 093 and main roads are dim gray in most cases. According to these prior knowl-
 094 edge, Ghaffarian *et al.* [14] split aerial scenes into three components (respec-
 095 tively, shadows and the vegetation, roads and the bare soil, buildings) using a
 096 group of manually established rules. Afterwards, a purposive fast independen-
 097 t component analysis (PFastICA) technique is employed to separate building
 098 area from remote sensing image. However, their results are significantly sensi-
 099 tive to parameter choice. A feasible alternative strategy is to learn the appear-
 100 ance representation using supervised learning algorithm. [8][15][10][9] designed
 101 a similar framework. Firstly, an aerial image is divided into superpixels using a
 102 certain over-segmentation algorithm. Secondly, hand-crafted features, such as,
 103 color histograms or local binary patterns (LBP), are extracted from each over-
 104 segmented regions. Finally, each region is classified using machine learning tools
 105 and a gallery of training descriptors. **Because it's inevitable for machine**
 106 **learning method to mislabel regions with close appearance, additional**
 107 **information is utilized to refine previous results.** [9] removed false
 108 rooftops using a assumption that buildings are surrounded by shadows because
 109 of illumination. [10] devised a “histogram method” to detect missed rooftops.
 110 [11] selected probable rooftops after pruning out blobs using shadows, light
 111 direction, a series of shape criteria, and then these rooftops is refined by high
 112 order conditional random field. The drawbacks of these algorithms are threefold.
 113 (1) It is problematic to recognize a over-segmented region as buildings because
 114 terrestrial objects have huge variational appearances in real aerial scene. (2)
 115 Hand-craft features are sensitive to input data, therefore, it is not robust to pro-
 116 cess large-scale remote sensing images. (3) Additional information is unreliable.
 117 For instance, some low buildings have no shadow in its neighbourhood, and some
 118 buildings have unique structures which are not satisfied to hand-coded criteria.
 119

120 As mentioned above, traditional methods are not capable of adapting to real
 121 scenes with huge variational appearance, occlusion or low contrast. Recently,
 122 deep neural networks have been widely deployed in general image segmentation
 123 or scene labelling tasks. Mnih, a pioneer, have presented a patch-based frame-
 124 work for learning to label aerial images[12]. A neural network architecture is
 125 carefully designed for predicting buildings in aerial imagery, and the output of
 126 this network is processed by conditional random fields (CRFs). Satito [13] im-
 127 proved Mnih’s networks for extracting multiple kinds of objects simultaneously,
 128 two techniques consisting of model averaging with spatial displacement (MA)
 129 and channel-wise inhibited softmax (CIS) are introduced to enhance the perfor-
 130 mance. However, these methods need to crop test image into a fixed size, which
 131 not only increases the time consuming, but also breaks the integrity of buildings.
 132

133 Mapping buildings from aerial image is essentially a problem of semantic
 134 segmentation. Recent work suggests a number of methods in processing natu-
 135 ral images. In [16], the coarse predictions from a late stage are upsampled and
 136 integrated with predictions from previous stage to yield finer ones, the whole
 137

procedure is repeated in this fashion until reaching a desired output. Chen et al. [17] present a system which combines the responses at the final convolutional layer with a fully connected conditional random field (CRF). The system is able to localize segment boundaries at a quite high level of accuracy. An end-to-end network [18] which integrates CRF modelling with CNNs avoids off-line post-processing methods for object delineation. Noh et al. [19] apply a deconvolution network to each proposal in an input image, and construct the final semantic segmentation map by combining the results from all proposals in a simple manner.

Although these methods show promise in segmenting natural images, they have components not suited for building extraction. First, buildings are frequently occluded by shadows or trees (see Fig.2 (a)). It is challengable to delineate building boundaries even for human experts. **Though [17][18] achieve excellent performance in processing boundary of natural image, neither of them reported that they have strong ability to handle occlusions.** Second, buildings have significantly variational appearance even in a single one. (see Fig.2 (b)). Moreover, a number of buildings are very close to the plot on the ground or road (see Fig.2 (c)). Based on our observation, there are few samples emerged in PASACAL VOC 2012 dataset. Last but not the least, the size of objects in a remote sensing image is in a wide range. For example, some images include a large number of tiny buildings (see Fig.2 (d)) and some ones are composed of moderate quantity of small-scale rooftops and a few of large-scale rooftops. On account of low resolution (eight resolution of input image) of output from [16], precise structures are sacrificed severely. It is reported that [16] have limited ability to locate small objects. [19] claimed it handles objects in multiple scales, but it only suitable to multi-class object segmentation.



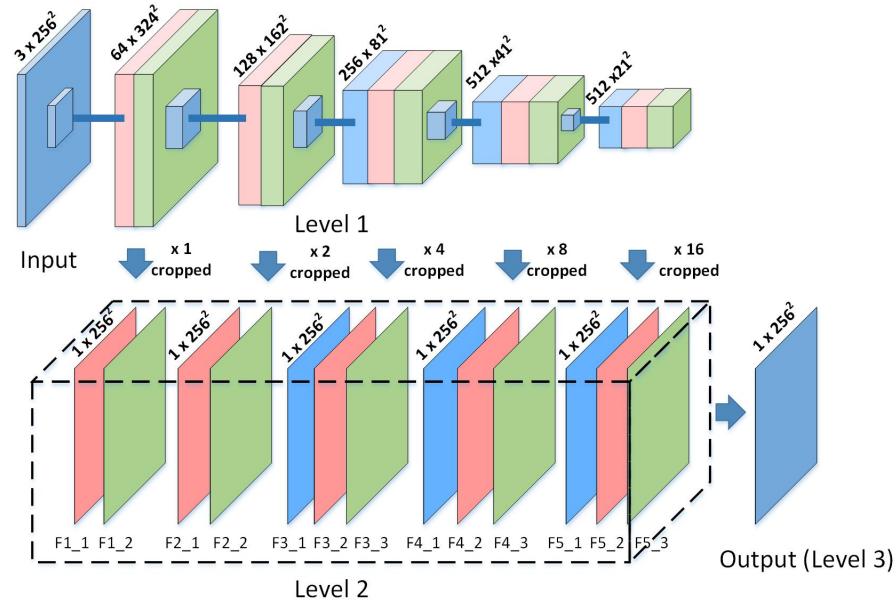
Fig. 1. Examples of aerial image

3 System Overview

In this section, we introduce a hierarchical fusion fully convolutional network (HF-FCN) for extracting rooftops, and then formulate our problem and loss function.

180 **3.1 Network Architecture** 180
181

182 We design our network based on VGG16 Net[20] and make some modification-
 183 s. The reasons for choose VGG16 Net are two-fold: (1) It has great depth (16
 184 convolutional layers), and multiple stages (five 2-stride down-sampling layer-
 185 s). We can acquire enough multi-level information from different stages and
 186 convolutional layers. (2) Network parameters pre-trained on very large image
 187 dataset(ImageNet) are helpful for initializing our network because our aerial data
 188 is essentially optical imagery. The modifications are listed as following: (1)
 189 Two fully connected layer $fc6$, $fc7$ and fifth pooling layer are cut, because they
 190 are at $\frac{1}{32}$ of input resolution. Meanwhile, the number of neurons in $fc6$, $fc7$ is
 191 too large to cost intensive computation. (2) **Feature maps from each convolu-**
 192 **tional layer in trimmed VGG16 Net (denote as level 1)** are fed
 193 **into a convolutional layer with a filter of 1×1 kernel and 1 neuron.** The
 194 outputs of these convolutional layers are upsampled and cropped to the same
 195 size of input image (denote as level 2). Upsampling is implemented via deconvolu-
 196 tion which is initialized by bilinear interpolation. Finally, All the feature maps
 197 in level 2 are stacked and put into a convolutional layer with a filter of 1×1
 198 kernel and 1 neuron to yield final predicted map (also denote as level 3). Our
 199 architecture is shown in Fig. 2.

222 **Fig. 2.** Our network architecture. 222
223 224

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

6 ACCV-16 submission ID ***

225

Table 1. The receptive field and stride size in level 2 of our architecture.

226

layer	F1_1	F1_2	F2_1	F2_2	F3_1	F3_2	F3_3	F4_1	F4_2	F4_3	F5_1	F5_2	F5_3
rf size	3	5	10	14	24	32	40	60	76	92	124	164	196
stride	1	1	2	2	4	4	4	8	8	8	16	16	16

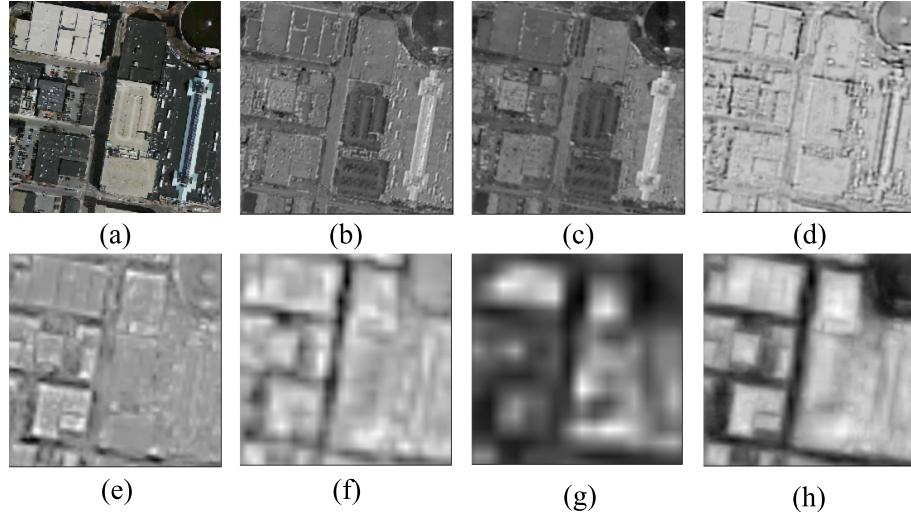
227

230

In level 2 of our architecture, feature maps with increasing receptive field (see Table 3.1) capture local information in different neighbourhood sizes and at different semantic levels. Therefore, if we integrate all these information together, it is helpful for extracting buildings with variational appearance or occlusion. We take a concrete instance to show how HF-FCN works for such cases. In this case, F1_1 with small receptive field generates fine spatial resolution and responds to low level features like edges and corners (see Fig. 3(b)). F1_2 functions like over-segmentation algorithm to grouping pixels with similar color or texture into a subregion (see Fig. 3(c)). In F2_1, color information is disappear, shape information is augmented (see Fig. 3(d)). In F3_3, it is surprised that regions with significantly varying appearance are merged into a integrated building **by considering an unknown high level features** (see Fig. 3(e)). In F4_2 and F5_2, our network learned strong semantic knowledge to distinguish dark rooftops with dim shadows and dark-green water (see Fig. 3(f)(g)). In level 3, we show that HF-FCN obtains reliable prediction by combining multi-level semantic information and spatial information (see Fig. 3(h)).

247

248



265

266

Fig. 3. (a) is input aerial image, feature maps generated from F1_1 (b), F1_2 (c), F2_1 (d), F3_3 (e), F4_2 (f), F5_2 (g), level 3 (h)

267

268

269

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270 **3.2 Formulation**

271

272 Our goal is to predict labelling image $\hat{\mathbf{M}}$ from an input aerial image \mathbf{S} . We
 273 directly learn a mapping from raw pixels in \mathbf{S} to a true label image $\tilde{\mathbf{M}}$ by training
 274 the whole network. Fig. 4 shows an example of \mathbf{S} , $\tilde{\mathbf{M}}$, $\hat{\mathbf{M}}$. Here we formulate
 275 our approach for building extraction. We denote our input training data set
 276 by $\mathbf{I} = \{(\mathbf{S}_n, \tilde{\mathbf{M}}_n), n = 1, \dots, |\mathbf{S}_n|\}$, where sample $\mathbf{S}_n = \{s_j^{(n)}, j = 1, \dots, |\mathbf{S}_n|\}$
 277 denotes the raw input image and $\tilde{\mathbf{M}}_n = \{\tilde{m}_j^{(n)}, j = 1, \dots, |\mathbf{S}_n|\}$, $\tilde{m}_j^{(n)} \in \{0, 1\}$
 278 denotes the corresponding ground truth binary labelling map for satellite image
 279 \mathbf{S}_n . Taking account of each image holistically and independently, thus, we adopt
 280 the subscript n for notational simplicity. Our goal is to have a network that
 281 learns features from which it is possible to produce building maps approaching
 282 the ground truth. In our image-to-image training, the loss function is computed
 283 over all pixels in a training image $\mathbf{S} = \{s_j, j = 1, \dots, |\mathbf{S}|\}$ and building map
 284 $\tilde{\mathbf{M}} = \{\tilde{m}_j, j = 1, \dots, |\mathbf{S}|\}$, $\tilde{m}_j \in \{0, 1\}$. For simplicity, we denote the collection of
 285 all standard network layer parameters as \mathbf{W} . For each pixel j in a training image,
 286 the possibility that assigns it to building is denoted as $\hat{m}_j = Pr(m_j = 1 | \mathbf{S}; \mathbf{W})$.
 287 the definition of sigmoid cross-entropy loss function is shown in Eq (1).

288

$$289 \quad \mathcal{L} = -\frac{1}{|\mathbf{S}|} \sum_{j \in \mathbf{S}} [\tilde{m}_j \log \hat{m}_j + (1 - \tilde{m}_j) \log (1 - \hat{m}_j)] \quad (1)$$

290
291

304 **Fig. 4.** An example of the resulting predicted image.310 **4 Experiments**

311

312 In this section, we discuss our detailed implementation and report the perfor-
 313 mance of our proposed algorithm.

314

315 **4.1 Dataset**

316 In our experiments, we use Massachusetts Buildings Dataset (*Mass. Buildings*)
 317 proposed by Mnih [12] and publicly available on website [http://www.cs.toronto.e
 318 du/~vnnih/data/](http://www.cs.toronto.edu/~vmnih/data/). The dataset consists of 151 aerial images of the Boston area,
 319 with each of the images being 1500×1500 pixels for an area of 2.25 square k-
 320 ilometers. Hence, the entire dataset covers roughly 340 square kilometers. The
 321 data is split into a training set of 137 images, a test set of 10 images and a vali-
 322 dation set of 4 images. To train the network, we create image tiles for train and
 323 validation by means of cropping entire image using a sliding window with size
 324 of 256×256 pixels and stride of 64 pixels. When scanning the whole dataset, im-
 325 age tiles which include more than 160 white pixels are removed. After scanning,
 326 train and validation dataset include 75938 tiles and 2500 tiles with corresponding
 327 building masks. For testing, we use ten 1500×1500 entire images covering area
 328 excluded from the training data. In our experiments, we find that it is benefit
 329 for prediction to scale the intensity of input image into range of [0,1].
 330

331 **4.2 Implementation**

332 The implementation of our networks are based on the publicly available *Caffe*
 333 [21] Library. These three networks are fine-tuned from an initialization with
 334 the pre-trained VGG16 Net model and trained in an end-to-end manner. All
 335 of our networks are trained using stochastic gradient descent with same hyper-
 336 parameters, including mini-batch size (18), initial learning rate (10^{-5}), learning
 337 rate is divided by 10 for each 5000 iterations, momentum (0.9), weight decay
 338 (0.02), clip_gradients (10000), number of training iterations (12000). We find that
 339 learned deconvolutions provide no noticeable improvements in our experiments,
 340 therefore, lr_mult is set to zero for all deconvolutional layers. It takes about six
 341 hours to train a network on a single NVIDIA Titan 12GB GPU.
 342

343 **4.3 Results**

344 To show the effectiveness of HF-FCN, we train and test our network on *Mass.*
 345 *Buildings*. In order to comparing our results with previous works [12][13], we use
 346 three metrics to evaluate our results: (1)relaxed precision and recall scores. (2)
 347 standard precision and recall scores. (3) time consuming. The relaxed precision
 348 is defined as the fraction of detected pixels that are within ρ pixels of a detected
 349 pixel, while the relaxed recall is defined as the fraction of true pixels that are
 350 within ρ pixels of a detected pixel. In our experiments, the slack parameter ρ is
 351 set to 3, which is the same value as used in [12][13]. Compared relaxed precision-
 352 recall curves are shown in Fig. 5(a). In order to evaluate our results more strictly,
 353 we set slack parameter ρ as 0, that is to say, it becomes a standard precision and
 354 recall scores. Compared standard precision-recall curves are shown in Fig. 5(b).
 355 Additionally, time consuming is another important index to evaluate system
 356 performance. We calculate the mean time of processing ten test images in the
 357 same computer using the same program. Table 4.3 shows that our method is able
 358
 359

360 to not only significantly improve the performance , but dramatically decrease
 361 the time-consuming.

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

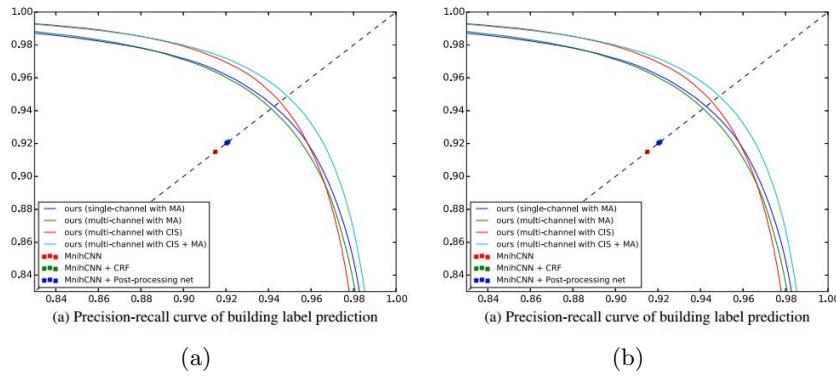
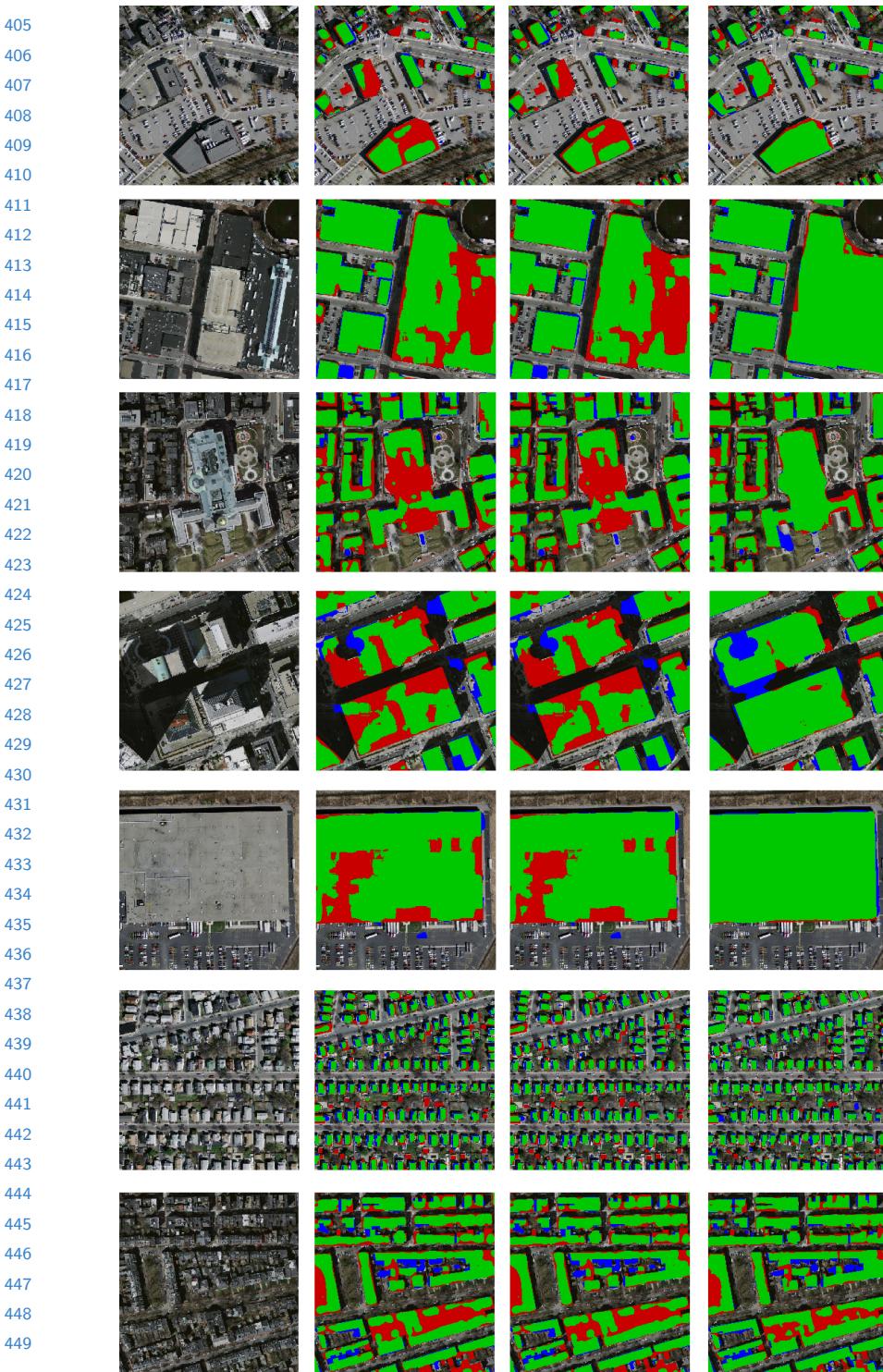


Fig. 5. (a) Relaxed precision-recall curve. (b) standard precision-recall curve.

To prove our network having strong ability in extracting buildings with variational appearances and occlusions, we perform further evaluation. We crop five 256×256 image patches that have large-size buildings with variational appearances or occlusions from test image of *Mass. Buildings*. And then, we directly crop corresponding predictions from predicted images generated by three models (Mnih-CNN+CRF[12], Saito-multi-MA&CIS[13] and ours). Here, we binarize the probability map using a threshold of 0.5. Seven groups of example are shown in Fig. 6. In addition, Table 4.3 shows the resulting recalls at breakeven points of standard precision recall curve for each patches.

Table 2. Performance is compared with [12][13]. Recall here means recall at breakeven points. Time is computed in the same computer with a single NVIDIA Titan 12GB GPU.

	Recall(relax = 3)	Recall(relax = 0)	Time(s)
Mnih-CNN[12]	0.9150	0.7661	8.70
Mnih-CNN+CRF[12]	0.9211		
Saito-multi-MA[13]	0.9426	0.7858	67.72
Saito-multi-MA&CIS[13]	0.9488	0.7857	67.84
Ours (HF-FCN)	0.9643	0.8424	1.07



(a)

(b)

(c)

(d)

Fig. 6. (a) Input image. (b) Results of Mnih-CNN+CRF[12]. (c) Results of Saito-multi-MA&CIS[13]. (d) Our results. Correct results (TP) are shown in green, false positives are shown in blue, and false negatives are shown in red.

Table 3. Recall at selected region of the test images

Image ID	01	02	03	04	05	06	07	mean
Mnih-CNN+CRF[12]								
Saito-multi-MA&CIS[13]	0.773	0.915	0.857	0.789	0.945	0.773	0.830	0.851
Ours (HF-FCN)	0.874	0.964	0.901	0.986	0.933	0.840	0.851	0.911

5 Conclusions

In this article, we proposed a improved fully convolutional network which is strongly capable of extracting buildings with variational appearance or occlusion-s. The network can take arbitrary-size image as input as long as GPU memory allowed. As a consequence, there is no need to label a whole image by cropping the image into small patches. Meanwhile, inconsistant border caused by cropped would not occurred in our system. Though a effective technique[13], namely, model averaging with spatial displacement, is proposed, it is troublesome to train a network eight times and predict labelling image with the same times. On the other hand, we demonstrate that our network is generally adapt to various types of aerial scenes selected from real-world data. Furthermore, our architecture can be easily extended to extract multi-objects in remote sensing imagery. Consequently, we believe that our technique potentially provides a generic solution to understand complex aerial scenes.

References

1. Noronha, S., Nevatia, R.: Detection and modeling of buildings from multiple aerial images. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* **23** (2001) 501–518
2. Nosrati, M.S., Saeedi, P.: A novel approach for polygonal rooftop detection in satellite/aerial imageries. In: *Image Processing (ICIP), 2009 16th IEEE International Conference on*, IEEE (2009) 1709–1712
3. Izadi, M., Saeedi, P.: Three-dimensional polygonal building model estimation from single satellite images. *Geoscience and Remote Sensing, IEEE Transactions on* **50** (2012) 2254–2272
4. Wang, J., Yang, X., Qin, X., Ye, X., Qin, Q.: An efficient approach for automatic rectangular building extraction from very high resolution optical satellite imagery. *Geoscience and Remote Sensing Letters, IEEE* **12** (2015) 487–491
5. Cote, M., Saeedi, P.: Automatic rooftop extraction in nadir aerial imagery of suburban regions using corners and variational level set evolution. *Geoscience and Remote Sensing, IEEE Transactions on* **51** (2013) 313–328
6. Sirmacek, B., Unsalan, C.: Building detection from aerial images using invariant color features and shadow information. In: *International Symposium on Computer and Information Sciences*. (2008) 1–5
7. Ok, A.O., Senaras, C., Yuksel, B.: Automated detection of arbitrarily shaped buildings in complex environments from monocular vhr optical satellite imagery. *Geoscience and Remote Sensing, IEEE Transactions on* **51** (2013) 1701–1717

CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

12 ACCV-16 submission ID ***

- 495 8. Chen, D., Shang, S., Wu, C.: Shadow-based building detection and segmentation
496 in high-resolution remote sensing image. *Journal of Multimedia* **9** (2014) 181–188
497 9. Ngo, T.T., Collet, C., Mazet, V.: (Automatic rectangular building detection from
498 vhr aerial imagery using shadow and image segmentation)
499 10. Baluyan, H., Joshi, B., Al Hinai, A., Woon, W.L.: Novel approach for rooftop
500 detection using support vector machine. *ISRN Machine Vision* **2013** (2013)
501 11. Li, E., Femiani, J., Xu, S., Zhang, X., Wonka, P.: Robust rooftop extraction from
502 visible band images using higher order crf. *Geoscience and Remote Sensing, IEEE*
Transactions on **53** (2015) 4483–4495
503 12. Mnih, V.: Machine learning for aerial image labeling. Doctoral (2013)
504 13. Saito, S., Yamashita, Y., Aoki, Y.: Multiple object extraction from aerial imagery
505 with convolutional neural networks. *Journal of Imaging Science & Technology* **60**
506 (2016)
507 14. Ghaffarian, S., Ghaffarian, S.: Automatic building detection based on purposive
508 fastica (pfica) algorithm using monocular high resolution google earth images. *IS-
509 PRS Journal of Photogrammetry and Remote Sensing* **97** (2014) 152–159
510 15. Dornaika, F., Moujahid, A., Bosaghzadeh, A., El Merabet, Y., Ruichek, Y.: Object
511 classification using hybrid holistic descriptors: Application to building detection in
aerial orthophotos. *Polibits* (2015) 11–17
512 16. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic
513 segmentation. *Computer Science* **79** (2014) 1337–1342
514 17. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Semantic
515 image segmentation with deep convolutional nets and fully connected crfs. In:
516 *ICLR*. (2015)
517 18. Zheng, S., Jayasumana, S., Romera-Paredes, B., Vineet, V.: Conditional random
518 fields as recurrent neural networks. (2015) 1529–1537
519 19. Noh, H., Hong, S., Han, B.: Learning deconvolution network for semantic segmen-
520 tation. (2015)
521 20. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale
522 image recognition. *Eprint Arxiv* (2014)
523 21. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama,
524 S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding.
Eprint Arxiv (2014) 675–678
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539