

# A hybrid MLP-CNN classifier for very fine resolution remotely sensed image classification

Ce Zhang<sup>1,\*</sup>, Xin Pan<sup>2,3</sup>, Huapeng Li<sup>2</sup>, Andy Gardiner<sup>4</sup>, Isabel Sargent<sup>4</sup>, Jonathon Hare<sup>5</sup>, Peter M. Atkinson<sup>1,\*</sup>

<sup>1</sup> Lancaster Environment Centre, Lancaster University, Lancaster LA1 4YQ, UK; <sup>2</sup> Northeast Institute of Geography and Agroecology, Chinese Academic of Science, Changchun 130102, China; <sup>3</sup> School of Computer Technology and Engineering, Changchun Institute of Technology, 130021 Changchun, China; <sup>4</sup> Ordnance Survey, Adanac Drive, Southampton SO16 0AS, UK; <sup>5</sup> Electronics and Computer Science (ECS), University of Southampton, Southampton SO17 1BJ, UK

**Abstract** The contextual-based convolutional neural network (CNN) with deep architecture and pixel-based multilayer perceptron (MLP) with shallow structure are well-recognized neural network algorithms, representing the state-of-the-art deep learning method and the classical non-parametric machine learning approach, respectively. The two algorithms, which have very different behaviours, were integrated in a concise and effective way using a rule-based decision fusion approach for the classification of very fine spatial resolution (VFSR) remotely sensed imagery. The decision fusion rules, designed primarily based on the classification confidence of the CNN, reflect the generally complementary patterns of the individual classifiers. In consequence, the proposed ensemble classifier MLP-CNN harvests the complementary results acquired from the CNN based on deep spatial feature representation and from the MLP based on spectral discrimination. Meanwhile, limitations of the CNN due to the adoption of convolutional filters such as the uncertainty in object boundary partition and loss of useful fine spatial resolution detail were compensated. The effectiveness of the ensemble MLP-CNN classifier was tested in both urban and rural areas using aerial photography together with an additional satellite sensor dataset. The MLP-CNN classifier achieved promising performance, consistently outperforming the pixel-based MLP, spectral and textural-based MLP, and the contextual-based CNN in terms of classification accuracy. This research paves the way to effectively address the complicated problem of VFSR image classification.

**Keywords:** convolutional neural network; multilayer perceptron; VFSR remotely sensed imagery; fusion decision; feature representation

## 33 1. Introduction

34 With the rapid development of modern remote sensing technologies, a large quantity of  
35 very fine spatial resolution (VFSR) images is now commercially available. These  
36 VFSR images, typically acquired at sub-metre spatial resolution, have opened up many  
37 opportunities for new applications (Zhong et al., 2014), for example, urban land use  
38 retrieval (Mathieu et al., 2007; Shi et al., 2015), precision agriculture (Ozdarici-Ok et  
39 al., 2015; Zhang and Kovacs, 2012), and tree crown delineation (Ardila et al., 2011;  
40 Yin et al., 2015). However, despite the presence of a rich spatial data content (Huang  
41 et al., 2014), the information conveyed by the imagery is conditional upon the quality  
42 of the processing (L ängkvist et al., 2016). With fewer spectral channels in comparison  
43 with coarse or medium spatial resolution remotely sensed data, it can be challenging to  
44 differentiate subtle differences amongst similar land cover types (Powers et al., 2015).  
45 Meanwhile, objects of the same class may exhibit strong spectral heterogeneity due to  
46 differences in age, level of maintenance and composition as well as illumination  
47 conditions (Demarchi et al., 2014), which might be further complicated by the  
48 scattering of peripheral ground objects (Chen et al., 2014). As a consequence, such high  
49 intra-class variability and low inter-class disparity make automatic classification of  
50 VFSR images a challenging task.

51 Ever since the advent of VFSR imagery, tremendous efforts have been made to develop  
52 robust and accurate, automatic image classification methods. Among these, machine  
53 learning is currently considered as the most promising and evolving approach (Zhang  
54 et al., 2015). Popular pixel-based machine learning algorithms, such as Multilayer  
55 Perceptron (MLP), Support Vector Machine (SVM) and Random Forest (RF), have  
56 drawn considerable attention in the remote sensing community (Attarchi and Gloaguen,  
57 2014; Yang et al., 2012; Zhang et al., 2015). The MLP, as a typical non-parametric  
58 neural network classifier, is designed to learn the nonlinear spectral feature space at the  
59 pixel level irrespective of its statistical properties (Atkinson and Tatnall, 1997; Foody  
60 and Arora, 1997; Mas and Flores, 2008). The MLP has been used widely in remote  
61 sensing applications, including VFSR-based land cover classification (e.g. Del Frate et  
62 al., (2007), Pacifici et al. (2009)). The MLP algorithm is mathematically complicated  
63 yet can be simple in model architecture (e.g., a shallow classifier with one or two feature  
64 representation levels). At the same time, a pixel-based MLP classifier does not consider,  
65 or make use of, the spatial patterns implicit in images, especially for VFSR imagery

66 with unprecedented spatial detail. In essence, the MLP (and related algorithms, e.g.  
67 SVM, RF, etc.) is a pixel-based classifier with shallow structure (one or two layers)  
68 (Chen et al., 2016), where the membership association of a pixel for each class is  
69 predicted.

70 Recent advances in neuroscience have shown that deep feature representations can be  
71 learned hierarchically from simple concepts such as oriented edges to higher-level  
72 complex patterns such as textures, segments, parts and objects (Arel et al., 2010). This  
73 discovery motivated the breakthrough of the so-called “deep learning” methods that  
74 represent the state-of-the-art in a variety of domains, including target detection (Chen  
75 et al., 2016; Tang et al., 2015), image recognition (Farabet et al., 2013; Krizhevsky et  
76 al., 2012) and robotics (Bezak et al., 2014; Lenz et al., 2015; Yu et al., 2013), amongst  
77 others. The convolutional neural network (CNN), a well-established deep learning  
78 approach, has produced excellent results in the field of computer vision and pattern  
79 recognition (Schmidhuber, 2015), such as for visual recognition (Farabet et al., 2013;  
80 Krizhevsky et al., 2012), image retrieval (X. Yang et al., 2015) and scene annotation  
81 (Othman et al., 2016).

82 In the remote sensing domain, CNNs have been studied actively and shown to produce  
83 state-of-the-art results over the past few years, focusing primarily on object detection  
84 (Dong et al., 2015) or scene classification (Hu et al., 2015a; Zhang et al., 2016). Recent  
85 studies further explored the feasibility of CNNs for the task of remotely sensed image  
86 classification. For example, Yue et al., (2016) utilized spatial pyramid pooling to learn  
87 multi-scale spatial features from hyperspectral data, Chen et al. (2016) introduced a 3D  
88 CNN to jointly extract spectral–spatial features, thus, making full use of the continuous  
89 hyperspectral and spatial spaces. In terms of the classification of multi- and  
90 hyperspectral imagery, a deep CNN model was formulated through a greedy layer-wise  
91 unsupervised pre-training strategy (Romero et al., 2016), whereas an image pyramid  
92 was built up through upscaling the original image to capture the contextual information  
93 at multiple scales (Zhao and Du, 2016). For VFSR image classification, CNN models  
94 with varying contextual input size were constructed to learn multi-scale features while  
95 preserving the original fine resolution information (L ängkvist et al., 2016). All of the  
96 above-mentioned work applied CNNs with contextual patches as their inputs, and  
97 demonstrated the robustness and effectiveness in spatial feature representations with  
98 excellent classification performance. However, the benefits and shortcomings of the

99 CNN as a classifier itself have not been studied thoroughly. In particular, the CNN, as  
100 a contextual classifier with deep structures (Szegedy et al., 2015), explores the complex  
101 spatial patterns hidden in the image that are not seen by representation in its shallow  
102 counterparts, whereas it may overlook certain information in spectral space observed  
103 by pixel-based classifiers. Moreover, uncertainties may appear in object boundaries due  
104 to the usage of convolutional filters of the CNN. These issues deserve further  
105 investigation.

106 Any single set of features (e.g., spectral only) or a specific classifier (e.g., pixel-based  
107 only) is unlikely to achieve the highest classification accuracies for VFSR imagery  
108 because the result is conditional upon both spectral and spatial information. In this  
109 context, two categories of spectral and spatial information were fused for classification  
110 or handled with a classifier ensemble. Information fusion can be realized by stacking  
111 the spatial and spectral information as feature bands. However, this does not allow the  
112 specification of the relative influence of the extracted features (Wang et al., 2016).  
113 Others proposed integrative algorithms considering the spatial and spectral features at  
114 the same time. For example, Fauvel et al., (2012) proposed a composite kernel-based  
115 SVM with spectral and spatial kernels applied simultaneously. However, the spatial  
116 kernel summarizes only basic information (e.g. median) of the spatial neighbourhood  
117 (Wang et al., 2016).

118 In terms of classifier ensemble technology, two strategies, namely “multiple classifier  
119 systems” (Benediktsson, 2009) and “decision fusion” (Fauvel et al., 2006) are  
120 employed. Multiple classifier systems are based on the manipulation of training sample  
121 sets, including boosting (Freund et al., 2003) and bagging (Breiman, 1996). This  
122 ensemble approach, however, usually requires a relatively large sample size and the  
123 computational complexity tends to be high. An alternative classifier ensemble is  
124 derived from decision fusion of the outputs of different classification algorithms  
125 according to a certain combination of approaches (Du et al., 2012; Löw et al., 2015).  
126 This decision fusion-based ensemble approach is preferable where the individual  
127 classifiers demonstrate complementary behaviour. For instance, different non-  
128 parametric classifiers are sometimes accurate in different locations in a classification  
129 map, thus, producing complementary results from the ensemble (Clinton et al., 2015;  
130 Löw et al., 2015). However, all the aforementioned fusion strategies are conducted

131 using pixel-based classifiers with shallow structures, whose complementary behaviours  
132 are insufficient to address the challenges of VFSR image classification.

133 In this paper, a hybrid classification system was proposed that combines the CNN (a  
134 contextual-based classifier with deep architectures) and MLP (a pixel-based classifier  
135 with shallow structures) using a rule-based decision fusion strategy. The hypothesis is  
136 that both MLP and CNN classifiers can provide different views or feature  
137 representations with strong complementarity. Thus, the classifier ensemble has the  
138 potential to enhance the final classification performance. The decision fusion rules were  
139 built up at the post-classification stage, primarily based on the confidence distribution  
140 of the contextual-based CNN classifier, such that the classified pixels with low  
141 confidence can be rectified by the MLP at the pixel level. The effectiveness of the  
142 proposed method was tested on images of both an urban scene and a rural area. A  
143 benchmark comparison was provided by the standard pixel-based MLP, spectral-  
144 texture based MLP as well as contextual-based CNN classifiers.

## 145 **2. Methodology**

### 146 *2.1 Brief review of multilayer perceptrons (MLP)*

147 A multilayer perceptron (MLP) is a network that maps sets of input data onto a set of  
148 outputs in a feedforward manner (Atkinson and Tatnall, 1997). The typical structure is  
149 that the MLP is composed of interconnected nodes in multiple layers (input, hidden and  
150 output layers), with each layer fully connected to the preceding layer as well as the  
151 succeeding layer (Del Frate et al., 2007). The outputs of each node are weighted units  
152 followed by a nonlinear activation function to distinguish the data that are not linearly  
153 separable (Pacifici et al., 2009). Formally, the output activation  $a^{(l+1)}$  at layer  $l+1$  is  
154 derived by the input activation  $a^{(l)}$ :

$$155 \quad a^{(l+1)} = \sigma(w^{(l)} a^{(l)} + b^{(l)}) \quad (1)$$

156 Where  $l$  corresponds to a specific layer,  $w^{(l)}$  and  $b^{(l)}$  denote the weight and bias at layer  
157  $l$ , and  $\sigma$  represents the nonlinear activation operation (e.g. sigmoid, hyperbolic tangent,  
158 rectified linear units) function. For an  $m$  layer multilayer perceptron, the first input layer  
159 is  $a^{(1)} = x$  while the last output layer is:

160 
$$h_{w,b}(x) = a^{(m)} \quad (2)$$

161 The weights  $w$  and bias  $b$  in equation (2) are learned by supervised training using a  
 162 backpropagation algorithm to approximate an unknown input-output relation (Del Frate  
 163 et al., 2007). The objective function is to minimize the difference between the predicted  
 164 outputs and the desired outputs:

165 
$$J(W, b; x, y) = \frac{1}{2} \|h_{w,b}(x) - y\|^2 \quad (3)$$

166 **2.2 Brief review of Convolutional Neural Networks (CNN)**

167 A Convolutional Neural Network (CNN), is a variant of the multilayer feed forward  
 168 neural networks, and is designed specifically to process large scale images or sensory  
 169 data in the form of multiple arrays by considering local and global stationary properties  
 170 (LeCun et al., 2015). Similar to the MLP, the CNN is a network stacked into a number  
 171 of *layers*, where the output of the previous layer is connected sequentially to the input  
 172 of the next one by a set of learnable weights and biases (Romero et al., 2016). The major  
 173 difference is that each layer is represented as input and output feature maps by capturing  
 174 different perspectives on features through the operation of convolution.

175 The CNN basically consists of three major operations: convolution, nonlinearity and  
 176 pooling/subsampling (Schmidhuber, 2015). The convolutional and pooling layers are  
 177 stacked together alternatively in the CNN framework, until obtaining the high-level  
 178 features on which a fully connected classification is performed (LeCun et al., 2015). In  
 179 addition, several feature maps may exist in each convolutional layer and the weights of  
 180 convolutional nodes in the same map are shared. This setting enables the network to  
 181 learn different features while keeping the number of parameters tractable.  
 182 Mathematically, the output feature map  $y_{i,j}^{(l)}$  at convolutional layer  $l$  is calculated as:

183 
$$y_{i,j}^{(l)} = \sigma^{(l)} \left( \sum_{n=1}^k \sum_{m=1}^k w_{n,m}^{(l)} \cdot x_{i+n,j+m}^{(l-1)} + b^{(l)} \right) \quad (4)$$

184 Where the  $w_{n,m}^{(l)}$  denotes the convolutional filter with size  $k \times k$  at layer  $l$ , and the  
 185  $x_{i+n,j+m}^{(l-1)}$  represents the spatial position of the corresponding feature map at the  
 186 preceding layer  $l-1$ . The algorithm passes the convolutional filter throughout the input

187 feature map using the dot product  $(\cdot)$  between them with an addition of a bias unit  $b^{(l)}$   
 188 (Arel et al., 2010). Moreover, a nonlinear activation function  $\sigma^{(l)}$  at layer  $l$  is taken  
 189 outside the dot product to strengthen the nonlinearity (Strigl et al., 2010).

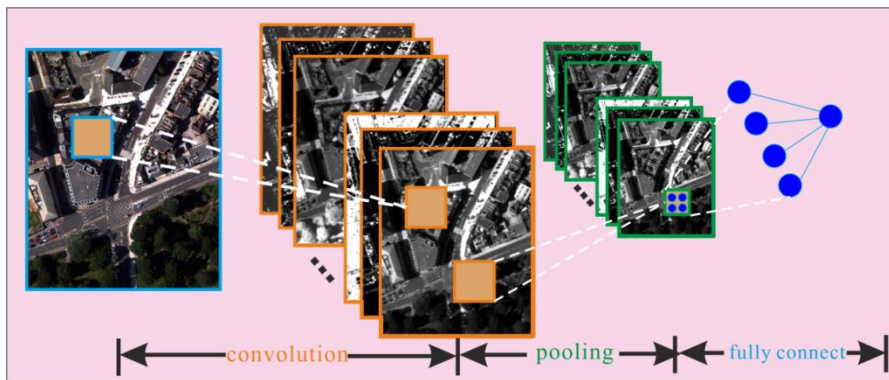
190 The pooling/subsampling layer can generalize the convolved features through down-  
 191 sampling and thereby reduce the computational complexity during the training process  
 192 (Zhao and Du, 2016). Given a pooling/subsampling layer  $q$ , the feature output  $F^q$  can  
 193 be derived from the preceding layer  $f^{(q-1)}$  through

$$194 \quad F_{i,j}^q = \max(f_{1+p(i-1),1+p(j-1)}^{q-1}, \dots, f_{pi,1+p(j-1)}^{q-1}, \dots, f_{1+p(i-1),pj}^{q-1}, \dots, f_{pi,pj}^{q-1})$$

195 (5)

196 Where  $p \times p$  is the size of the local spatial region, and  $1 \leq i, j \leq (m - n + 1) / p$ , here the  
 197  $m$  refers to the size of input feature map, while  $n$  corresponds to the size of filter  
 198 (Längkvist et al., 2016). The *max* simply summarizes the input features within local  
 199 spatial region using the maximum value (Figure 1: Pooling). By doing this, the learnt  
 200 features become robust and abstract with certain sparseness and translation invariance.

201 Once the higher level features are extracted, the output feature maps are flattened into  
 202 a one-dimensional vector, followed by a fully connected output layer (Figure 1: fully  
 203 connect). This operation is exactly a simple logistic regression, which is equivalent to  
 204 the standard MLP discussed in section 2.1, but without any hidden layer.



205  
 206 Figure 1 A schematic illustration of the three core layers within the CNN architecture, including the  
 207 convolutional layer (convolution), pooling layer (pooling) and fully connected layer (fully connect).

208 **2.3 Hybrid MLP-CNN Classification Approach**

209 Suppose the predictive outputs of the MLP and CNN at each pixel are  $n$ -dimensional  
 210 vectors  $C = (c_1, c_2, \dots, c_n)$ , where  $n$  represents the number of classes and each dimension  
 211  $i \in [1, n]$  corresponds to the probability of a specific class ( $i$ -th class) with certain  
 212 membership association. Ideally, the probability of the classification prediction would  
 213 be 1 for the target class and 0 for the others. However, due to the uncertainty in the  
 214 process of remotely sensed image classification, the probability value  $c$  is denoted as  
 215  $f(x) = \{c_x \mid x \in (1, 2, \dots, n)\}$ , where  $c_x \in [0, 1]$  and  $\sum_1^n c_x = 1$ . The classification model  
 216 simply takes the maximum membership association as the predicted output label  
 217 (denoted as  $\text{class}(C)$ ):

$$218 \quad \text{class}(C) = \arg \max(\{f(x) = c_x \mid x \in (1, 2, \dots, n)\}) \quad (6)$$

219 The confidence  $\text{conf}$  of such membership association is defined here as:

$$220 \quad \text{conf} = \text{Max}(C) - \text{Mean}(C) \quad (7)$$

221 In equation (7),  $\text{Max}(C)$  represents the maximum value of vector  $C$ , while  $\text{Mean}(C)$   
 222 denotes the average of all the values of  $C$ . The  $\text{conf}$ , quantified by the difference  
 223 between  $\text{Max}(C)$  and  $\text{Mean}(C)$ , measures the confidence or reliability of the class  
 224 membership allocation (i.e. classification confidence map). Since the CNN takes  
 225 contextual image patches as its inputs instead of image pixels, it has the following  
 226 properties:

227 (1). If the input image patch is located at the central homogeneous region, its class  
 228 purity is relatively high with large difference between the membership association of  
 229 the predicted class and those of the other classes, and the  $\text{conf}$  tends to be large (White  
 230 regions in Figure 2(c)).

231 (2). If the image patch contains other land cover classes as contextual information, the  
 232 resulting distinction between the membership association of prediction and those of the  
 233 others is relatively low, and the  $\text{conf}$  tends to be small (Dark regions in Figure 2(c)).

234 However, the MLP (spectral feature only) is based on per-pixel spectral information,  
 235 thereby ruling out the difference of membership association between central and  
 236 boundary regions of the classified objects (Figure 1(b)). According to the

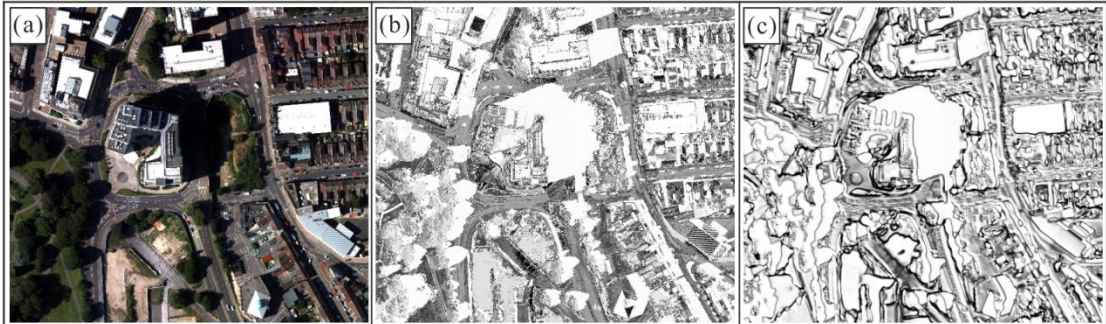


237 aforementioned properties, the fusion decision rules are constructed primarily based on  
 238 CNN confidence. To be more specific, the fusion output gives credit to the CNN when  
 239 its confidence is larger than a predefined threshold ( $\alpha_1$ ), while the MLP is trusted given  
 240 that the CNN confidence is lower than another threshold ( $\alpha_2$ ); once the confidence of  
 241 the CNN lies in-between the two thresholds ( $\in (\alpha_1, \alpha_2)$ ), the fusion output chooses the  
 242 CNN or MLP classification result with a larger confidence. Therefore, for a given image  
 243 pixel at location  $(h, v)$ , a rule-based decision fusion approach to determining the class  
 244 label ( $class(h, v)$ ) of the corresponding pixel is formulated as follows:

$$245 \quad class(h, v) = \left. \begin{cases} class_{mlp} & conf_{cnn} < \alpha_1 \\ class_{mlp} & (\alpha_1 \leq conf_{cnn} < \alpha_2 \ \& \ conf_{cnn} < conf_{mlp}) \\ class_{cnn} & (\alpha_1 \leq conf_{cnn} < \alpha_2 \ \& \ conf_{cnn} > conf_{mlp}) \\ class_{cnn} & conf_{cnn} \geq \alpha_2 \end{cases} \right\} (8)$$

246 Where the  $class_{mlp}$  and  $class_{cnn}$  represent the classification results of the MLP and CNN  
 247 respectively; the  $conf_{mlp}$  and  $conf_{cnn}$  denote the classification confidence of the MLP  
 248 and CNN accordingly.

249 Estimation of the two thresholds ( $\alpha_1, \alpha_2$ ) is conducted using grid search with cross-  
 250 validation (Min and Lee, 2005; Zhang et al., 2015) based on the CNN classification  
 251 confidence map (as illustrated by Figure 2(c)). Specifically, the  $\alpha_1$  was searched from  
 252 0.1 to 0.5 to detect those regions with low confidence as predicted by the CNN, while  
 253 the  $\alpha_2$  was chosen from 0.5 to 0.9 to discover the high confidence regions. By initially  
 254 fixing  $\alpha_1$  as 0.1,  $\alpha_2$  was tuned with step size of 0.05 (i.e.  $\alpha_2=0.5, 0.55, 0.6, \dots, 0.9$ ) to  
 255 cross-validate the classification accuracy influenced by the selected thresholds;  $\alpha_1$  was  
 256 then increased to further tune  $\alpha_2$  in a similar way until the optimal  $\alpha_1$  and  $\alpha_2$  were found  
 257 with the best classification accuracy.



258

259 Figure 2 (a) A subset of the original imagery with RGB spectral bands, (b) the classification confidence  
260 of the MLP and (c) the classification confidence of the CNN. The dark pixels represent low confidence,  
261 while white pixels signify high confidence.

### 262 3. Experiment

#### 263 3.1 Study area and data source

264 For this study, the city of Southampton, UK and its surrounding environment, which  
265 lies on the south coast of England, was chosen as a case study area (Figure 3). The  
266 urban and suburban areas in Southampton are strongly heterogeneous with a mixture of  
267 anthropogenic urban surface (e.g. roof materials, asphalt, concrete) and semi-natural  
268 environment (e.g. vegetation, bare soil), thereby representing a good test for  
269 classification algorithms.

270 A scene of aerial imagery of Southampton was captured on 22 July 2012 using a Vexcel  
271 UltraCam Xp digital aerial camera with 50 cm spatial resolution and four multispectral  
272 bands (Red, Green, Blue and Near Infrared). Two study sites S1 (3087×2750 pixels)  
273 and S2 (2022×1672 pixels) were selected to investigate the effectiveness of the  
274 proposed algorithm. S1 is located in the city centre of Southampton, which consists of  
275 eight dominant land cover classes, including Clay roof, Concrete roof, Metal roof,  
276 Asphalt, Grassland, Trees, Bare soil and Shadow, with detailed descriptions listed in  
277 Table 1. S2, on the other hand, is situated in a suburban and rural area of Southampton  
278 comprised of large patches of forest, grassland and bare soil speckled with small  
279 buildings and roads. There are six land cover categories in this study site, namely,  
280 Buildings, Road-or-track, Grassland, Trees, Bare soil and Shadow (Table 1).



281

282 Figure 3 Southampton, UK Location of study area and aerial imagery with two study sites S1 and S2.

283 Sample points were collected using a stratified random scheme from ground data  
 284 provided by local surveyors at Southampton, and split into 50% training samples and  
 285 50% testing samples for each class (Table 1). Field land cover survey was conducted  
 286 throughout the study area on July 2012 to further check the validity and precision of  
 287 the selected samples. In addition, a highly detailed vector map from Ordnance Survey,  
 288 namely the MasterMap Topographic Layer (Regnauld and Mackaness, 2006), was fully  
 289 consulted and cross-referenced to gain a comprehensive appreciation of the land cover  
 290 and land use within the study area.

291 Table 1 Detailed description of land covers at two study sites with training and testing sample size per  
 292 class.

Study Sites	Class	Train	Test	Description
S1	Clay roof	144	144	Predominantly residential buildings in red clay tiles
	Concrete roof	132	132	Predominantly residential buildings in grey clay tiles
	Metal roof	134	134	Predominantly industrial buildings in white metal panels
	Asphalt	136	136	Urban road or car parks covered by asphalt
	Grassland	126	126	Areas of grass covering the urban park or lawn
	Trees	137	137	Patches of tree species
	Bare soil	118	118	Open areas covered by bare soil
	Shadow	123	123	Areas of shadow cast from buildings and trees
S2	Building	82	82	Predominantly small buildings at rural areas
	Road-or-track	85	85	Asphalt road or small path
	Grassland	86	86	Large areas of wild grass or lawn
	Trees	98	98	Large patches of deciduous trees
	Bare soil	84	84	Open areas covered by bare soil
	Shadow	86	86	Areas of shadow cast from buildings and trees

293

294 To further test the applicability of the proposed method, another scene of Worldview-  
 295 2 satellite sensor imagery was acquired on 24 July 2013 in the same region of  
 296 Southampton with urban (S1') and rural (S2') study sites close to the Northwest of S1  
 297 and S2. The Worldview-2 image was geometrically and atmospherically corrected, and  
 298 pan-sharpened at 50 cm spatial resolution to be consistent with the aerial imagery.  
 299 Figure 4 demonstrates the WorldView-2 satellite sensor image together with two  
 300 subsets S1' and S2'.



301

302 Figure 4 Additional WorldView-2 satellite sensor image covering the same region of Southampton  
 303 with the S1' and S2' study sites to the northwest of S1 and S2, respectively.

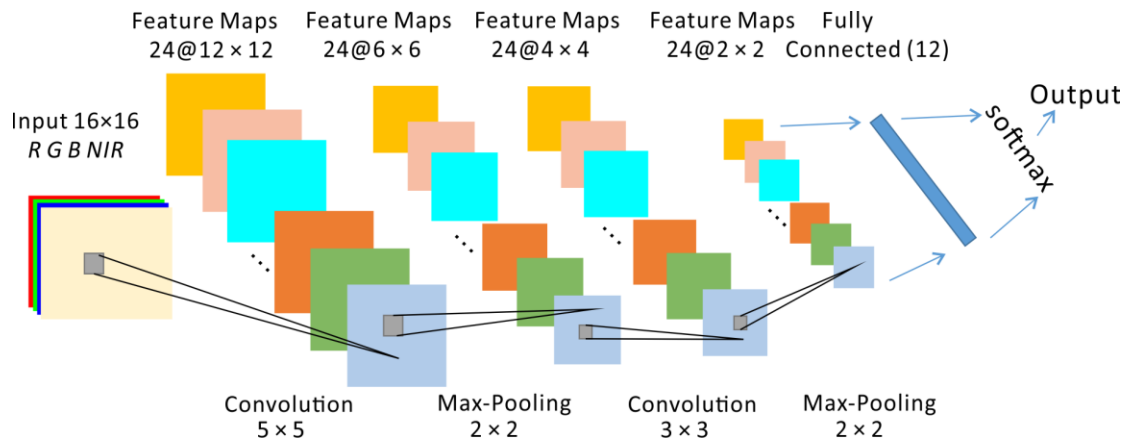
### 304 **3.2 Model input variables and parameters**

305 Model inputs: the standard pixel-based MLP (hereafter, MLP) and CNN take only the  
 306 four spectral bands as their input variables, whereas the pixel-based texture MLP based  
 307 on the standard Grey Level Co-occurrence Matrix (hereafter, GLCM-MLP)  
 308 simultaneously makes use of both the four spectral bands and the texture features  
 309 derived from GLCM textural features including the Mean, Variance, Homogeneity,  
 310 Contrast, Dis-similarity, Entropy, Second moment and Correlation (Haralick et al.,  
 311 1973; Rodriguez-Galiano et al., 2012; Xia et al., 2010; Zhang et al., 2003). Three  
 312 window sizes for each spectral band, including  $3 \times 3$  ( $1.5 \times 1.5$  m),  $5 \times 5$  ( $2.5 \times 2.5$  m), and  
 313  $7 \times 7$  ( $3.5 \times 3.5$  m), were optimally chosen to perform multi-scale texture feature  
 314 representation, thus generating 96 GLCM texture features in total. It should be noted  
 315 that both the MLP and the CNN as well as the GLCM-MLP were trained to predict all  
 316 pixels within the images. Although the CNN was designed to predict a single label from  
 317 a small image patch, the sliding window was densely overlapping to cover the entire  
 318 image at the inference phase.

319 Both the MLP (also including GLCM-MLP) and CNN models require a series of  
 320 predefined parameters to optimize the learning accuracy and generalization capability.  
 321 Following the recommendations of Mas and Flores, (2008), the MLPs with one, two  
 322 and three hidden layers were tested, using a varying number of {4, 8, 12, 16, 20, and  
 323 24} nodes in each layer. The learning rate was chosen optimally as 0.2 and the  
 324 momentum factor was set as 0.7. In addition, the number of iterations was set as 1000  
 325 to fully converge to a stable state. Through cross-validation with different numbers of

326 nodes and hidden layers, the best predicting MLP was found using two hidden layers  
327 with 8 nodes in each layer. Similar parameters were also set for the GLCM-MLP,  
328 except that two hidden layers with 20 nodes in each layer were found to be the optimal  
329 solution in this case.

330 For the CNN, a range of parameters including the number of layers, the input image  
331 patch size, the number and size of convolutional filter, as well as other predefined  
332 parameters, such as the learning rate and number of epochs (iterations), need to be tuned  
333 (Romero et al., 2016). Following the discussion by L ängkvist et al., (2016), the input  
334 image size was chosen from  $\{8 \times 8, 10 \times 10, 12 \times 12, 14 \times 14, 16 \times 16, 18 \times 18, 20 \times 20, 22 \times 22$   
335  $\text{and } 24 \times 24\}$  to evaluate the influence of context area on classification performance. In  
336 general, a small-sized contextual area results in overfitting of the model, whereas a  
337 large one often leads to under-segmentation. In consideration of the image object size  
338 and contextual relationship coupled with a small amount of trial and error, the optimal  
339 input image patch size was set to  $16 \times 16$  in this research. Besides, as discussed by Chen  
340 et al., (2014) and L ängkvist et al., (2016), the depth plays a key role in classification  
341 accuracy because the quality of learnt feature is highly influenced by the level of  
342 abstraction and representation. As suggested by Chen et al. (2016), the number of CNN  
343 layers was chosen as four to balance the network complexity and robustness. Other  
344 parameters were set based on standard practice in the field of computer vision. For  
345 example, the filter size was set to  $5 \times 5$  for the first convolution layer and  $3 \times 3$  for the  
346 rest with stride of 1, and the number of the filters was set to 24 to extract multiple  
347 convolutional features at each level. The fully connected layer was tuned as 12 nodes  
348 followed by a softmax classification. The learning rate was set to 0.01 and the number  
349 of epochs (iterations) was chosen as 600 to fully learn the features through  
350 backpropagation. The detailed architecture of the CNN and its parameter configurations  
351 is illustrated in Figure 5.



352

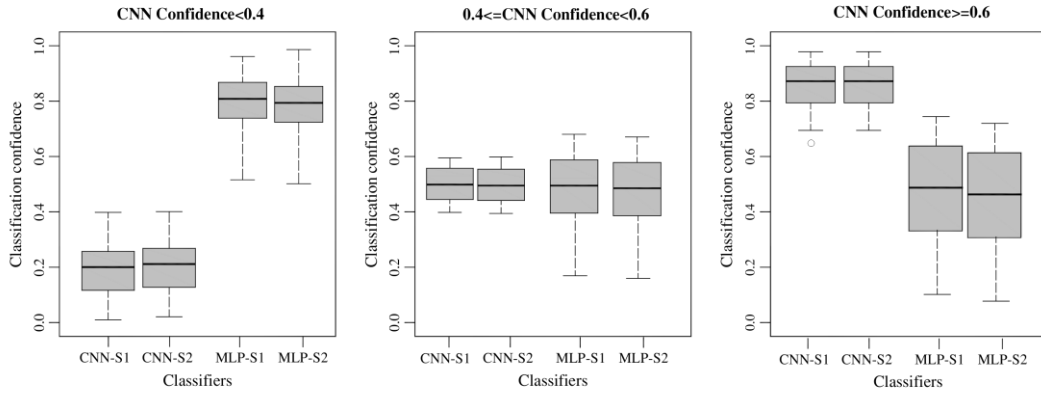
353

Figure 5. The architecture of the CNN and its configurations.

354 **3.3 Decision Fusion Parameter Setting and analysis**

355 A rule-based decision fusion approach was implemented based on the classification  
 356 confidence maps of the CNN (e.g. Figure 2(b)) and MLP (e.g. Figure 2(c)). The  
 357 parameters of decision fusion, including two thresholds  $\alpha_1$  and  $\alpha_2$ , were determined by  
 358 grid search with cross-validation using 10% of the randomly chosen samples. In this  
 359 study, the optimal thresholds  $\alpha_1=0.4$  and  $\alpha_2=0.6$  were found that reported the greatest  
 360 classification accuracy.

361 For the sake of visual interpretation, the confidence distribution of the CNN and MLP  
 362 influenced by the chosen thresholds is shown in Figure 6. Clearly, the CNN and MLP  
 363 demonstrated individually consistent, but mutually converse distribution patterns in the  
 364 two study sites: along with the increase in the CNN's confidence, the MLP inversely  
 365 exhibited a decreasing trend. Specifically, for low CNN confidence ( $<0.4$ ), the MLP  
 366 confidence was around 0.75, significantly higher than that of the CNN, thus outputting  
 367 the results of MLP in the final decision; once the CNN confidence ranged from 0.4 to  
 368 0.6, no significant difference was shown between the two classifiers, thereby, optimally  
 369 choosing the classification results based on the competitive "winner-takes-all"  
 370 approach; while for large CNN confidence ( $>0.6$ ), the MLP was, in contrast, much less  
 371 reliable ( $<0.45$ ), thus, taking the classification results of the CNN only in this situation.



372

373 Figure 6 Classification confidence distributions of the CNN and MLP at two study sites (S1 and S2)  
 374 under different fusion thresholds.

### 375 3.4 Classification results and analysis

#### 376 3.4.1 Classification results and visual assessment

377 By integrating the classification results of the MLP and CNN using the above-  
 378 mentioned fusion parameters, the final classification of the proposed MLP-CNN was  
 379 obtained at both study sites, S1 (city centre with complex urban scene) and S2 (rural  
 380 areas with natural phenomena). To provide a better visualization, Figure 7 (three  
 381 subsets of S1) and Figure 8 (three subsets of S2) highlights the correct or incorrect  
 382 classification results of different classifiers marked in yellow or red circles, respectively.

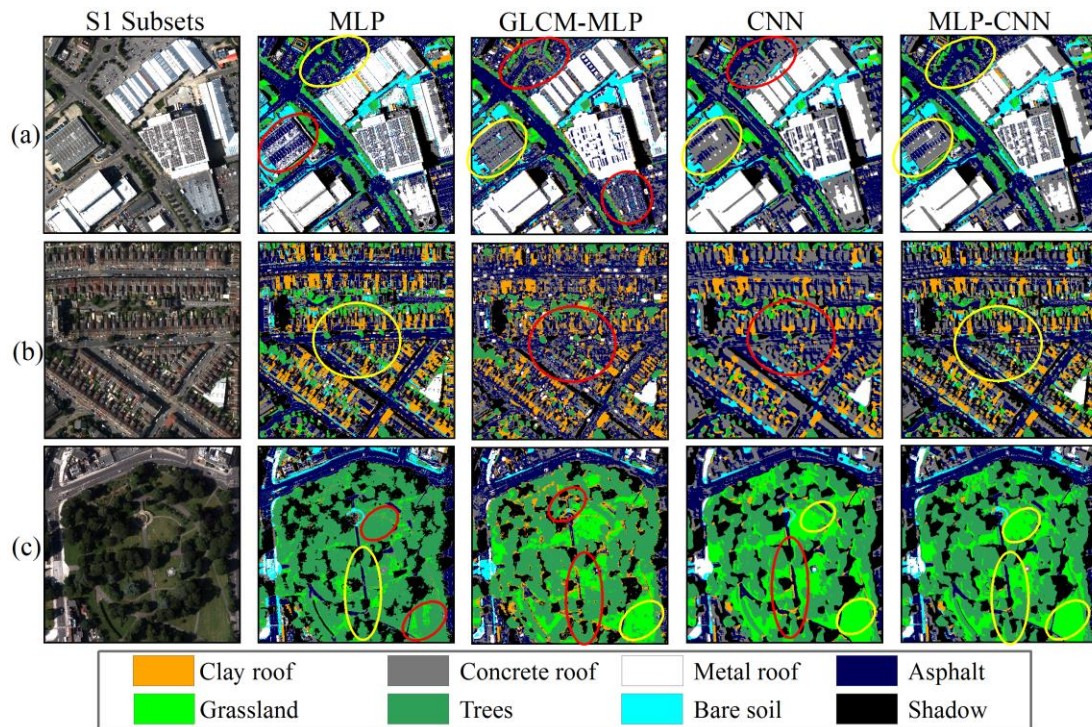
383 From Figure 7, it can be seen that the MLP classification results consist of undesirable  
 384 noise (marked in red circle), such as a severe salt-and-pepper effect in Figure 7(a) and  
 385 7(b), and linear noisy textures in Figure 8(c). Besides, Trees and Grassland are seriously  
 386 confused with each other as illustrated by Figure 7(c) and Figure 8(a) and 8(b).  
 387 However, as shown by Figure 7(b), the MLP has certain advantages over CNN in  
 388 identifying the Clay roof class with spectrally distinctive features (marked in yellow  
 389 circle). With the addition of the GLCM textures, the GLCM-MLP achieved certain  
 390 improvements in both spectral and spatial pattern differentiation. For example, Trees  
 391 and Grassland are better distinguished to some extent compared with the pixel-based  
 392 MLP results, as illustrated in Figure 7(c) and Figure 8(b). Besides, the clear linear noisy  
 393 textures in Figure 8(c) are much reduced, and primarily turned into small speckles due  
 394 to the introduction of texture features. Yet, the GLCM-MLP falsely identifies some  
 395 edges or boundaries as Clay Roof, as shown in Figure 7(c) and Figure 8(a) and 8(b)

396 (marked in red circle). Additionally, some geometrical distortions of building roof tops,  
397 e.g. the Metal Roof and Concrete Roof in Figure 7(b), are shown in the GLCM-MLP  
398 classification results caused by the GLCM texture filters.

399 In contrast to the pixel-based MLP and the GLCM-MLP, the classification results of  
400 the CNN in both study sites exhibit smoothed visual effects with the least speckle noise  
401 as shown by Figure 7 and 8. Additionally, the classes of green vegetation including  
402 Grassland and Trees are accurately distinguished as demonstrated by the yellow circles  
403 in Figure 7(c) and Figure 8(a) and 8(b) in spite of their spectral similarity. Moreover,  
404 the CNN is able to discriminate the Concrete roof from Asphalt with a moderate  
405 accuracy, as highlighted by the yellow circle in Figure 7(a). Nevertheless, the CNN  
406 delivers some uncertainties in partitioning object boundaries. For example, the regular  
407 shapes of some buildings (e.g. the geometries of some Clay roof and Concrete roof  
408 areas) are distorted with false boundary partitions, as marked by the red circle in Figure  
409 7(b). In addition, small or linear features are either merged into a large object or  
410 discarded by over-smoothness. For instance, some Clay roof buildings (small objects)  
411 are falsely connected together, while Asphalt is sometimes misclassified as Clay roof  
412 (Figure 7(c)) and the small paths covered by Bare soil are discarded (Figure 8(b)).

413 With respect to the results of the MLP-CNN, all of the aforementioned  
414 misclassifications produced by MLP or CNN are resolved with a higher resulting  
415 accuracy. Thus, the incorrect classifications (marked by red circles) which appeared in  
416 Figure 7 and 8 are revised accordingly, with no red circles appearing in the  
417 classification results of MLP-CNN. The MLP-CNN modifies the classification errors  
418 of the CNN for Asphalt, as illustrated by the red circles in Figure 7(c) and Figure 8(b),  
419 thanks to the correct classification results of the MLP. Moreover, the linear-shaped Bare  
420 Soil area missed by the CNN in Figure 8(a) is brought back correctly without losing  
421 useful information. In addition, the original shapes of the Clay roof and Concrete roof  
422 areas shown in Figure 7(b) are accurately restored. Most importantly, some mutual  
423 misclassifications between the MLP and CNN are successfully rectified. For example,  
424 the MLP-CNN correctly differentiates some Asphalt (with spectrally distinctive but  
425 spatially confusing characteristics) and Concrete roof (distinctive in texture and  
426 geometry but vague in spectrum) areas that are mutually misclassified by the MLP and  
427 CNN respectively (see the regions marked by red circles in Figure 7(a)).



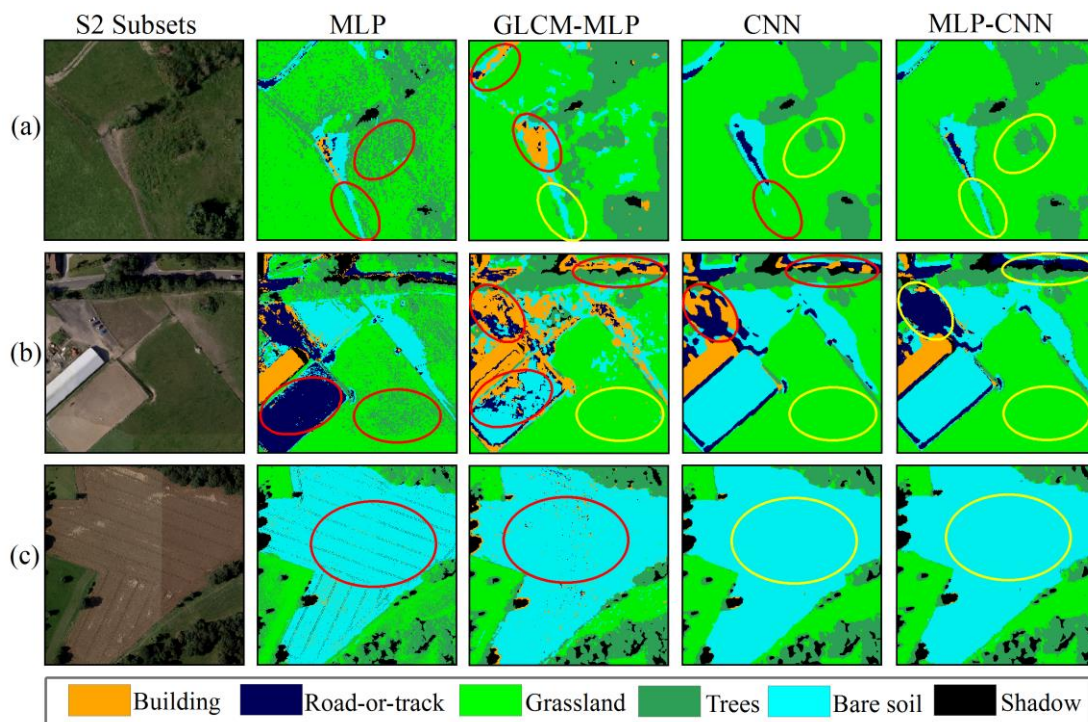


428

429 Figure 7 Three typical image subsets (a, b and c) in study site S1 with their classification results.

430 Columns from left to right represent the original images (R G B bands), the MLP classification, the  
 431 GLCM-MLP classification, the CNN classification and the MLP-CNN classification correspondingly.

432 The red and yellow circles denote incorrect and correct classification, respectively.



433

434 Figure 8 Three typical image subsets (a, b and c) in study site S2 with their classification results.

435 Columns from left to right represent the original images (R G B bands), the MLP classification, the

436 GLCM-MLP classification, the CNN classification and the MLP-CNN classification correspondingly.  
437 The red and yellow circles denote incorrect and correct classification, respectively.

### 438 ***3.4.2 Classification accuracy assessment***

439 The classification performance of the proposed MLP-CNN approach was further  
440 investigated through benchmark comparison with the MLP, GLCM-MLP and the CNN.  
441 Table 2 lists the classification accuracy assessment, including the overall accuracy  
442 (OA), Kappa coefficient ( $\kappa$ ), and the class-wise mapping accuracy. From the table, it  
443 can be seen that the decision fusion approach (MLP-CNN) consistently reports the best  
444 classification OA with up to 90.93% for S1 and 89.64% for S2, higher than that of the  
445 CNN (85.39% and 86.56%, respectively) and GLCM-MLP (83.12% and 82.63%,  
446 respectively) as well as MLP (81.62% and 80.73%, respectively) (Table 2). Moreover,  
447 a Kappa  $z$ -test for pair-wise comparison also shows that a significant increase in  
448 classification accuracy has been achieved by the proposed MLP-CNN classifier over  
449 the MLP, GLCM-MLP and CNN in S1, with  $z$ -value=3.68, 3.12 and 2.25, respectively.  
450 For S2, the MLP-CNN also revealed a significant increase over the MLP with  $z$ -  
451 value=3.71 as well as GLCM-MLP with  $z$ -value=3.18, but no significant difference in  
452 comparison with the CNN ( $z = 1.59$ , smaller than 1.96 at 95% confidence level)  
453 (Congalton, 1991), despite the obvious improvement shown in Table 2.

454 The increase in classification accuracy was also checked by class-wise accuracy  
455 assessment (Table 3). As illustrated by the table, MLP-CNN outperforms CNN for all  
456 classes at both study sites in terms of classification accuracy. The largest increase is up  
457 to 9.77% for the class of Concrete roof in S1 and 7.16% for the class of Road-or-track  
458 in S2. Similar patterns were found such that the MLP-CNN was constantly superior to  
459 GLCM-MLP at the class-wise level, where the greatest increase in accuracy was shown  
460 up to 11.56% for the class of Concrete Roof in S1 and 11.74% for the class of Grassland  
461 in S2. When compared with the MLP, most classes in the two sites except for Asphalt  
462 and Shadow in S1 are classified with higher accuracy by the MLP-CNN. Here,  
463 Grassland exhibits the highest increase in classification accuracy, up to 33.51% and  
464 18.83% for S1 and S2, respectively. For the classes of Asphalt and Shadow, the  
465 accuracy of the MLP is slightly larger than that of the MLP-CNN, but without a  
466 statistically significant difference. Thus, they can be regarded as similar to each other.

467 With respect to the three benchmark classifiers themselves (i.e. MLP, GLCM-MLP and  
 468 CNN), it can be seen from Table 2 that their classification accuracies are ordered as:  
 469 MLP <GLCM-MLP < CNN. While the accuracy of CNN is remarkably higher (3%-  
 470 5%) than that of the MLP and GLCM-MLP, the GLCM-MLP is just slightly higher  
 471 (<2%) than the MLP. The Kappa  $z$ -tests (Table 3) further demonstrate that the CNN is  
 472 statistically significantly more accurate than MLP and GLCM-MLP in both urban and  
 473 rural areas, whereas a significant increase in accuracy of the GLCM-MLP over the MLP  
 474 appears only in the rural area rather than the urban area.

475 Table 2 Classification accuracy comparison amongst MLP, GLCM-MLP, CNN and the proposed MLP-  
 476 CNN approach for study sites S1 and S2 using the per-class mapping accuracy, overall accuracy (OA)  
 477 and Kappa coefficient ( $\kappa$ ). The bold font highlights the greatest classification accuracy per row.

Study Sites	Class	MLP	GLCM-MLP (Benchmark)	CNN	MLP-CNN
S1	Clay roof	92.26%	91.43%	90.11%	<b>95.03%</b>
	Concrete roof	67.06%	62.44%	64.23%	<b>74.00%</b>
	Metal roof	91.13%	90.36%	94.19%	<b>94.63%</b>
	Asphalt	<b>92.72%</b>	88.67%	85.98%	91.26%
	Grassland	60.51%	82.58%	90.73%	<b>94.02%</b>
	Trees	63.88%	78.46%	82.28%	<b>88.83%</b>
	Bare soil	79.63%	83.05%	86.16%	<b>92.49%</b>
	Shadow	<b>92.33%</b>	91.06%	91.14%	91.52%
	Overall Accuracy (OA)	81.62%	83.12%	85.39%	<b>90.93%</b>
Kappa Coefficient ( $\kappa$ )	0.78	0.81	0.84	<b>0.89</b>	
S2	Building	82.83%	80.79%	83.08%	<b>88.48%</b>
	Road or track	83.02%	80.14%	82.42%	<b>89.58%</b>
	Grassland	71.11%	78.20%	88.34%	<b>89.94%</b>
	Trees	79.31%	84.55%	90.70%	<b>92.86%</b>
	Bare soil	74.07%	76.32%	81.36%	<b>86.86%</b>
	Shadow	89.41%	88.25%	88.37%	<b>90.17%</b>
	Overall Accuracy (OA)	80.73%	82.63%	86.56%	<b>89.64%</b>
Kappa Coefficient ( $\kappa$ )	0.78	0.79	0.84	<b>0.87</b>	

478

479 Table 3 Kappa  $z$ -test ( $p$ -value) comparing the performance of the three classifiers for two study sites S1  
 480 and S2. Significantly different accuracies with confidence of 95% ( $z$ -value > 1.96 with  $p$ -value < 0.05)  
 481 are indicated by \*.

Study sites	Classifiers	Kappa Z-test (p-value)			
		MLP	GLCM-MLP (Benchmark)	CNN	MLP- CNN
S1	MLP	—	—	—	—
	GLCM-MLP	1.56 (0.1188)	—	—	—
	CNN	2.64* (0.0083)	2.44* (0.0147)	—	—
	MLP-CNN	3.68* (0.0002)	3.12* (0.0018)	2.25* (0.0244)	—
S2	MLP	—	—	—	—
	GLCM-MLP	2.05* (0.0404)	—	—	—
	CNN	2.51* (0.0121)	2.36* (0.0183)	—	—
	MLP-CNN	3.71* (0.0002)	3.18* (0.0015)	1.59 (0.1118)	—

482

483 The proposed MLP-CNN method and the other three benchmarks (MLP, GLCM-MLP  
484 and the CNN) were also validated using an additional WorldView-2 satellite sensor  
485 dataset at the S1' and S2' study sites. The OA and  $\kappa$  of both study sites are in accordance  
486 with the results of aerial photo classification, where the decision fusion approach (MLP-  
487 CNN) acquires the largest OA of 90.56% at S1' and 89.77% at S2', consistently higher  
488 than the CNN (86.15% and 86.39%), the GLCM-MLP (83.26% and 82.52%) and the  
489 MLP (81.42% and 80.32%) (Table 4). Such coherency of classification results further  
490 demonstrates the wide applicability of the proposed method with different datasets.

491 Table 4 Classification accuracy comparison amongst MLP, GLCM-MLP (Benchmark), CNN and the  
492 proposed MLP-CNN approach for study sites S1' and S2' from the WorldView-2 satellite sensor image  
493 using overall accuracy (OA) and Kappa coefficient ( $\kappa$ ). The bold font highlights the greatest  
494 classification accuracy per row.

WorldView-2	Classification	MLP	GLCM-MLP (Benchmark)	CNN	MLP- CNN
S1'	OA	81.42%	83.26%	86.15%	<b>90.56%</b>
	$\kappa$	0.77	0.80	0.82	<b>0.89</b>
S2'	OA	80.32%	82.52%	86.39%	<b>89.77%</b>
	$\kappa$	0.77	0.79	0.83	<b>0.87</b>

495

#### 496 4. Discussion

497 In this research, a rule-based decision fusion approach (MLP-CNN) was proposed to  
498 integrate classifiers of the pixel-based MLP with shallow structures and the contextual-  
499 based CNN with deep architectures for the classification of VFSR remotely sensed  
500 imagery. The MLP-CNN takes advantage of the merits of the two classifiers and  
501 overcomes their individual shortcomings as discussed below.

#### 502 ***4.1 Characteristics of MLP and GLCM-MLP classification***

503 In principle, the MLP builds the decision boundaries among classes in feature space  
504 based on per-pixel spectral information (Mokhtarzade and Zoej, 2007). Such  
505 classification boundaries are very sensitive to the class with salient spectral properties  
506 that are spectrally distinctive from other classes (Berberoglu et al., 2000). For example,  
507 classes like Clay roof, Asphalt and Shadow in Site 1 are spectrally exclusive to other  
508 classes, leading to high classification accuracies, up to 92.26%, 92.72% and 92.33%,  
509 respectively (Table 2). However, the MLP relies on the pixel-based spectral information  
510 in the classification process without exploiting the abundant spatial information  
511 appearing in the VFSR imagery (e.g. texture, geometry or contextual relationship)  
512 (Wang et al., 2016). These limitations often result in unsatisfactory classification  
513 performance; for example, confusion and misclassification between the Trees and  
514 Grassland classes that are spectrally similar. Even for those correctly identified objects,  
515 severe salt and pepper effects still exist (Dark and Bram, 2007), for example, the linear  
516 texture noise appearing for Bare soil in Figure 8(c). For these reasons, the classification  
517 accuracy of MLP is generally statistically significantly lower than that of the CNN and  
518 the proposed MLP-CNN. However, objects in VFSR imagery are mostly depicted by  
519 pure pixels, especially for human-made features with crisp boundaries, such as  
520 buildings, residential houses and cultivated land. The membership association of a pixel  
521 deduced by MLP is, therefore, not affected by its relative position (e.g. lying on or close  
522 to boundaries), as long as the corresponding spectral space is separable.

523 The inclusion of GLCM texture features in the GLCM-MLP classifier enables the  
524 model to process spectral and spatial information simultaneously. Those GLCM texture  
525 descriptors are handcrafted features that are designed to capture statistical co-  
526 occurrence information (Xia et al., 2010). However, the GLCM textures are essentially  
527 first or second order feature transformations instead of feature learning. Such hand-  
528 coded features might be effective for a particular region and/or season, but are often

529 challenging to generalize to other domains and datasets. Besides, the addition of 96  
530 GLCM textures results in a dramatically increased number of input variables, which  
531 leads to a relatively high dimensional feature space. The so-called “curse of  
532 dimensionality” (Hughes, 1968) and collinearity make the GLCM-MLP hard to  
533 parameterize and potentially leads to texture overfitting. That is why the GLCM-MLP  
534 cannot substantially increase the classification accuracy compared to the MLP. That is,  
535 the spectral and spatial information cannot be effectively exploited by the GLCM-MLP.  
536 For example, some spectrally different classes but with similar textures such as Clay  
537 Roof, Concrete Roof and Asphalt are confused to some degree.

#### 538 *4.2 Characteristics of CNN classification*

539 Spatial features in remotely sensed data like VFSR imagery are intrinsically local  
540 (especially in lower layers) and spatially invariant (Masi et al., 2016). The MLP,  
541 however, assumes that the location of the data in the input is irrelevant to the model  
542 construction and it is, thus, incapable of learning spatial features of remote sensing data.  
543 In contrast, by using multiple convolution and pooling operations, CNN models the  
544 way that the human visual cortex works and enforces weight sharing with translation  
545 invariance that enables the extraction of high-level spatial features from image patches.  
546 It should be mentioned that the pooling operations play an important role in dimension  
547 reduction, thus, avoiding “the curse of dimensionality” present in the GLCM-MLP  
548 classifier. Thanks to these superior characteristics, the CNN classifier outperforms the  
549 MLP and GLCM-MLP classifiers in both the urban scene and rural areas. Especially,  
550 classes like Concrete roof and Road-or-track that are difficult to distinguish from their  
551 backgrounds with only spectral or low-level features (e.g. distance between the  
552 prediction and the target class at spectral space), are identified with relatively high  
553 accuracies. In addition, classes with heavy spectral confusion in both study sites (e.g.  
554 Trees and Grassland), are accurately differentiated due to their obvious spatial pattern  
555 differences; for example, the texture of tree canopies is generally much rougher than  
556 for grassland. As a contextual classifier with deep architectures, the CNN could reveal  
557 the spatial patterns hidden in the image data that cannot be perceived by its shallow  
558 counterparts (e.g. MLP classifier or even the GLCM-MLP classifier). The higher layers  
559 in CNN models provide more semantically meaningful information concentrating on  
560 global semantics rather than local or pixel-level information, making the CNN  
561 classification work well for classes with spectral confusion (Hu et al., 2015a, 2015b;

562 Yang et al., 2015). Therefore, the CNN shows an impressive stability and effectiveness  
563 in spatial feature representation, which is crucial for VFSR image classification (Zhao  
564 and Du, 2016).

565 However, according to the “no free lunch” theorem (Wolpert and Macready, 1997), any  
566 elevated performance in one aspect of a problem will be paid for through others, and  
567 the CNN is no exception. Using contextual image patches as inputs and learning deep  
568 spatial features, the CNN demonstrates power in spatial pattern recognition but also  
569 weakness in spatial partition. Boundary uncertainties (over-smoothness) often appear  
570 in the classified object and small useful features are erased, somewhat similar to  
571 morphological or Gabor filter methods (Pingel et al., 2013; Reis and Tasdemir, 2011).  
572 For example, the human-made objects in urban scenes like buildings and asphalt are  
573 often geometrically enlarged with distortion to some degree (See Figure 7(b)). As for  
574 natural objects in rural areas (S2), edges or porosities of a landscape patch are simplified  
575 or ignored, and even worse, linear features like river channels or dams that are of  
576 ecological importance, are erroneously erased. One may argue that the reduction of  
577 image patch size might be able to detect small features by multiple CNNs by varying  
578 the contextual filter size as adopted in Längkvist et al. (2016). However, objects,  
579 whether large or small in size, all have boundaries, thus, retaining the problem of  
580 smoothing edges. In addition, the adoption of convolution and pooling operations  
581 intrinsically reduces the image contextual size but strengthens the spatial feature  
582 representation. Thus, a far too small initial image patch size can limit the network depth  
583 of a CNN model. In fact, the currently used  $16 \times 16$  window size is close to the minimum  
584 requirements for a deep CNN with four hidden layers in total. Moreover, certain  
585 spectrally distinctive features without obvious spatial patterns are poorly differentiated.  
586 For example, some Asphalt pixels are wrongly identified as Concrete roofs as illustrated  
587 in Figure 7(a). This further demonstrates the necessity of introducing spectral features  
588 for VFSR image classification.

#### 589 ***4.3 fusion decision of MLP-CNN classification***

590 Huge uncertainty and inconsistency exists inherently in any remotely sensed data  
591 (including VFSR imagery), and this runs through the training and the testing samples.  
592 In fact, different classification algorithms vary in terms of remote sensing data  
593 processing strategies. Thus there is no ‘one-algorithm-fits-all’ solution (Löw et al.,

594 2015) to various applications of VFSR image classification, even for the powerful CNN  
595 classifier with deep spatial feature representations. It is therefore especially important  
596 to make use of the complementarities of different classifiers. It should be mentioned  
597 that, the more heterogeneous the classification algorithms' behaviours, the more that  
598 different places might be accurately classified by each individual classifier, and the  
599 more accurate the ensemble classifier might be (Löw et al., 2015). An ideal ensemble  
600 classifier, thereby, should be established using individual classifiers that are very  
601 differently behaved.

602 The experimental results show that the pixel-based MLP classifier with shallow  
603 structures and the contextual-based CNN classifier with deep architectures can provide  
604 complementary information, leading to a more accurate classification result than either  
605 classifier alone. In addition to the elimination of heavy noise, the CNN can accurately  
606 identify classes with rich spatial information implicit in VFSR data. Such  
607 characteristics of the CNN emphasize the limitations of the MLP classifier for VFSR  
608 image classification. At the same time, the CNN might lose some useful details, and it  
609 has difficulties in utilizing spectral information and delineating object boundaries and  
610 is, thus, incapable of maintaining geometric fidelity. The MLP classifier, however,  
611 compensates directly with regard to the limitations of the CNN. The aforementioned  
612 complementary properties between the CNN and MLP are well reflected from the  
613 inverse confidence trends of the two classifiers (Figure 2). Specifically, in the case of  
614 the CNN with the highest confidence, the MLP has the least confidence and *vice versa*,  
615 which further indicates that the proposed MLP-CNN ensemble classifier can take  
616 advantage of the MLP and CNN.

617 The proposed fusion decision rules were derived primarily on the basis of the CNN's  
618 confidence distribution, in consideration of the superiority of CNN classification  
619 performance and the regularity of its confidence distribution. Such a decision fusion  
620 strategy captures the patterns of the complementarities between the two individual  
621 classifiers in general, thus, achieving a desirable classification result. At the same time,  
622 the MLP-CNN classifier demonstrates great utility and wide applicability for both  
623 aerial photography and WorldView-2 satellite sensor imagery with consistent and  
624 competitive classification performance. However, in comparison with MLP, the  
625 classification accuracies of Asphalt and Shadow were slightly higher than for the  
626 proposed MLP-CNN. This means that there is still room for improvement of the



627 decision fusion rules at the class-wise level for VFSR image classification. It might be  
628 better to incorporate the spectral separability differentiated by MLP to achieve the best  
629 classification performance at class level. Besides, no significant improvement was  
630 acquired for rural areas (S2) by the MLP-CNN compared with the CNN. This is mainly  
631 due to the ineffectiveness of the MLP in classifying natural features that dominate in  
632 the rural environment. This shortcoming might be overcome by the replacement of the  
633 MLP by other non-parametric machine learning classifiers (e.g. SVM, RF, etc.).  
634 Moreover, incorporating other data sources (e.g. digital surface model) might be needed  
635 to increase the accuracy of the MLP-CNN for both the CNN and MLP with very low  
636 confidence simultaneously. These aforementioned issues will be investigated in future  
637 research.

## 638 **5. Conclusion**

639 Due to its high intra-class variability and low inter-class disparity, VFSR image  
640 classification poses great challenges to any single machine learning algorithm, even for  
641 the powerful deep learning convolutional neural network (CNN). In this paper, two  
642 neural network classifiers with strong heterogeneous behaviours (i.e. pixel-based MLP  
643 with shallow structures and contextual-based CNN with deep architectures), were  
644 integrated in a concise and effective way using a rule-based decision fusion strategy.  
645 The decision fusion rules, designed primarily on the basis of the classification  
646 confidence of the CNN, reflect the general complementary patterns of both the MLP  
647 and CNN. In consequence, the proposed ensemble classifier MLP-CNN harvests the  
648 complementary results acquired from the CNN with deep spatial feature representations  
649 (CNN) and from the MLP based on spectral discrimination. Meanwhile, limitations of  
650 the CNN such as uncertainty in object boundary partition and loss of useful fine  
651 resolution detail were compensated. The effectiveness of the new MLP-CNN algorithm  
652 was tested in both urban and rural areas using aerial and satellite sensor images. The  
653 MLP-CNN algorithm consistently outperformed both of the individual classifiers (MLP  
654 and CNN) as well as the GLCM-MLP that includes the GLCM texture features, with a  
655 statistically significant difference in the majority of cases. This research paves the way  
656 to an effective solution to the complicated problem of automatic VFSR image  
657 classification.

## 658 **Acknowledgement**

659 This research was funded by PhD studentship “Deep Learning in massive area, multi-  
660 scale resolution remotely sensed imagery” (NO. EAA7369), sponsored by Ordnance  
661 Survey and Lancaster University. The authors thank to the staff from the Ordnance  
662 Survey for the supply of aerial imagery and supporting ground data. The authors also  
663 thank to the two anonymous referees for their constructive comments on this  
664 manuscript.

## 665 **Reference**

666 Ardila, J.P., Tolpekin, V.A., Bijker, W., Stein, A., 2011. Markov-random-field-based  
667 super-resolution mapping for identification of urban trees in VHR images.  
668 *ISPRS J. Photogramm. Remote Sens.* 66, 762–775.  
669 doi:10.1016/j.isprsjprs.2011.08.002

670 Arel, I., Rose, D.C., Karnowski, T.P., 2010. Deep machine learning - A new frontier  
671 in artificial intelligence research. *IEEE Comput. Intell. Mag.* 5, 13–18.  
672 doi:10.1109/MCI.2010.938364

673 Atkinson, P.M., Tatnall, A.R.L., 1997. Introduction Neural networks in remote  
674 sensing. *Int. J. Remote Sens.* 18, 699–709. doi:10.1080/014311697218700

675 Attarchi, S., Gloaguen, R., 2014. Classifying complex mountainous forests with L-  
676 Band SAR and landsat data integration: A comparison among different machine  
677 learning methods in the Hyrcanian forest. *Remote Sens.* 6, 3624–3647.  
678 doi:10.3390/rs6053624

679 Benediktsson, J.A., 2009. Ensemble classification algorithm for hyperspectral remote  
680 sensing data. *IEEE Geosci. Remote Sens. Lett.* 6, 762–766.  
681 doi:10.1109/LGRS.2009.2024624

682 Berberoglu, S., Lloyd, C.D., Atkinson, P.M., Curran, P.J., 2000. The integration of  
683 spectral and textural information using neural networks for land cover mapping  
684 in the Mediterranean. *Comput. Geosci.* 26, 385–396. doi:10.1016/S0098-  
685 3004(99)00119-3

686 Bezak, P., Bozek, P., Nikitin, Y., 2014. Advanced robotic grasping system using deep  
687 learning. *Procedia Eng.* 96, 10–20.  
688 doi:http://dx.doi.org/10.1016/j.proeng.2014.12.092

- 689 Breiman, L., 1996. Bagging Predictors. *Mach. Learn.* 24, 123–140.
- 690 Chen, S., Member, S., Wang, H., Xu, F., Member, S., 2016. Target classification  
691 using the deep Convolutional Networks for SAR images. *IEEE Trans. Geosci.*  
692 *Remote Sens.* 54, 4806–4817.
- 693 Chen, Y., Jiang, H., Li, C., Jia, X., Member, S., 2016. Deep feature extraction and  
694 classification of hyperspectral images based on Convolutional Neural Networks.  
695 *IEE Trans. Geosci. Remote Sens.* 54, 6232–6251.  
696 doi:10.1109/TGRS.2016.2584107
- 697 Chen, Y., Lin, Z., Zhao, X., Wang, G., Gu, Y., 2014. Deep learning-based  
698 classification of hyperspectral data. *IEEE J. Sel. Top. Appl. Earth Obs. Remote*  
699 *Sens.* 7, 2094–2107. doi:10.1109/JSTARS.2014.2329330
- 700 Clinton, N., Yu, L., Gong, P., 2015. Geographic stacking: Decision fusion to increase  
701 global land cover map accuracy. *ISPRS J. Photogramm. Remote Sens.* 103, 57–  
702 65. doi:10.1016/j.isprsjprs.2015.02.010
- 703 Congalton, R.G., 1991. A review of assessing the accuracy of classifications of  
704 remotely sensed data. *Remote Sens. Environ.* 37, 35–46.
- 705 Dark, S.J., Bram, D., 2007. The modifiable areal unit problem (MAUP) in physical  
706 geography. *Prog. Phys. Geogr.* 31, 471–479. doi:10.1177/0309133307083294
- 707 Del Frate, F., Pacifici, F., Schiavon, G., Solimini, C., 2007. Use of neural networks  
708 for automatic classification from high-resolution images. *IEEE Trans. Geosci.*  
709 *Remote Sens.* 45, 800–809. doi:10.1109/TGRS.2007.892009
- 710 Demarchi, L., Canters, F., Cariou, C., Licciardi, G., Chan, J.C.W., 2014. Assessing  
711 the performance of two unsupervised dimensionality reduction techniques on  
712 hyperspectral APEX data for high resolution urban land-cover mapping. *ISPRS*  
713 *J. Photogramm. Remote Sens.* 87, 166–179. doi:10.1016/j.isprsjprs.2013.10.012
- 714 Dong, Z., Pei, M., He, Y., Liu, T., Dong, Y., Jia, Y., 2015. Vehicle type classification  
715 using unsupervised Convolutional Neural Network. *IEEE Trans. Intell. Transp.*  
716 *Syst.* 16, 2247–2256. doi:10.1109/ICPR.2014.39
- 717 Du, P., Xia, J., Zhang, W., Tan, K., Liu, Y., Liu, S., 2012. Multiple classifier system

718 for remote sensing image classification: A review. *Sensors* 12, 4764–4792.  
719 doi:10.3390/s120404764

720 Farabet, C., Couprie, C., Najman, L., Lecun, Y., 2013. Learning hierarchical features  
721 for scene labeling. *IEEE Trans. Pattern Anal. Mach. Intell.* 35, 1915–1929.  
722 doi:10.1109/TPAMI.2012.231

723 Fauvel, M., Chanussot, J., Benediktsson, J.A., 2012. A spatial-spectral kernel-based  
724 approach for the classification of remote-sensing images. *Pattern Recognit.* 45,  
725 381–392. doi:10.1016/j.patcog.2011.03.035

726 Fauvel, M., Chanussot, J., Benediktsson, J.A., 2006. Decision fusion for the  
727 classification of urban remote sensing images. *IEEE Trans. Geosci. Remote*  
728 *Sens.* 44, 2828–2838. doi:10.1109/TGRS.2006.876708

729 Foody, G.M., Arora, M.K., 1997. An evaluation of some factors affecting the  
730 accuracy of classification by an artificial neural network. *Int. J. Remote Sens.* 18,  
731 799–810. doi:10.1080/014311697218764

732 Freund, Y., Iyer, R., Schapire, R.E., Singer, Y., 2003. An efficient boosting algorithm  
733 for combining preferences. *J. Mach. Learn. Res.* 4, 933–969.

734 Haralick, R.M., Shanmugam, K., Dinstein, I., 1973. Textural Features for Image  
735 Classification. *IEEE Trans. Syst. Man. Cybern.* 3, 610–621.  
736 doi:10.1109/TSMC.1973.4309314

737 Hu, F., Xia, G.-S., Hu, J., Zhang, L., 2015a. Transferring Deep Convolutional Neural  
738 Networks for the Scene Classification of High-Resolution Remote Sensing  
739 Imagery. *Remote Sens.* 7, 14680–14707. doi:10.3390/rs71114680

740 Hu, F., Xia, G.S., Wang, Z., Huang, X., Zhang, L., Sun, H., 2015b. Unsupervised  
741 Feature Learning Via Spectral Clustering of Multidimensional Patches for  
742 Remotely Sensed Scene Classification. *IEEE J. Sel. Top. Appl. Earth Obs.*  
743 *Remote Sens.* 8, 2015–2030. doi:10.1109/JSTARS.2015.2444405

744 Huang, X., Lu, Q., Zhang, L., 2014. A multi-index learning approach for  
745 classification of high-resolution remotely sensed images over urban areas. *ISPRS*  
746 *J. Photogramm. Remote Sens.* 90, 36–48. doi:10.1016/j.isprsjprs.2014.01.008

- 747 Hughes, G.F., 1968. On the Mean Accuracy of Statistical Pattern Recognizers. IEEE  
748 Trans. Inf. Theory 14, 55–63. doi:10.1109/TIT.1968.1054102
- 749 Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. ImageNet classification with deep  
750 Convolutional Neural Networks, in: NIPS2012: Neural Information Processing  
751 Systems. Lake Tahoe, Nevada, pp. 1–9.
- 752 L ängkvist, M., Kiselev, A., Alirezaie, M., Loutfi, A., 2016. Classification and  
753 segmentation of satellite orthoimagery using Convolutional Neural Networks.  
754 Remote Sens. 8, 1–21. doi:10.3390/rs8040329
- 755 LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. Nature 521, 436–444.  
756 doi:10.1038/nature14539
- 757 Lenz, I., Lee, H., Saxena, A., 2015. Deep learning for detecting robotic grasps. Int. J.  
758 Rob. Res. 34, 705–724. doi:10.1177/0278364914549607
- 759 Löw, F., Conrad, C., Michel, U., 2015. Decision fusion and non-parametric classifiers  
760 for land use mapping using multi-temporal RapidEye data. ISPRS J.  
761 Photogramm. Remote Sens. 108, 191–204.  
762 doi:http://dx.doi.org/10.1016/j.isprsjprs.2015.07.001
- 763 Mas, J.F., Flores, J.J., 2008. The application of artificial neural networks to the  
764 analysis of remotely sensed data. Int. J. Remote Sens. 29, 617–663.  
765 doi:10.1080/01431160701352154
- 766 Masi, G., Cozzolino, D., Verdoliva, L., Scarpa, G., Wang, L., Zhou, G., Thenkabail,  
767 P.S., 2016. Pansharpening by Convolutional Neural Networks. Remote Sens. 8,  
768 1–22. doi:10.3390/rs8070594
- 769 Mathieu, R., Freeman, C., Aryal, J., 2007. Mapping private gardens in urban areas  
770 using object-oriented techniques and very high-resolution satellite imagery.  
771 Landsc. Urban Plan. 81, 179–192. doi:10.1016/j.landurbplan.2006.11.009
- 772 Min, J.H., Lee, Y.-C., 2005. Bankruptcy prediction using support vector machine with  
773 optimal choice of kernel function parameters. Expert Syst. Appl. 28, 603–614.
- 774 Mokhtarzade, M., Zoj, M.J.V., 2007. Road detection from high-resolution satellite  
775 images using artificial neural networks. Int. J. Appl. Earth Obs. Geoinf. 9, 32–40.

776           doi:10.1016/j.jag.2006.05.001

777   Othman, E., Bazi, Y., Alajlan, N., Alhichri, H., Melgani, F., 2016. Using  
778           convolutional features and a sparse autoencoder for land-use scene classification.  
779           *Int. J. Remote Sens.* 37, 2149–2167. doi:10.1080/01431161.2016.1171928

780   Ozdarici-Ok, A., Ok, A., Schindler, K., 2015. Mapping of agricultural crops from  
781           single high-resolution multispectral images—Data-driven smoothing vs. Parcel-  
782           based smoothing. *Remote Sens.* 7, 5611–5638. doi:10.3390/rs70505611

783   Pacifici, F., Chini, M., Emery, W.J., 2009. A neural network approach using multi-  
784           scale textural metrics from very high-resolution panchromatic imagery for urban  
785           land-use classification. *Remote Sens. Environ.* 113, 1276–1292.  
786           doi:10.1016/j.rse.2009.02.014

787   Pingel, T.J., Clarke, K.C., McBride, W.A., 2013. An improved simple morphological  
788           filter for the terrain classification of airborne LIDAR data. *ISPRS J.*  
789           *Photogramm. Remote Sens.* 77, 21–30. doi:10.1016/j.isprsjprs.2012.12.002

790   Powers, R.P., Hermosilla, T., Coops, N.C., Chen, G., 2015. Remote sensing and  
791           object-based techniques for mapping fine-scale industrial disturbances. *Int. J.*  
792           *Appl. Earth Obs. Geoinf.* 34, 51–57. doi:10.1016/j.jag.2014.06.015

793   Regnauld, N., Mackaness, W. a., 2006. Creating a hydrographic network from its  
794           cartographic representation: a case study using Ordnance Survey MasterMap  
795           data. *Int. J. Geogr. Inf. Sci.* 20, 611–631. doi:10.1080/13658810600607402

796   Reis, S., Tasdemir, K., 2011. Identification of hazelnut fields using spectral and gabor  
797           textural features. *ISPRS J. Photogramm. Remote Sens.* 66, 652–661.  
798           doi:10.1016/j.isprsjprs.2011.04.006

799   Rodriguez-Galiano, V.F., Chica-Olmo, M., Abarca-Hernandez, F., Atkinson, P.M.,  
800           Jeganathan, C., 2012. Random Forest classification of Mediterranean land cover  
801           using multi-seasonal imagery and multi-seasonal texture. *Remote Sens. Environ.*  
802           121, 93–107. doi:10.1016/j.rse.2011.12.003

803   Romero, A., Gatta, C., Camps-valls, G., Member, S., 2016. Unsupervised deep feature  
804           extraction for remote sensing image classification. *IEEE Trans. Geosci. Remote*

805 Sens. 54, 1349–1362. doi:10.1109/TGRS.2015.2478379.

806 Schmidhuber, J., 2015. Deep Learning in neural networks: An overview. *Neural*  
807 *Networks*. doi:10.1016/j.neunet.2014.09.003

808 Shi, H., Chen, L., Bi, F., Chen, H., Yu, Y., 2015. Accurate urban area detection in  
809 remote sensing images. *IEEE Geosci. Remote Sens. Lett.* 12, 1948–1952.

810 Strigl, D., Kofler, K., Podlipnig, S., 2010. Performance and scalability of GPU-based  
811 Convolutional Neural Networks, in: 2010 18th Euromicro Conference on  
812 Parallel, Distributed and Network-Based Processing. pp. 317–324.  
813 doi:10.1109/PDP.2010.43

814 Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D.,  
815 Vanhoucke, V., Rabinovich, A., 2015. Going deeper with convolutions, in:  
816 Proceedings of the IEEE Computer Society Conference on Computer Vision and  
817 Pattern Recognition. pp. 1–9. doi:10.1109/CVPR.2015.7298594

818 Tang, J., Deng, C., Huang, G.-B., Zhao, B., 2015. Compressed-domain ship detection  
819 on spaceborne optical image using Deep Neural Network and Extreme Learning  
820 Machine. *IEEE Trans. Geosci. Remote Sens.* 53, 1174–1185.  
821 doi:10.1109/TGRS.2014.2335751

822 Wang, L., Shi, C., Diao, C., Ji, W., Yin, D., 2016. A survey of methods incorporating  
823 spatial information in image classification and spectral unmixing. *Int. J. Remote*  
824 *Sens.* 37, 3870–3910. doi:10.1080/01431161.2016.1204032

825 Wolpert, D.H., Macready, W.G., 1997. No free lunch theorems for optimization.  
826 *IEEE Trans. Evol. Comput.* 1, 67–82. doi:10.1109/4235.585893

827 Xia, G.S., Delon, J., Gousseau, Y., 2010. Shape-based invariant texture indexing. *Int.*  
828 *J. Comput. Vis.* 88, 382–403. doi:10.1007/s11263-009-0312-3

829 Yang, W., Dai, D., Triggs, B., Xia, G.S., 2012. SAR-based terrain classification using  
830 weakly supervised hierarchical Markov aspect models. *IEEE Trans. Image*  
831 *Process.* 21, 4232–4243. doi:10.1109/TIP.2012.2199127

832 Yang, W., Yin, X., Xia, G.S., 2015. Learning high-level features for satellite image  
833 classification with limited labeled samples. *IEEE Trans. Geosci. Remote Sens.*

834 53, 4472–4482. doi:10.1109/TGRS.2015.2400449

835 Yang, X., Qian, X., Mei, T., 2015. Learning salient visual word for scalable mobile  
836 image retrieval. *Pattern Recognit.* 48, 3093–3101.  
837 doi:10.1016/j.patcog.2014.12.017

838 Yin, W., Yang, J., Yamamoto, H., Li, C., 2015. Object-based larch tree-crown  
839 delineation using high-resolution satellite imagery. *Int. J. Remote Sens.* 36, 822–  
840 844. doi:10.1080/01431161.2014.999165

841 Yu, J., Weng, K., Liang, G., Xie, G., 2013. A vision-based robotic grasping system  
842 using deep learning for 3D object recognition and pose estimation, in: 2013  
843 IEEE International Conference on Robotics and Biomimetics (ROBIO). pp.  
844 1175–1180. doi:10.1109/ROBIO.2013.6739623

845 Yue, J., Mao, S., Li, M., 2016. A deep learning framework for hyperspectral image  
846 classification using spatial pyramid pooling. *Remote Sens. Lett.* 7, 875–884.  
847 doi:10.1080/2150704X.2016.1193793

848 Zhang, C., Kovacs, J.M., 2012. The application of small unmanned aerial systems for  
849 precision agriculture: A review. *Precis. Agric.* 13, 693–712. doi:10.1007/s11119-  
850 012-9274-5

851 Zhang, C., Wang, T., Atkinson, P.M., Pan, X., Li, H., 2015. A novel multi-parameter  
852 support vector machine for image classification. *Int. J. Remote Sens.* 36, 1890–  
853 1906. doi:10.1080/01431161.2015.1029096

854 Zhang, F., Du, B., Zhang, L., 2016. Scene classification via a gradient boosting  
855 random convolutional network framework. *IEEE Trans. Geosci. Remote Sens.*  
856 54, 1793–1802. doi:10.1109/TGRS.2015.2488681

857 Zhang, Q., Wang, J., Gong, P., Shi, P., 2003. Study of urban spatial patterns from  
858 SPOT panchromatic imagery using textural analysis. *Int. J. Remote Sens.* 24,  
859 4137–4160. doi:10.1080/0143116031000070445

860 Zhao, W., Du, S., 2016. Learning multiscale and deep representations for classifying  
861 remotely sensed imagery. *ISPRS J. Photogramm. Remote Sens.* 113, 155–165.  
862 doi:10.1016/j.isprsjprs.2016.01.004



863 Zhong, Y., Zhao, J., Zhang, L., 2014. A hybrid object-oriented conditional random  
864 field classification framework for high spatial resolution remote sensing imagery.  
865 IEEE Trans. Geosci. Remote Sens. 52, 7023–7037.  
866 doi:10.1109/TGRS.2014.2306692

867