

An Ensemble Architecture of Deep Convolutional Segnet and Unet Networks for Building Semantic Segmentation from High-resolution Aerial Images

Abstract

Building objects is one of the principal features that are essential for updating the geospatial database. Extracting building features from high-resolution imagery automatically and accurately is challenging because of the existence of some obstacles in these images, such as shadows, trees, and cars. Although deep learning approaches have shown significant improvements in the results of image segmentation in recent years, most deep neural networks still cannot achieve highly accurate results with correct segmentation map when processing high-resolution remote sensing images. Therefore, we implemented a new deep neural network named Seg–Unet method, which is a composition of Segnet and Unet techniques, to exploit building objects from high-resolution aerial imagery. Results obtained 92.73% accuracy carried on the Massachusetts building dataset. The proposed technique improved the performance to 0.44%, 1.17%, and 0.14% compared with fully convolutional neural network (FCN), Segnet, and Unet methods, respectively. Results also confirmed the superiority of the proposed method in building extraction.

Keywords: building extraction; image segmentation; remote sensing; Seg–Unet approach

1. Introduction

Highly accurate feature extraction from high-resolution remote sensing imagery produces reliable information for various applications (Shrestha and Vanneschi 2018). The extraction of small ground objects, such as building objects from the imagery of the surface of the earth, can be a potential application (Krizhevsky *et al.* 2012). High-precision building extraction from high-resolution satellite images can perform an essential task in several applications, such as disaster management, geospatial database updating, urban planning, and navigation (Mayer 1999). Raw data should be converted into sensible information by using geospatial information system (GIS) to enable the quantification process. The time-consuming and labor-intensive data interpretation and digitization are often required for this transformation. Although Yuan (2017) introduced a source called volunteered geographic information (VGI) as an alternative option, its availability is restricted due to the differences in positional and completeness accuracy. Participation inequality,

30 in terms of varying impressions, cultures, and judgments, can be the principal reasons for the
31 aforementioned issue (Shrestha and Vanneschi 2018), thereby restricting the accessibility of
32 dependable and up-to-date building maps. Automatic building extraction using remote sensing
33 imagery needs a promising approach that remains underdeveloped in spite of a decade of research
34 in this field (Marcu and Leordeanu 2016). The main elements that make this process challenging
35 are the wide changes in building appearances in images because of various building features, such
36 as shadows, cars, structures, various roofing materials, and illumination statuses, which are formed
37 by buildings (Yuan and Cheriyyadat 2014). Traditional methods have been mixed with genetic
38 algorithms (Sumer and Turker 2013) and support vector machine (SVM) method (Inglada 2007)
39 to detect buildings. Other characteristics, such as multi-spectral features, textures (Levitt and
40 Aghdasi 1998), and shadow properties (Peng and Liu 2005); local structures, such as corners, lines,
41 and edges (Huertas and Nevatia 1988) of remote sensing images, have been utilized as main factors
42 for extracting building objects. The efficiency of these types of approaches is restricted due to the
43 dependence of the method performance on low-level local characteristics. Thus, to well distinguish
44 the features, the utilization and exploitation of representative high-level features that play a
45 principal role in image segmentation are favorable.

46 In recent studies, feature-based deep convolutional approaches, such as convolutional neural
47 network (CNN), have demonstrated that they can achieve reliable results in image classification
48 for computer vision (He *et al.* 2015, Szegedy *et al.* 2015) and feature semantic segmentation
49 (Vakalopoulou *et al.* 2015, Alshehhi *et al.* 2017, Abdollahi *et al.* 2020). The CNN model is
50 efficient in image processing because of its capability to learn from raw images without following
51 pre-processing steps. In addition, deep convolutional network (DCNN) has become a promising
52 technique in image processing because of its ability to efficiently mix spatial and spectral features
53 on the basis of raw input data without preprocessing (Alshehhi *et al.* 2017). Recent works have
54 revealed that different kinds of deep learning approaches, which are based on CNNs, such as deep
55 convolutional encode–decoder architecture and fully convolutional network (FCN), have shown
56 significant improvements in the remote sensing field. In terms of computational proficiency and
57 accuracy, FCN is the most proficient approach for pixel-wise semantic segmentation. However,
58 several problems restrict model performance in detection, leading to failure in generating
59 inadequate or redundant prediction detection and in identifying numerous objects (Shrestha and
60 Vanneschi 2018, Abdollahi *et al.* 2020). In the next section, previous studies related to applying

61 promising CNN methods for remote sensing image classification and building semantic
62 segmentation are discussed.

63 Deep neural network features have illustrated their ability in semantic segmentation (Long *et al.*
64 *2015*, Chen *et al.* 2017), object detection (Girshick *et al.* 2014), and visual identification (Sharif
65 Razavian *et al.* 2014, Audebert *et al.* 2016). Deep convolutional frameworks can be utilized in
66 different remote sensing tasks, such as data merging (Kussul *et al.* 2016), image classification
67 (Yang *et al.* 2018), and detection (Audebert *et al.* 2017). These networks have been successfully
68 utilized to label and classify high-resolution remote sensing images (Penatti *et al.* 2015). Marmanis
69 *et al.* (2018) introduced a deep neural network on the basis of an end-to-end trainable network
70 (DCNN) for detecting boundaries and improving semantic image segmentation. Farabet *et al.*
71 (2012) mixed conditional random fields (CRFs) with multi-scale CNNs to classify dense street
72 scenes. Vakalopoulou *et al.* (2015) implemented a deep convolutional model to identify building
73 features from high-resolution multi-spectral images. Previous works have confirmed that the
74 results of remote sensing imagery classification cannot be decisive (Wilkinson 2005) because
75 improving the resolution of remote sensing images is more useful in the identification and
76 detection of different features on the ground. However, the separation of certain objects with the
77 same spectral values has become difficult due to these improvements, leading to the decrease of
78 the inter-class difference and increase of the intra-class difference of objects such as cars, shadows,
79 streets, and buildings (Paisitkriangkrai *et al.* 2016). That is, extracting sensible spatial features to
80 solve the pixel classification in building extraction has become challenging because various
81 objects may represent similar spectral classes in remote sensing images. Reliable results have been
82 recently achieved by FCN for semantic image segmentation (Fu *et al.* 2017). The method can
83 identify various object classes, including their shapes, such as trees, road objects, and building
84 curves. The model can not only identify the structures of spatial objects but also learn how to
85 categorize pixels and detect what they are (Audebert *et al.* 2016). However, the outcomes are
86 visually degraded during image classification and segmentation when using FCN. The reason is
87 that the model cannot detect objects with multiple borders or small objects because object
88 boundaries are blurred (Maggiori *et al.* 2017). The structures of deep convolutional frameworks
89 have been developed in certain research either by utilizing CRFs mixed with dilated convolution
90 (Chen *et al.* 2014) or by appending skip-layer structure after up-sampling to regenerate high-
91 frequency and comprehensive image information (Marmanis *et al.* 2016), thereby leading to the

92 performance improvement of semantic segmentation and accuracy improvement of image
93 classification (Sherrah 2016).

94 Recent works have attempted to boost precision in areas such as pixel labeling; feature
95 extraction from raw data; image encoding, specifically for high-resolution remote sensing imagery
96 on the basis of deep convolutional techniques, such as FCN and CNN (Volpi and Tuia 2016).
97 However, impervious and building objects extracted from high-resolution remote sensing images
98 are difficult to handle due to the presence of various geometric shapes and spatial and spectral
99 features. That is, similar objects in urban areas have various spectral values because high-
100 resolution remote sensing images are usually restricted to three or four channels, and these spectral
101 characteristics may lack the capability to recognize objects. Various objects may also present the
102 same spectral values (e.g., roofs and roads) (Bakhtiari *et al.* 2017, Abdollahi *et al.* 2018).

103 Although prior scholars have presented helpful insights into different approaches, which can be
104 utilized in pixel labeling, these approaches misclassify certain pixels with the same spectral values
105 and lack the capability to eliminate salt-and-pepper classification noise and to clearly identify
106 object boundary. To solve these issues, we present a new deep neural network called SegUnet, a
107 combination of Segnet and Unet architectures for building objects extraction by using high-
108 resolution aerial imagery. The proposed network is dedicated to restoring pixel position
109 information and produces a high-resolution segmentation map. The model has an encoder–decoder
110 architecture that incorporates index pooling (Segnet) and skip connection (Unet) to generate and
111 disseminate image spatial information. As can be seen in the aforementioned literature review, the
112 proposed method has not been used before, and this study is the first to propose this kind of
113 approach for a given task. The proposed approach is compared with other state-of-the-art deep
114 learning-based techniques, such as FCN (Long et al. 2015), Segnet (Badrinarayanan et al. 2017),
115 and Unet (Ronneberger et al. 2015) on the basis of a similar dataset to demonstrate the ability of
116 the method in building extraction. Such outcomes prove that the new proposed network is efficient
117 in building extraction. The remainder of the paper is organized in the following manner. Section 2
118 outlines the methodology of the suggested SegUnet approach. Section 3 highlights the results and
119 discussion. Section 4 provides the conclusion.

120 **2. Materials and methodology**

121

122 In this section, we explain the overall framework of Unet, Segnet, and SegUnet models (Figure
123 1). Subsequently, the prepared high-resolution remote sensing aerial dataset for applying the
124 proposed approach is explained. Finally, the common metrics for calculating the performance of
125 state-of-the-art techniques applied for building extraction are described.

126 **Figure 1.** around here

127 **2.1. Unet architecture**

128 The Unet model is an elegant DCNN that can yield accurate image segmentations. The main
129 concept of the Unet model is the replacement of pooling layers with up-sampling operators to
130 complete a typical contracting network by continuous layers, followed by the enhancement of
131 output layer resolution. For localization, the high-resolution features of the contracting part are
132 mixed with up-sampled output. Finally, continuous convolution layer can be used to assemble an
133 accurate outcome on the basis of this information (Long *et al.* 2015). One significant factor in the
134 Unet model is the several feature channels in the up-sampling section where the network can
135 spread context information to layers with high resolution. The Unet deep learning model comprises
136 two principal sections: expansive part (right side) and contracting part (left side). Given that the
137 contracting and expansive parts are symmetric, a U-shaped CNN is formed. The model only
138 utilizes the right part of every convolution and does not have any fully connected layers. For
139 example, considering that the segmentation map only contains pixels, the entire context is
140 accessible in the input image. Therefore, an overlap-tile strategy is utilized to provide a monolithic
141 and random segmentation of large images. For extrapolating the missing context and foretelling
142 pixels in the border section of images, input image mirroring is also utilized. The resolution can
143 be restricted by the GPU memory unless the tiling strategy implements the network to extensive
144 images.

145 The generic framework of a Unet model is followed by a contracting path that includes two
146 repeated convolution layers of 3×3 window size, followed by a down-sampling layer of 2×2
147 window size. Activation function (1), which is a kind of transformation function, is used in the
148 convolution process. Assuming that a weight vector is w ; a bias vector is b ; $x_k(ii, jj)$ is the input of
149 activation function and the output of convolution operation, respectively.

150
$$Z(x_k(ii, jj)) = f\left(\sum_{k=1}^k x_k(ii, jj) \cdot w_k + b_k\right) \Leftrightarrow Z = f(X \cdot W + b). \quad (1)$$

151 For $f(\cdot)$, activation function (2), that is, rectified linear unit (ReLU) is used in the Unet model.
 152 Neurons do not confront the gradient vanishing issue, which arises when the gradient norm declines
 153 after sequential updates in the back-propagation process. Neurons also efficiently operate with
 154 rectified function because this function encourages sparsity in the hidden layers and prevents
 155 saturation during the learning process (Zhou *et al.* 2014). In each down-sampling stage, the number
 156 of feature channels is doubled. As previously mentioned, max-pooling layers are utilized to
 157 decrease image size, parameter number, and network computing. In the down-sampling method,
 158 images are sampled using their principle local correlations. This approach retains efficient
 159 information while lessening data processing and allowing the features taken through convolution to
 160 have spatial uniformity (Maggiori *et al.* 2017).

161 An up-sampling, followed by a convolution with a stride of 2×2 that halves the number of
 162 feature channels, is used in each step of the expansive path. Two convolution layers of 3×3 kernel
 163 size, followed by the ReLU activation function and a concatenation with the correspondingly
 164 cropped feature map from the contracting path, are utilized in the expansive path. Eventually, a
 165 convolution layer of 1×1 window size and a sigmoid function (3) are utilized for mapping every
 166 32-component feature vector to the desired number of classes (road and non-road) and for mapping
 167 the predicted values to probabilities, respectively (Hu *et al.* 2015). The generic framework of the
 168 Unet model is illustrated in Figure 2.

169
$$A(x_k(ii, jj)) = \max(0, Z(x_k(ii, jj))), \quad (2)$$

170 where $x_k(ii, jj)$ is utilized as the input to the activation function and the output of convolution
 171 operation, respectively.

172
$$S(z) = \frac{1}{1 + e^{-z}}, \quad (3)$$

173
 174 where S is the output between 0 and 1, and z is the input.

175
 176 **Figure 2.** around here

177

178 2.2. Segnet architecture

179 The Segnet model consists of encoder and corresponding decoder parts, followed by the last layer
180 of pixel-wise classification (Badrinarayanan *et al.* 2017). The overall architecture of deep
181 convolutional Segnet model is illustrated in Figure 3. Each layer in the encoder part has a
182 corresponding layer in the decoder part, and both sections include 13 convolutional layers that
183 correspond to the initial 13 convolutional layers in the model named the VGG16 network
184 (Simonyan and Zisserman 2014), which is outlined for feature classification. A multi-class
185 classifier named Softmax (Equation 4) is fed into the last decoder network to generate independent
186 classification possibilities for individual pixels. Softmax output transforms into possibility
187 dispensation as it always ranges [0–1] and adds up to 1. The n channel of image possibility is the
188 output of the Softmax classifier, where n presents the number of classes, x is the output vector of

189 the model, and index i is in the range $(0, \dots, n-1)$.
$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} . \quad (4)$$

190 For producing and batch normalizing a collection of feature maps, every encoder in the encoder
191 network implements a filter bank with a convolution. Subsequently, ReLU is utilized as an
192 activation function, followed by max-pooling layers with a kernel size of 2×2 . Using a factor of 2,
193 the outcoming output is sub-sampled. For achieving translation invariance over tiny spatial
194 changes in the input data, max-pooling layers are utilized. Although additional translation
195 invariance for strong classification can be obtained by multiple map-pooling layers, a
196 corresponding spatial resolution loss of feature maps occurs. Therefore, before implementing sub-
197 sampling, storing, and capturing, boundary information is essential in encoder feature maps. For
198 up-sampling the input feature maps in the decoder network, the memorized sub-sampling indices
199 from the corresponding encoder feature maps are utilized. Dense feature maps are produced by
200 convolving a trainable decoder filter back with these feature maps. Subsequently, a batch
201 normalization step is implemented to each map. The whole feature map in the Unet model
202 (Ronneberger *et al.* 2015) is first transferred to the corresponding decoders, and then is
203 concatenated to up-sample decoder feature maps (using deconvolution), whereas the Segnet model
204 reutilizes pooling indices. In addition, the Segnet model utilizes the whole weights of the pre-

205 trained convolutional layer from the VGG network as pre-trained weights, whereas no max-pool
206 5 block and conv 5 exist in the Unet model as in the architecture of the VGG network.

207

208 **Figure 3.** around here

209

210

211 **2.3 Seg-Unet architecture**

212 Similar to the Unet architecture, the SegUnet model comprises three sections (Do *et al.* 2019):
213 1) The encoder or contracting part, which is similar to the VGG network, has four blocks. In every
214 block, two convolution layers are followed by batch normalization and max-pooling layers. After
215 every max-pooling index, the number of features is doubled in the convolutional layer. 2)
216 Bottleneck, which only comprises two convolution layers, is a place for storing sparse feature
217 maps. 3) The decoder or expanding part restores the input image resolution by using up-sampling
218 layers. For transferring local contextual information into the decoder part, each encoder layer is
219 connected to the corresponding decoder layer. Unlike the Unet model, the same padding is utilized
220 instead of valid padding. For classifying each pixel and generating the segmentation map, a 1×1
221 convolution layer with sigmoid function is utilized at the last decoder block. The loss function of
222 binary cross-entropy is also applied to quantify the contrast between two possibility spreads and
223 assess the efficiency of the technique whose output value may be between 0 and 1. The over-fitting
224 issue can be prevented from using the new network because data normalization uses the batch
225 normalization layer, which is placed after the convolutional layer. Moreover, the sparse feature
226 map can be well restored using an up-sampling layer on the basis of the max-pooling index in the
227 decoder network. The overall framework of the proposed SegUnet deep neural network is
228 demonstrated in Figure 4.

229

230 **Figure 4.** around here

231 **2.4. Dataset**

232 To apply the proposed SegUnet model for building extraction, the Massachusetts building
 233 dataset (Mnih 2013) is used. Given the computational restriction, the original dataset that contains
 234 10, 4, and 137 aerial images for test, validation, and training with a spatial dimension of 1500×1500
 235 pixel dimension is divided into the size of 384×384, respectively. The total number of images used
 236 in this study is 1,564, where 1,532, 24, and 8 images are considered for training, validation, and
 237 test, respectively. Certain samples of a building dataset with various scenes are depicted in Figure
 238 5.

239

Figure 5. around here

240

241 **2.5. Evaluation metrics**

242 In this study, four principal calculation measurements, namely, overall accuracy (OA) (5), F1
 243 score (6), recall (7), and precision (8) are utilized on the basis of the confusion matrix
 244 (Ghasemkhani *et al.* 2020) with four main factors, such as false negative (FN), false positive (FP),
 245 true negative (TN), and true positive (TP), to assess the model performance for extracting building
 246 features from high-resolution aerial imagery. OA is specified as the sum of rightly identified pixels
 247 divided by the entire number of pixels. Precision is calculated as a percentage of precisely
 248 identified pixels among the identified pixels of the building. Meanwhile, F1 score is the
 249 combination of recall and precision metrics. Recall is specified as a percentage of correctly
 250 predicted pixels among all the actual pixels of building, whereas F1 score is the combination of
 251 recall and precision (Wang *et al.* 2020).

252
$$OA = \frac{TP + TN}{TP + TN + FP + FN} \quad (5)$$

253
$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall} \quad (6)$$

254
$$Recall = \frac{TP}{TP + FN} \quad (7)$$

255
$$Precision = \frac{TP}{TP + FP} \quad (8)$$

256

257 3. Results and performance evaluation

258 In this part, the quantitative and visual results of the proposed SegUnet model and other state-
259 of-the-art building extraction approaches, such as the Segnet model for semantic pixel-wise
260 segmentation, the FCN model for image semantic segmentation, and the deep convolutional Unet
261 model, are discussed.

262 3.1. Experiment results

263 By considering a representative section for the images with a specific attention on various
264 surroundings and building features, the visual inspection of classification maps achieved by the
265 proposed SegUnet model was implemented for qualitative analysis. For training the proposed model,
266 the ground truth labels, and all the prepared samples were treated as inputs to the model. The
267 parameters and framework of the proposed approach, such as the number of blocks and size of each
268 block, are illustrated in Figure 4. For updating the parameters of the proposed model and minimizing
269 the energy function while training the network, an exceptional optimization algorithm is needed.
270 Therefore, in our network, we utilized one of the most common optimizers called adaptive moment
271 estimation (Adam) to update parameters such as biases and weights and to lessen the losses. We set
272 the learning rate of the SegUnet network to $1e-4$ during training to speed up the processing and
273 achieve an improved performance. In this study, the whole process of the introduced network for
274 extracting building features from aerial imagery was performed on a GPU Nvidia Quadro P5000
275 with a computation capacity of 6.1 and a memory of 16 GB under the framework of Keras with
276 TensorFlow backend.

277 Figure 6 depicts the results of two images obtained by the proposed approach for building
278 extraction. The figure is presented in three columns and four rows. The first, second, and third
279 columns respectively represent the original image, the ground truth image, and the building
280 segmentation results, which were obtained by the SegUnet model. Meanwhile, the second and fourth
281 rows represent the zoomed results. Figure 6 shows that the proposed SegUnet model achieves the
282 OA of 92.33% and 91.3% for Image 1 and Image 2, respectively, proving that the model can
283 generally extract buildings from high-resolution aerial images accurately. However, the FN
284 (illustrated as blue pixel) and FP (illustrated as green pixel) of the identified pixels illustrate multiple
285 failures for our suggested approach and show multiple issues with the data. The proposed approach
286 can identify a building where tiny nearby buildings emerge as a joined area, which increases the FP

287 between the spaces of the building. However, the proposed model cannot make a right building
288 prediction where no building is found in the label image, but one exists in the original image that
289 appears as an FP prediction.

290

291 **Figure 6.** around here

292

293 **3.2. Discussion**

294 To verify the performance of the proposed SegUnet technique for extracting building objects
295 from high-resolution remotely sensing aerial imagery, we compared the method with other
296 DCNNs. Specifically, we compared the suggested SegUnet model with the deep convolutional
297 encoder–decoder approach called the Segnet model, FCN technique, and deep convolutional Unet
298 model. By comparing the results achieved via Segnet and Unet models with the results of the
299 proposed SegUnet model, the difference in the accuracy for building extraction can be witnessed.

300 The visual outcomes of building extraction by using the suggested SegUnet model and other
301 comparative techniques for calculating the efficiency of the SegUnet approach in building
302 extraction are illustrated in Figure 7. The obtained outcomes demonstrate that the influence of
303 shortcomings can be reduced to a specific degree by using the proposed methods because these
304 methods consider the spatial information for semantic segmentation. However, FCN and Segnet
305 approaches predict additional FNs and FPs, which are depicted by blue and green colors,
306 respectively. Thus, these methods cannot precisely preserve and achieve boundary information,
307 leading to the detection of FNs and FPs and production of a low-resolution segmentation map. The
308 Unet model, which utilizes deconvolution layers and skip connection, can also achieve and
309 preserve boundary information with higher accuracy than FCN and Segnet methods, thus obtaining
310 a correct segmentation map. By contrast, the proposed SegUnet model, which utilizes skip
311 connection (Unet) and index pooling (Segnet), can predict fewer FNs and FPs, preserve boundary
312 information, and produce a correct segmentation map compared to other comparative approaches.

313

314 **Figure 7.** around here

315 To test the efficiency of the introduced SegUnet approach for building extraction in comparison
316 with other DCNNs, we demonstrated the quantitative results of the techniques in Table 1. The first
317 eight rows of Table 1 present the quantitative accuracy of the four main metrics achieved by the
318 comparative approaches for the eight images, whereas the last row presents the average accuracy of
319 the metrics. As shown in Table 1, the FCN model can obtain higher accuracy for the recall factor
320 than other methods because the model predicts many FNs. By contrast, the Unet method can obtain
321 higher accuracy for precision and OA factors than FCN and Segnet methods. Moreover, the Unet
322 method is the second-best approach in building extraction and can obtain a correct segmentation
323 map. Finally, the average accuracy for F1 score and OA factors achieved by the proposed SegUnet
324 model is higher than those by other techniques with almost 0.14%, 1.17%, and 0.44% higher than
325 the Unet, Segnet, and FCN approaches, respectively. These results indicate that the proposed model
326 can improve the results and exceed other state-of-the-art techniques in building extraction from high-
327 resolution remote sensing imagery. Figure 8 plots the clear differences between the introduced
328 SegUnet model and other deep learning approaches for building object segmentation. Figure 8 also
329 illustrates that the proposed SegUnet network achieves higher precision for the OA factor than other
330 techniques.

331

332

Figure 8. around here

333

Table 1 around here

334 **4. Conclusion**

335 For extracting building objects from high-resolution aerial imagery, we presented a new deep
336 neural network called the SegUnet model, which is a combination of Segnet and Unet techniques, in
337 this work. We applied the proposed model on the Massachusetts building dataset. After training and
338 validating the method, we utilized four accuracy metrics to assess the efficiency of the indicated
339 technique in building extraction, which achieved a 92.73% accuracy on average for OA. This result
340 indicated that the proposed model can produce a correct segmentation map and can accurately extract
341 building objects. Furthermore, we compared the visual and quantitative results of the proposed
342 SegUnet model with those of other deep learning techniques, such as Segnet, FCN, and Unet models,
343 to show its effectiveness. The results confirmed that the proposed method obtained the best

344 quantitative and visual performances and outperformed other DCNNs in building extraction from
345 high-resolution aerial imagery.

346 **Author Contributions:** Conceptualization, A.A. and B.P.; methodology and formal analysis,
347 A.A.; data curation, A.A.; writing—original draft preparation, A.A.; writing—review and editing,
348 B.P.; supervision, B.P.; funding, B.P. and A.A.

349 **Funding:** This research is supported by the Centre for Advanced Modelling and Geospatial
350 Information Systems (CAMGIS), Faculty of Engineering and IT, the University of Technology
351 Sydney (UTS). This research is also supported by Researchers Supporting Project (RSP) number
352 RSP-2020/14, King Saud University, Riyadh, Saudi Arabia.

353 **Conflict of Interest:** The authors declare no conflict of interest.

354 **References**

- 355 1. Abdollahi, A., Bakhtiari, H.R.R. & Nejad, M.P., 2018. Investigation of svm and level set
356 interactive methods for road extraction from google earth images. *Journal of the Indian*
357 *Society of Remote Sensing*, 46 (3), 423-430.
- 358 2. Abdollahi, A., Pradhan, B., Shukla, N., Chakraborty, S. & Alamri, A., 2020. Deep learning
359 approaches applied to remote sensing datasets for road extraction: A state-of-the-art review.
360 *Remote Sensing*, (12), 1444.
- 361 3. Alshehhi, R., Marpu, P.R., Woon, W.L. & Dalla Mura, M., 2017. Simultaneous extraction of
362 roads and buildings in remote sensing imagery with convolutional neural networks. *ISPRS*
363 *Journal of Photogrammetry Remote Sensing*, 130, 139-149.
- 364 4. Audebert, N., Boulch, A., Lagrange, A., Le Saux, B. & Lefevre, S., 2016. Deep learning for
365 remote sensing. Technical Report. DOI: 10.1109/JURSE.2017.7924536.
- 366 5. Audebert, N., Boulch, A., Randrianarivo, H., Le Saux, B., Ferecatu, M., Lefèvre, S. & Marlet,
367 R., 2017. Deep learning for urban remote sensing. *Joint Urban Remote Sensing Event*
368 *(JURSE)IEEE*, 1-4. DOI: 10.1109/JURSE.2017.7924536.
- 369 6. Audebert, N., Le Saux, B. & Lefèvre, S., 2017. Semantic segmentation of earth observation
370 data using multimodal and multi-scale deep networks. *Asian Conference on Computer Vision*,
371 180-196. DOI: 10.1007/978-3-319-54181-5_12.
- 372 7. Badrinarayanan, V., Kendall, A. & Cipolla, R., 2017. Segnet: A deep convolutional encoder-
373 decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis*
374 *Machine Intelligence*, 39 (12), 2481-2495.
- 375 8. Bakhtiari, H.R.R., Abdollahi, A. & Rezaeian, H., 2017. Semi automatic road extraction from
376 digital images. *The Egyptian Journal of Remote Sensing and Space Science*, 20 (1), 117-123
377 Available from: <http://www.sciencedirect.com/science/article/pii/S1110982317300820>.

- 378 9. Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K. & Yuille, A.L., 2014. Semantic image
379 segmentation with deep convolutional nets and fully connected crfs. 834 - 848. DOI:
380 10.1109/TPAMI.2017.2699184.
- 381 10. Chen, L.-C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.J.I.T.O.P.A. &
382 Intelligence, M., 2017. Deeplab: Semantic image segmentation with deep convolutional nets,
383 atrous convolution, and fully connected crfs. 40 (4), 834-848.
- 384 11. Do, N.-T., Joo, S.-D., Yang, H.-J., Jung, S.T. & Kim, S.-H., 2019. Knee bone tumor
385 segmentation from radiographs using seg-unet with dice loss. *25th International Workshop*
386 *on Frontiers of Computer Vision (IW-FCV), Gangneung, South Korea.*
- 387 12. Farabet, C., Couprie, C., Najman, L. & Lecun, Y., 2012. Learning hierarchical features for
388 scene labeling. *IEEE Transactions on Pattern Analysis Machine Intelligence*, 35 (8), 1915-
389 1929.
- 390 13. Fu, G., Liu, C., Zhou, R., Sun, T. & Zhang, Q., 2017. Classification for high resolution remote
391 sensing imagery using a fully convolutional network. *Remote Sensing*, 9 (5), 498.
- 392 14. Ghasemkhani, N., Vayghan, S.S., Abdollahi, A., Pradhan, B. & Alamri, A., 2020. Urban
393 development modeling using integrated fuzzy systems, ordered weighted averaging (owa),
394 and geospatial techniques. *Sustainability*, 12 (3), 809.
- 395 15. Girshick, R., Donahue, J., Darrell, T. & Malik, J., 2014. Rich feature hierarchies for accurate
396 object detection and semantic segmentation. *Proceedings of the IEEE Conference on*
397 *Computer Vision and Pattern Recognition*, 580-587. Available from:
398 <https://arxiv.org/abs/1311.2524>.
- 399 16. He, K., Zhang, X., Ren, S. & Sun, J., 2015. Delving deep into rectifiers: Surpassing human-
400 level performance on imagenet classification. *Proceedings of the IEEE International*
401 *Conference on Computer Vision*, 1026-1034. Available from:
402 <https://arxiv.org/abs/1502.01852>.
- 403 17. Hu, F., Xia, G.-S., Hu, J. & Zhang, L., 2015. Transferring deep convolutional neural networks
404 for the scene classification of high-resolution remote sensing imagery. *Remote Sensing*, 7 (11),
405 14680-14707.
- 406 18. Huertas, A. & Nevatia, R., 1988. Detecting buildings in aerial images. *Computer Vision,*
407 *Graphics, Image Processing*, 41 (2), 131-152.
- 408 19. Inglada, J., 2007. Automatic recognition of man-made objects in high resolution optical
409 remote sensing images by svm classification of geometric image features. *ISPRS Journal of*
410 *Photogrammetry Remote Sensing*, 62 (3), 236-248.
- 411 20. Krizhevsky, A., Sutskever, I. & Hinton, G.E., 2012. Imagenet classification with deep
412 convolutional neural networks. *Advances in Neural Information Processing Systems*, 1097-
413 1105. DOI: 10.1145/3065386.
- 414 21. Kussul, N., Shelestov, A., Lavreniuk, M., Butko, I. & Skakun, S., 2016. Deep learning
415 approach for large scale land cover mapping based on remote sensing data fusion. *IEEE*
416 *International Geoscience and Remote Sensing Symposium (IGARSS)*, 198-201. DOI:
417 10.1109/IGARSS.2016.7729043.

- 418 22. Levitt, S. & Aghdasi, F., 1998. An investigation into the use of wavelets and scaling for the
419 extraction of buildings in aerial imaged. *Proceedings of the 1998 South African Symposium*
420 *on Communications and Signal Processing-COMSIG'98 (Cat. No. 98EX214)*, 133-138. DOI:
421 10.1109/COMSIG.1998.736936.
- 422 23. Long, J., Shelhamer, E. & Darrell, T., 2016. Fully convolutional networks for semantic
423 segmentation. *Proceedings of the IEEE Conference on Computer Vision and Pattern*
424 *Recognition*, 3431-3440. DOI: 10.1109/TPAMI.2016.2572683.
- 425 24. Maggiori, E., Tarabalka, Y., Charpiat, G. & Alliez, P., 2017. Convolutional neural networks
426 for large-scale remote-sensing image classification. *IEEE Transactions on Geoscience*
427 *Remote Sensing*, 55 (2), 645-657.
- 428 25. Marcu, A. & Leordeanu, M., 2016. Dual local-global contextual pathways for recognition in
429 aerial imagery. Available from: <https://arxiv.org/abs/1605.05462>.
- 430 26. Marmanis, D., Schindler, K., Wegner, J.D., Galliani, S., Datcu, M. & Stilla, U., 2018.
431 Classification with an edge: Improving semantic image segmentation with boundary
432 detection. *ISPRS Journal of Photogrammetry Remote Sensing*, 135, 158-172.
- 433 27. Marmanis, D., Wegner, J.D., Galliani, S., Schindler, K., Datcu, M. & Stilla, U., 2016.
434 Semantic segmentation of aerial images with an ensemble of cnns. *ISPRS Annals of the*
435 *Photogrammetry, Remote Sensing Spatial Information Sciences*, 3, 473-480.
- 436 28. Mayer, H., 1999. Automatic object extraction from aerial imagery—a survey focusing on
437 buildings. *Computer Vision Image Understanding*, 74 (2), 138-149.
- 438 29. Mnih, V., 2013. *Machine learning for aerial image labeling, ph.D. Dissertation, dept.*
439 *Comput. Sci., univ. Toronto, Canada*.
- 440 30. Paisitkriangkrai, S., Sherrah, J., Janney, P. & Van Den Hengel, A., 2016. Semantic labeling
441 of aerial and satellite imagery. *IEEE Journal of Selected Topics in Applied Earth Observations*
442 *Remote Sensing*, 9(7), 2868-2881.
- 443 31. Penatti, O.A., Nogueira, K. & Dos Santos, J.A., 2015. Do deep features generalize from
444 everyday objects to remote sensing and aerial scenes domains? *Proceedings of the IEEE*
445 *Conference on Computer Vision and Pattern Recognition Workshops*, 44-51. DOI:
446 10.1109/CVPRW.2015.7301382.
- 447 32. Peng, J. & Liu, Y., 2005. Model and context-driven building extraction in dense urban aerial
448 images. *International Journal of Remote Sensing*, 26 (7), 1289-1307.
- 449 33. Ronneberger, O., Fischer, P. & Brox, T., 2015. U-net: Convolutional networks for biomedical
450 image segmentation. *International Conference on Medical Image Computing and*
451 *Computer-assisted Intervention*, 234-241. DOI: 10.1007/978-3-319-24574-4_28.
- 452 34. Sharif Razavian, A., Azizpour, H., Sullivan, J. & Carlsson, S., 2015. Cnn features off-the-
453 shelf: An astounding baseline for recognition. *Proceedings of the IEEE Conference on*
454 *Computer Vision and Pattern Recognition Workshops*, 806-813. DOI:
455 10.1109/CVPRW.2014.131.
- 456 35. Sherrah, J., 2016. Fully convolutional networks for dense semantic labelling of high-
457 resolution aerial imagery. Available from: <https://arxiv.org/abs/1606.02585>.

- 458 36. Shrestha, S. & Vanneschi, L., 2018. Improved fully convolutional network with conditional
459 random fields for building extraction. *Remote Sensing*, 10 (7), 1135.
- 460 37. Simonyan, K. & Zisserman, A., 2014. Very deep convolutional networks for large-scale image
461 recognition. 1-14. Available from: <https://arxiv.org/abs/1409.1556>.
- 462 38. Sumer, E. & Turker, M., 2013. An adaptive fuzzy-genetic algorithm approach for building
463 detection using high-resolution satellite images. *Computers, Environment Urban Systems*, 39,
464 48-62.
- 465 39. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke,
466 V. & Rabinovich, A., 2015. Going deeper with convolutions. *Proceedings of the IEEE
467 Conference on Computer Vision and Pattern Recognition*, 1-9. DOI:
468 10.1109/CVPR.2015.7298594.
- 469 40. Vakalopoulou, M., Karantza, K., Komodakis, N. & Paragios, N., 2015. Building detection
470 in very high resolution multispectral data with deep learning features. *IEEE International
471 Geoscience and Remote Sensing Symposium (IGARSS)IEEE*, 1873-1876. DOI:
472 10.1109/IGARSS.2015.7326158.
- 473 41. Volpi, M. & Tuia, D., 2016. Dense semantic labeling of subdecimeter resolution images with
474 convolutional neural networks. *IEEE Transactions on Geoscience Remote Sensing*, 55 (2),
475 881-893.
- 476 42. Wang, S., Hou, X. & Zhao, X., 2020. Automatic building extraction from high-resolution
477 aerial imagery via fully convolutional encoder-decoder network with non-local block. *IEEE
478 Access*, 8, 7313-7322.
- 479 43. Wilkinson, G.G., 2005. Results and implications of a study of fifteen years of satellite image
480 classification experiments. *IEEE Transactions on Geoscience Remote Sensing*, 43 (3), 433-
481 440.
- 482 44. Yang, X., Ye, Y., Li, X., Lau, R.Y., Zhang, X. & Huang, X., 2018. Hyperspectral image
483 classification with deep learning models. *IEEE Transactions on Geoscience Remote Sensing*,
484 56 (9), 5408-5423.
- 485 45. Yuan, J., 2017. Learning building extraction in aerial scenes with convolutional networks.
486 *IEEE Transactions on Pattern Analysis Machine Intelligence*, 40 (11), 2793-2798.
- 487 46. Yuan, J. & Cheryadat, A.M., 2014. Learning to count buildings in diverse aerial scenes.
488 *Proceedings of the 22nd ACM SIGSPATIAL International Conference on Advances in
489 Geographic Information Systems*, 271-280. DOI: 10.1145/2666310.2666389.
- 490 47. Zhou, B., Lapedriza, A., Xiao, J., Torralba, A. & Oliva, A., 2014. Learning deep features for
491 scene recognition using places database. *27th International Conference on Neural
492 Information Processing Systems, Montreal, Canada, 1*, 487-495.
- 493
494