

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

# Very fine spatial resolution urban land cover mapping using an explicable sub-pixel mapping network based on learnable spatial correlation

Da He<sup>a,b</sup>, Qian Shi<sup>a,b,\*</sup>, Jingqian Xue<sup>a</sup>, Peter M. Atkinson<sup>c</sup>, Xiaoping Liu<sup>a,b</sup>

<sup>a</sup>School of Geography and Planning, Sun Yat-sen University, Guangzhou 510275, China

<sup>b</sup>Guangdong Provincial Key Laboratory for Urbanization and Geo-simulation, Guangzhou 510275, China

<sup>c</sup>Faculty of Science and Technology, Lancaster University, Lancaster LA1 4YR, UK

**Abstract:**

Sub-pixel mapping is the prevailing approach for dealing with the mixed pixel effect in urban land use/land cover classification, by reconstructing the sub-pixel-scale distribution inside each mixed-pixel based on spatial autocorrelation. However, 1) traditional spatial autocorrelation is limited to a local window, which cannot model the teleconnection between two locations or objects that are far apart and 2) autocorrelation is based on the idea of “the more proximate, the more similar”, which relies on a distance-weight decay parameter and cannot characterize the rich variety of mutual information in spatially heterogenous areas in urban. In this research, we develop and demonstrate a learnable correlation-based sub-pixel mapping (LECOS) method. 1) We use the “mutual retrieval” mechanism of the self-attention operation to model teleconnections that enable more distant locations or objects to be mutually correlated and 2) we design a parameter-free “self-attention in self-attention” operation to learn

---

\* Corresponding author: shixi5@mail.sysu.edu.cn

26 adaptively the diverse global correlation patterns between pixel and sub-pixel. The  
27 learned spatial correlations are then used for reasoning the sub-pixel-scale distribution  
28 of each class. We validated our method on the most challenging public datasets of urban  
29 scenes, which exhibit considerable spatial heterogeneity with complex structures and  
30 broken objects. The learned building-tree, building-road and road-tree correlation  
31 patterns contributed most to the sub-pixel reconstruction result of the urban scenes,  
32 consistent with *in-situ* reference data. We further explored the model's explicability in  
33 a large-area of several metropolises in China, by mapping land cover in these cities at  
34 a 2 m very fine spatial resolution using 10 m Sentinel-2 input images, and found that  
35 the derived result not only revealed rich urban spatial heterogeneity, but also that the  
36 learned correlation was indicative of urban pattern dynamics, suggesting the potential  
37 for greater understanding of issues such as urban fairness, accessibility, human  
38 exposure and sustainability.

39

40 Keywords: land use/land cover classification, urban spatial pattern, sub-pixel mapping,  
41 spatial teleconnection, self-attention mechanism

42

## 43 **1. Introduction**

44 Urban land use/land cover can provide fundamental information for understanding  
45 the impact of changes in urban composition on the urban heat island phenomenon (Xia  
46 et al., 2022), air pollution exposure (IPCC Climate Change 2013), carbon losses (Foley  
47 et al., 2005) and urban sustainability (Stokes and Seto, 2018). Accurate urban land  
48 use/land cover monitoring and the corresponding regulation can alleviate the impact of

49 human activities on the climate, land surface and biodiversity (Liu et al., 2020; Liu et  
50 al., 2019). However, urban scenes are highly dynamic and heterogenous, which makes  
51 the widely used moderate spatial resolution remote sensing images (e.g.,  
52 Landsat/Sentinel-2) unsuitable for fine-grained monitoring, since the edges of buildings,  
53 green parcels and other small-sized urban components will inevitably cut through pixels,  
54 leading to a jagged boundary or even disappearance at a coarse spatial resolution,  
55 known as the mixed pixel effect. Traditional hard classification methods that assign  
56 each pixel to a single label may result in small-sized and elongated objects being missed.  
57 Although many high-resolution images are available now, such as World-View series  
58 or GF series satellite, they generally have a long revisiting period and narrow width,  
59 and are often obscured by clouds with year-round rainfall in southern China. It is  
60 difficult to mosaick a complete city-level or provincial-level image without cloud  
61 contamination, which is difficult to meet the needs of periodic continuous monitoring.

62         Researchers have addressed the mixed pixel problem with the spectral unmixing  
63 technique. Instead of hard classification, spectral unmixing aims to decompose each  
64 mixed pixel into a combination of multiple pure land covers weighted by their  
65 proportions (Adam et al., 1995), thereby, estimating the percentage of each land cover  
66 within each mixed pixel. Spectral unmixing has been applied widely to derive many  
67 land cover fractions, such as fractional forest cover (Rashed et al., 2003; Small, 2003;  
68 Xiao and Moody, 2005) and fractional urban impervious cover (Wu and Murray, 2003;  
69 Deng et al., 2019). However, the spatial resolution of the fractional images is still coarse,  
70 providing no indication of how each pure land cover is located inside each mixed pixel

71 (Atkinson, 2009).

72 Sub-pixel mapping (SPM) is an effective method for estimating the location of  
73 each pure land cover inside each mixed pixel, which is equivalent to improving the  
74 observation resolution. It has been applied widely for the identification of small-sized  
75 objects that would otherwise require a very fine resolution, for example, 0.6 m spatial  
76 resolution for urban tree reconstruction from 2.4 m QuickBird MS images (Ardila et al.,  
77 2011), 2 m resolution for urban tree reconstruction from 10 m Sentinel-2A MSI images  
78 (He et al., 2022a), and 0.5 m resolution for individual potato plant reconstruction from  
79 2 m WorldView-2 images (Poudyal, 2013).

80 The SPM grid divides each mixed pixel into several sub-pixels, and allocates land  
81 cover labels (derived from spectral unmixing) to each sub-pixel location to realize a  
82 finer land cover map. Considering that the allocation process is ill-posed, a spatial prior  
83 is necessary to describe the spatial correlation of each sub-pixel and constrain the SPM  
84 solution. According to the formulation of the prior, SPM can be categorized into three  
85 types:

86 1) Spatial dependence prior. This assumes that the adjacent sub-pixel/pixel  
87 locations are likely to belong to the same class, and each sub-pixel's label can be  
88 predicted according to the abundances of its eight-neighboring mixed-pixels or sub-  
89 pixels, such as the sub-pixel swapping model (Atkinson, 2005), spatial attraction model  
90 (Mertens et al., 2006), genetic algorithm based SPM (He et al., 2016a), Hopfield neural  
91 network based SPM (Su, 2019) and random field based SPM (Kasetkasem et al., 2005).

92 2) Geostatistical prior. This assumes that the sub-pixel distribution is subject to

93 empirical geostatistical models that can be regressed with given samples (Boucher and  
94 Kyriakidis 2006), such as area-to-point kriging-based SPM (Wang et al., 2015) and  
95 spatial distribution pattern-based SRM (Ge et al. 2016).

96 3) Regularization prior. This assumes that handcrafted filters like the Laplacian  
97 model, non-local model or sparse assumption can regularize the sub-pixel distribution,  
98 which includes maximum *a posteriori* model-based SPM (Zhong et al., 2015), sparse  
99 representation model-based SPM (Song et al., 2019), spectral-spatial fusion-based SPM  
100 (Xu et al. 2018), and spatial-temporal-spectral fusion-based SPM (Li et al., 2017).  
101 However, the above model-driven spatial priors are essentially based on the idea that  
102 “more proximate data points are expected to be more similar”, which is not so readily  
103 extensible to frequent spatial variation as commonly encountered in urban scenarios.

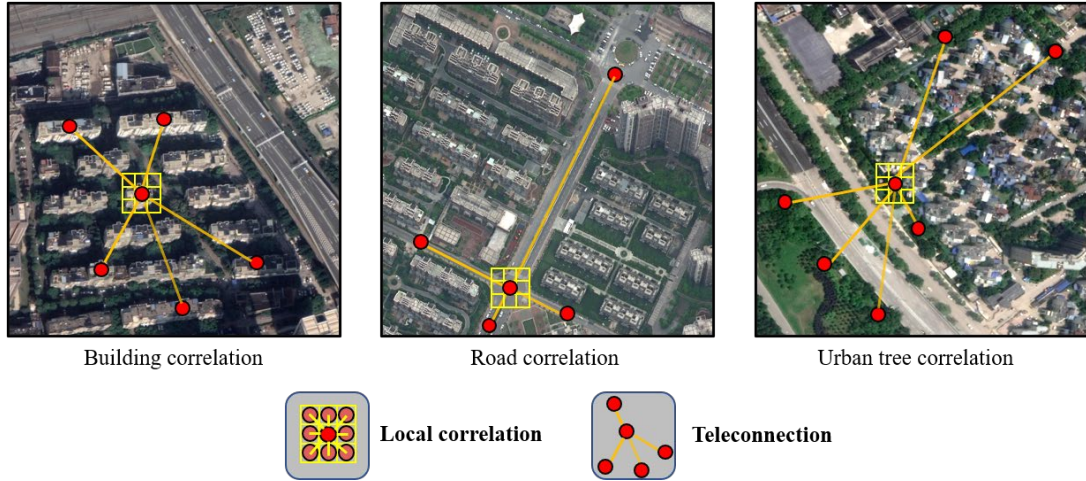
104 Deep learning theory is an alternative method that can alleviate the requirement  
105 for an empirical spatial prior, focused on “allowing the data to speak for themselves”  
106 (Ling et al., 2019; He et al., 2021a). Deep learning approaches assume that the prior  
107 can be learned from exemplar pairs of coarse resolution (CR) image and fine resolution  
108 (FR) annotation images in actual geographical scenarios. Deep learning-based SPM  
109 generally uses the convolution and deconvolution operations to reconstruct an FR result,  
110 and a FR ground reference is compared with the estimates and the gradient error is  
111 backpropagated to update the network parameters, and learn the underlying spatial  
112 correlation. Many prevalent networks were applied for SPM, such as the super-  
113 resolution reconstruction network (Ling et al. 2019a; 2019b; Ma et al. 2019; He et al.  
114 2021a; Shang et al. 2020), semantic segmentation network (Arun et al. 2018; Jia et al.

115 2019), attention mechanism-based network (He et al. 2021b; 2022), graph convolution  
116 network (Zhang et al. 2021b) and spatiotemporal fusion network (Chen et al. 2021b).  
117 However, the convolutional layer usually has a limited receptive field due to the local  
118 connection property (e.g.,  $3 \times 3$  kernel size), which cannot learn the non-local spatial  
119 correlation between urban compositions that is essential for sub-pixel location  
120 reasoning (Zhang et al., 2019; Verdonck et al., 2017). Besides this, the learning process  
121 of the convolution network cannot explain what spatial correlation pattern the network  
122 has learned, and how it guides inference of the sub-pixel-scale location.

123       Based on the above, two challenges follow: 1) all the above spatial correlation  
124 approaches are limited to a local region and, thus, cannot exploit the teleconnections  
125 that are more characteristic of spatial heterogeneity (e.g., buildings and trees repeatedly  
126 occur at a certain distance, exhibiting a correlation between each other remotely (**Fig.**  
127 **1**)). Although non-local convolution was developed to capture the global context  
128 information by measuring each point of the feature map with the weighted average of  
129 the other points (Wang et al., 2018), and dilated convolution can also expand the  
130 receptive field of the traditional kernels by injecting blank pixels into  $3 \times 3$  kernels to  
131 achieve larger kernel size, while maintaining the calculation efficiency (Fisher et al.,  
132 2015). However, they only aim to enlarge the receptive field, and are still limited to  
133 convolution framework, i.e., the learning process of the kernel parameters is  
134 uncontrollable, which is difficult to explain the meaning of the learned parameters; 2)  
135 classical spatial autocorrelation relies on the empirical distance decay assumption,  
136 which cannot characterize the various sources of mutual information between urban

137 components, while the implicitly learned correlation in a convolutional network cannot  
138 provide a reasonable explanation and, thus, is unreliable.

139



**Fig. 1.** Comparison of local correlation and teleconnection between objects in urban an scenario.

140

141 To overcome the above two challenges and reveal urban patterns at high-resolution  
142 for planning or decision-making, we combine the learning ability of the data-driven  
143 idea with the spatial correlation modeling process, and develop a learnable correlation  
144 based sub-pixel mapping network (LECOS). 1) We use the “mutual retrieval”  
145 mechanism of the self-attention operation to excavate the diverse contextual correlation  
146 patterns, which models the global teleconnection between objects like trees, buildings,  
147 and roads that are far apart; 2) The mutual retrieval process is established in the end-to-  
148 end network architecture, which enables the spatial correlation patterns to be learned  
149 explicitly through data-driven regression; 3) Considering that the typical self-attention  
150 operation cannot establish the correlation (between pixels and sub-pixels) that is  
151 essential for sub-pixel location inference, we further designed a “self-attention in self-  
152 attention (SNS)” operation to enable exploration of the hierarchical correlations

153 between the two scales; 4) Since self-attention operation explicitly learns the spatial  
154 dependency between each pixel and sub-pixel in the global scene, so that it can figure  
155 out, in a given image patch, which object is related to which one, therefore, the learned  
156 spatial correlation is explicable. Besides, the learned correlation was found to be able  
157 to quantitatively reflect the urban spatial patterns, such as accessibility of greenspace  
158 to impervious area, and was used to evaluate urban environmental planning and  
159 resource distribution of several main cities in China.

160 The aims of this paper were to 1) develop a learnable correlation-based sub-pixel  
161 mapping network architecture for reconstructing a fine-grained urban distribution with  
162 rich spatial heterogeneity; 2) validate the proposed method on public urban scenario  
163 datasets, examining its performance in reconstructing spatial detail and the underlying  
164 correlation; 3) evaluate the method over large areas across several main metropolises  
165 in China, and reveal our findings on urban spatial patterns as reflected by the correlation  
166 between each urban composition at very fine spatial resolution, and validate the  
167 findings with multi-source products.

168

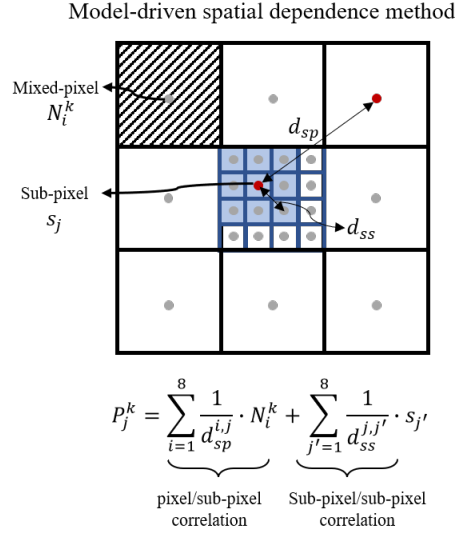
## 169 **2. Method**

### 170 *2.1 Foundation of spatial correlation modeling in sub-pixel mapping*

171 In model-driven SPM methods, given a CR image, the spectral unmixing approach  
172 is firstly used to decompose each mixed-pixel into pure land cover classes (i.e.,  
173 endmember) and estimate their pixel-level proportions. Then, each mixed-pixel is  
174 divided into  $S \times S$  sub-pixels ( $S$  is the scale factor), and a possible label is assigned to  
175 each sub-pixel location, based on the spatial correlation among the mixed-pixels and



176 the sub-pixels within the neighborhood. **Fig. 2** demonstrates the specific estimation  
 177 process with an example in a  $3 \times 3$  mixed-pixel area.



**Fig. 2.** The formulation of the model-driven spatial correlation modeling.

178

179 Supposing there are  $NC$  endmembers, the eight-neighborhood mixed-pixel is  
 180 denoted as  $N_i$ ,  $i = 1 \dots 8$ , and the proportion of the  $k$ th endmember in the mixed-pixel  
 181  $N_i$  is denoted as  $N_i^k$ ,  $k = 1 \dots NC$ . The sub-pixel in the center mixed-pixel area is  
 182 denoted as  $s_j$ ,  $j = 1 \dots S^2$ . The probability of assigning the  $k$ th endmember label to the  
 183  $j$ th sub-pixel (i.e.,  $P_j^k$ ) can be estimated by the proportions of the  $k$ th endmember in  
 184 each mixed-pixel and sub-pixel within the eight-neighborhood, as well as the spatial  
 185 correlation between them, which can be given as:

186

$$P_j^k = \sum_{i=1}^8 \frac{1}{d_{sp}^{i,j}} \cdot N_i^k + \sum_{j'=1}^8 \frac{1}{d_{ss}^{j,j'}} \cdot S_{j'} \quad (1)$$

187

188 where the spatial correlation is usually simplified to the inverse distance weighting in  
 189 Eq. (1), that is,  $d_{sp}^{i,j}$  represents the distance between the  $i$ th mixed-pixel and the  $j$ th

190 sub-pixel, and  $d_{ss}^{j,j'}$  represents the distance between the  $j$ th sub-pixel and the  $j'$ th sub-  
191 pixel in the eight-pixel neighborhood. After deriving all the class probabilities, the class  
192  $k$  corresponding to the largest probability  $P_j^k$  will be assigned to the  $j$ th sub-pixel  
193 location.

194 However, the spatial correlation is variable in practice, and may not be  
195 characterized sufficiently by the simple inverse distance assumption. Deep learning is  
196 a potential way of learning diverse spatial correlation patterns from exemplar datasets.  
197 However, how to model the spatial correlation in an end-to-end network remains an  
198 open question.

199

## 200 *2.2 Learnable correlation sub-pixel mapping*

201 To make full use of the learning ability of the data-driven idea and the correlation  
202 modelling ability of the model-driven idea, this research developed a learnable  
203 correlation sub-pixel mapping (LECOS) method. Instead of using convolution  
204 operation, the LECOS models the global spatial correlation between sub-pixel and  
205 mixed-pixel based on the “mutual retrieval” mechanism of the self-attention operation  
206 in an end-to-end network structure, to enable the spatial correlation learnable.  
207 Furthermore, considering that the typical self-attention operation can only establish the  
208 single pixel-level correlation, we further designed a “self-attention in self-attention”  
209 (SNS) operation to establish explicitly the spatial correlation rule between pixels and  
210 sub-pixels in a feed forward network fashion, which can be iteratively rectified by data-  
211 driven, back-propagation learning. Finally, the learned spatial correlation rule is used  
212 to infer the class label at each sub-pixel location.

213 The LECOS consists of the preparation work, encoder part and decoder part (Fig.  
 214 3). In the preparation work stage, *inductive bias* is designed to divide each mixed-pixel  
 215 of the input CR image into sub-pixels at the beginning of the network, which analogizes  
 216 the practice of the model-driven SPM that enables the spatial correlation calculation at  
 217 each sub-pixel location. Then, *image tokenization* is designed to convert each mixed-  
 218 pixel and sub-pixel into a 1-dimensional sequential mixed-pixel token and sub-pixel  
 219 token, which aims to accommodate the self-attention operation. In the encoder part, a  
 220 4-stage multi-scale residual structure is designed to consider the large scale-variation  
 221 of the urban compositions, which analogizes the structure of the typical ResNet (He et  
 222 al., 2016b). Each stage is composed of a stacked layer of *syntax-dependence builder*  
 223 and *context-dependence builder* based on the designed SNS operation, which is the  
 224 main part of the LECOS for learning the spatial correlation rule amongst the mixed-  
 225 pixels and sub-pixels. In the decoder part, the output features of the 4-stage are  
 226 reformed from sequential token to a 2-dimensional feature map, and are integrated and  
 227 recovered to the target spatial resolution. Finally, a classifier is used to estimate the  
 228 class probability from the integrated feature maps.

229

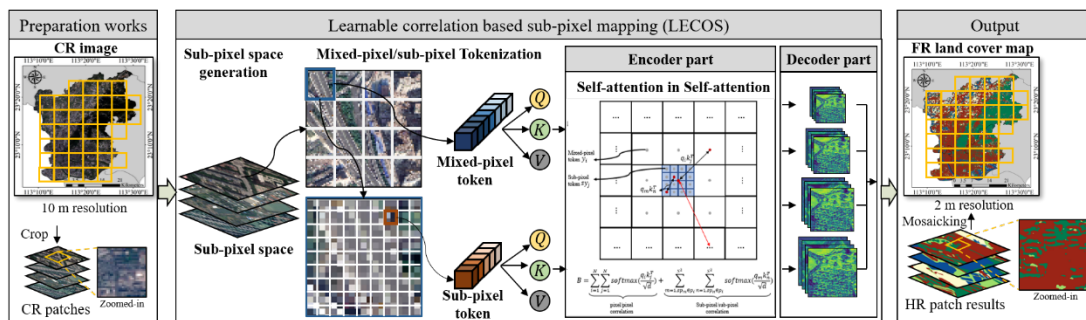


Fig. 3. Overall architecture of LECOS.

230

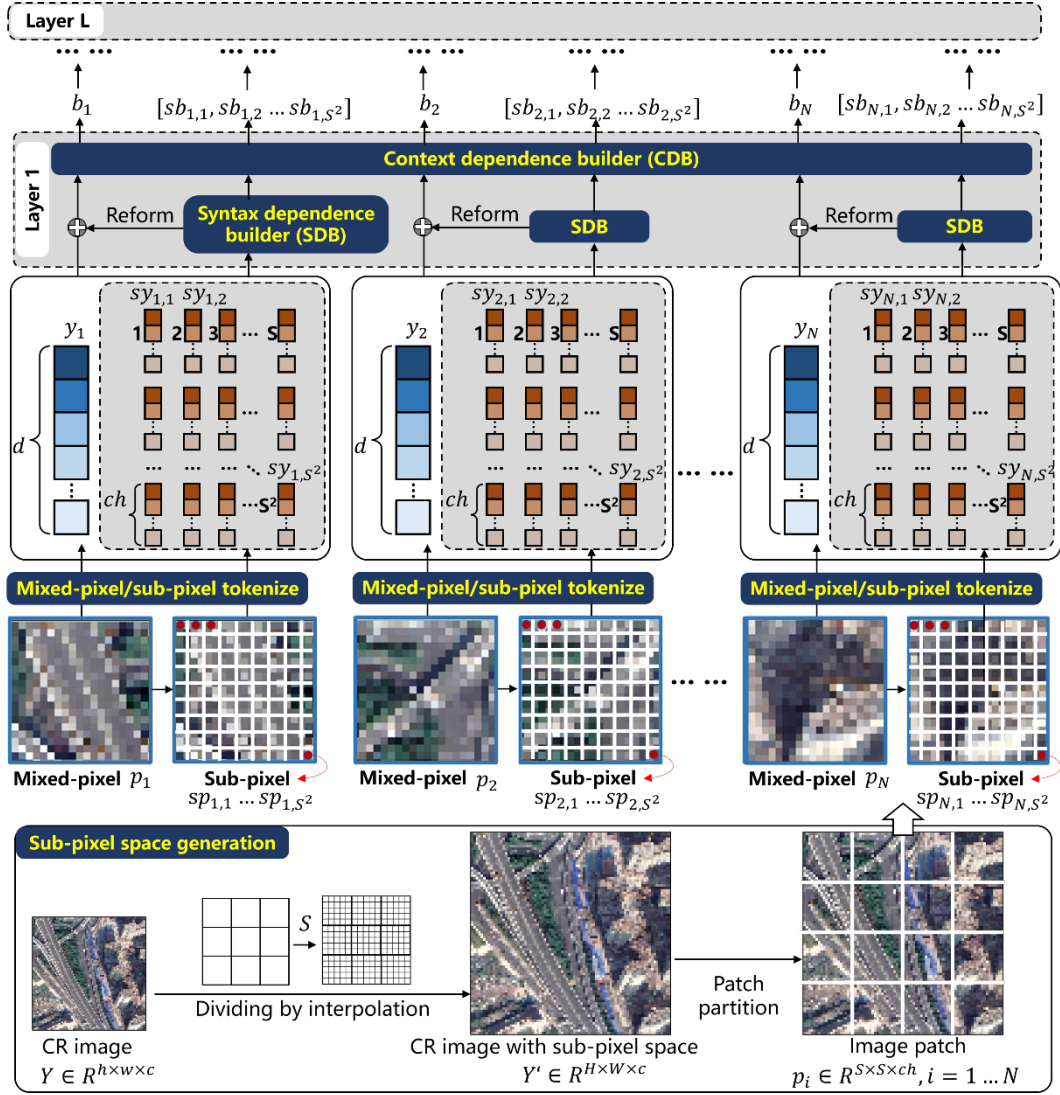
231 *2.2.1 Sub-pixel space generation*

232 Given an CR image  $Y \in R^{h \times w \times c}$ , it is upsampled by a scale factor  $S$  (denoted as  
233  $Y' \in R^{H \times W \times c}$ ) in the first layer of the network to generate the sub-pixel space, which  
234 is equivalent to the division of the mixed-pixel into sub-pixels in the classical spatial  
235 dependence method (**Fig. 4**). In this way, the sub-pixels can be manipulated explicitly  
236 in the subsequent “mutual retrieval” process to explore the potential correlation. This  
237 design is different from the conventional learning-based SPM that usually deploys the  
238 upsampling in the last layer, which is suitable to the purpose of establishing the  
239 correlation at each sub-pixel location.

240

241 *2.2.2 Mixed-pixel/sub-pixel tokenization*

242 Image tokenization is constructed to transform the image into sequential tokens  
243 (**Fig. 4**). Specifically, the input image is grid partitioned to patches  $p_i \in R^{S \times S \times ch}$ ,  $i =$   
244  $1 \dots N$  with size of  $S \times S$  (where  $N = H \times W / S \times S$ ,  $ch$  is the channel number), which  
245 is equivalent to the area of one mixed-pixel. Each patch is then flattened to a one-  
246 dimension vector to generate a mixed-pixel token  $y_i \in R^{d \times 1}$ ,  $i = 1 \dots N$  (where  $d =$   
247  $S \times S \times ch$ ). Each sub-pixel  $s_{i,m} \in R^{ch \times 1}$ ,  $m = 1 \dots S^2$  in the patch  $p_i$  can be also  
248 viewed as a 1-dimensional vector, which forms a sub-pixel token  $sy_{i,m} \in R^{ch \times 1}$ ,  $m =$   
249  $1 \dots S^2$ .



**Fig. 4.** Demonstration of the specific process in the preparation work.

250

### 251 2.2.3 Self-attention in self-attention (SNS) operation

#### 252 2.2.3.1 Syntax and context in urban scenario

253 How to establish the hierarchical correlation amongst mixed-pixels and sub-pixels

254 in the network is key to our method. In this research, we constructed a multi-scale

255 observation framework inspired by the language principle. Specifically, sub-pixels can

256 be regarded as words, which can be aggregated to form a mixed-pixel (sentence) based

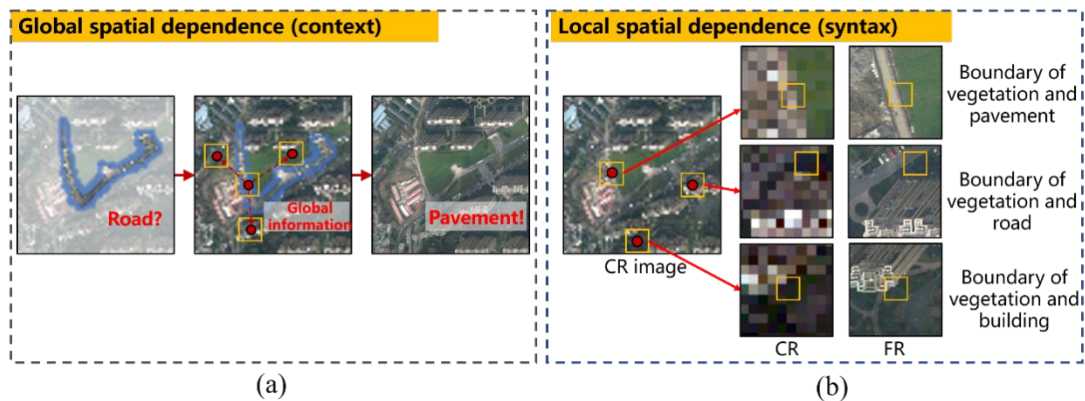
257 on syntax. Syntax can indicate the location of the sub-pixels inside each mixed-pixel,

258 such as linear distribution, sparse distribution, etc., according to the morphological

259 shape of the object. On the other hand, mixed-pixels can be aggregated to form a scene  
 260 (paragraph) based on context, which indicates the distribution rule of the mixed-pixel  
 261 in the scene, that is, trees located along the road, plane located inside the airport, etc.

262 An intuitive example of how context and syntax work in the urban scene is  
 263 provided in **Fig. 5**. In **Fig. 5a**, for example, it is difficult to recognize the highlighted  
 264 object with its context masked, which could be mistaken for a road. When providing  
 265 the context information, such as grassland and neatly arranged buildings, it is easy to  
 266 determine that here is a residential area, and the highlighted object is the pavement. On  
 267 consideration that the mixed-pixel is too coarse to describe the tiny pavement, the  
 268 syntax information can be used to determine the sub-pixel location at the place of  
 269 intersection, according to the morphological shape of each object, and guide the  
 270 network to reconstruct linearly distributed sub-pixels for the jagged pavement (**Fig. 5b**).

271



**Fig. 5.** Demonstration of the role of context and syntax in inferring urban compositions.

272

### 273 2.2.3.2 Syntax- and context-dependence builder

274 To learn the syntax and context, respectively, for mixed-pixels and sub-pixels, we

275 designed a syntax-dependence builder (SDB) and context-dependence builder (CDB)  
276 (**Fig. 6**), using the self-attention in self-attention (SNS) operation. The SDB aims to  
277 learn the syntax between sub-pixels, which represents local features, such as  
278 morphological shape, detail edge, etc. After the syntax is learned, the sub-pixels are  
279 aggregated by syntax to formulate a complete mixed-pixel (sentence), and the CDB is  
280 used to learn the context amongst mixed-pixels. Finally, all the mixed-pixels are  
281 aggregated by context to formulate a complete scene.

282 The self-attention operation is the critical technique to build the relationship  
283 amongst mixed-pixel tokens and sub-pixel tokens, inspired by the human vision system  
284 that is able to guide attention towards interested objects through “mutual retrieval”  
285 between two components within a scene. For example, given an object to be retrieved  
286 as a *Query* (e.g., pavement), the pavement-related objects (e.g., grassland, building, etc.)  
287 in the scene are taken as *Key*, and the intrinsic features of each object are taken as *Value*.  
288 The self-attention operation builds the relationship between *Key* and *Query*, to measure  
289 the correlation intensity of each component related to the pavement, and the intensity  
290 is used as the attention weight to rescale the *Value* of all the objects, which are finally  
291 aggregated to obtain the global feature embedded with correlation information.  
292 However, the classical self-attention operation cannot establish the correlation between  
293 pixel-level and sub-pixel-level, which inspired us to design the SNS operation to  
294 support SDB and CDB.

295 The architectures of the SDB and the CDB are shown in **Fig. 6**. Taking the first  
296 layer of the encoder part as an example, given the mixed-pixel token  $y_i \in R^{d \times 1}, i =$

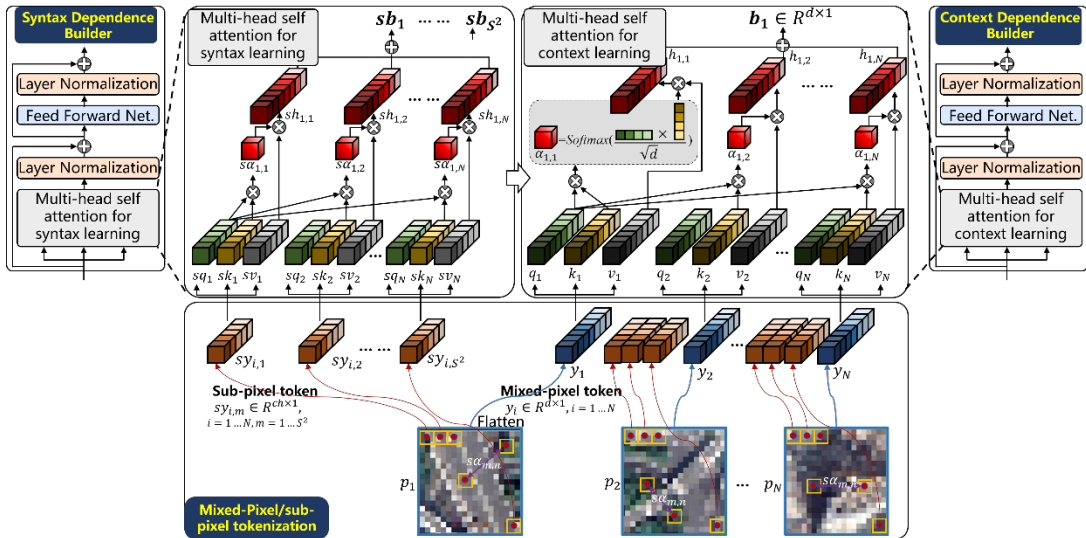
297  $1 \dots N$  and the sub-pixel token  $sy_{i,m} \in R^{ch \times 1}, m = 1 \dots S^2$  derived from image  
 298 tokenization, the SNS operation is used to investigate the context correlation of the  $ith$   
 299 token  $y_i$  with the  $jth$  token  $y_j$  ( $j = 1 \dots N$ ), and the syntax correlation of the  $mth$  token  
 300  $sy_{i,m}$  with the  $nth$  token  $sy_{i,n}$  ( $n = 1 \dots S^2$ ). To enable the sub-pixels to retrieve each  
 301 other, and to be retrieved, the *Query* ( $SQ = \{sq_{i,m}\} \in R^{ch \times S^2}$ ), *Key* ( $SK = \{sk_{i,m}\} \in$   
 302  $R^{ch \times S^2}$ ) and *Value* ( $SV = \{sv_{i,m}\} \in R^{ch \times S^2}$ ) of the sub-pixel token  $sy_{i,m}$  within  $y_i$ , are  
 303 generated by linear transformation of the  $y_i$  and  $sy_{i,m}$ , respectively:

$$\begin{aligned}
 sk_{i,m} &= sW^k \cdot sy_{i,m}, \\
 sq_{i,m} &= sW^q \cdot sy_{i,m}, \\
 sv_{i,m} &= sW^v \cdot sy_{i,m}, \quad m = 1 \dots S^2
 \end{aligned} \tag{2}$$

304

305  
 306 where  $W^q, W^k, W^v \in R^{d \times d}$ ,  $sW^k, sW^q, sW^v \in R^{c \times c}$  represent the transformation  
 307 matrices of *Query*, *Key* and *Value*, respectively.

308



**Fig. 6.** Syntax-dependence and context-dependence builder based on self-attention in self-attention operation.

309



310 Firstly, the *Query* of the sub-pixel  $sy_{i,m}$  is taken to retrieve the *Key* of the global  
 311 sub-pixel  $sy_n$  to build the syntax correlation  $s\alpha_{i,m,n} \in R^{1 \times 1}$ , which can be given as:

$$312 \quad s\alpha_{i,m,n} = \sum_{m=1, sy_{i,m} \in y_i}^{S^2} \sum_{n=1, sy_{i,n} \in y_i}^{S^2} \text{softmax}\left(\frac{q_{i,m} k_{i,n}^T}{\sqrt{ch}}\right) \quad (3)$$

313

314 where  $\sqrt{ch}$  normalizes the  $q_{i,m} k_{i,n}^T$ .

315 The derived syntax correlation is used to rescale the *Value* of each sub-pixel token  
 316  $sy_{i,n}$  within the  $y_i$ , and the output syntax feature of  $sy_{i,m}$  is obtained by summarizing  
 317 all the rescaled *Values*:

318

$$319 \quad sb_{i,m} = \sum_{n=1}^{S^2} sv_{i,m} \cdot s\alpha_{i,m,n}, m = 1 \dots S^2 \quad (4)$$

319

320 where  $sb_{i,m} \in R^{c \times 1}$  is the derived syntax feature of the sub-pixel token  $sy_{i,m}$ .

321 The above syntax-dependence modeling process can be given in matrix form:

322

$$323 \quad SB_i = \mathbf{SDB}(\{sy_{i,m}\}_{m=1 \dots S^2}) = SV \cdot \text{SoftMax}(SK^T \cdot SQ / \sqrt{ch}) \quad (5)$$

323

324 where  $SB_i \in R^{ch \times S^2}$  is the aggregation of the sub-pixel token rescaled by the syntax.

325 Then, the  $SB_i$  are flattened to a 1-dimensional vector with size  $(ch \times S^2) \times 1$ , and  
 326 propagated to the mixed-pixel token  $y_i$  by summarizing:

327

$$328 \quad y'_i = y_i + \text{flatten}(SB_i) \quad (6)$$

328

329 In the same way, the *Query* ( $Q = \{q_i\} \in R^{d \times N}$ ), *Key* ( $K = \{k_i\} \in R^{d \times N}$ ), and

330 *Value* ( $V = \{v_i\} \in R^{d \times N}$ ) of the pixel token  $y'_i$  is firstly generated:

331

$$\begin{aligned} q_i &= W^q \cdot y'_i, \\ k_i &= W^k \cdot y'_i, \\ v_i &= W^v \cdot y'_i, \quad i = 1 \dots N \end{aligned} \quad (7)$$

332

333 The context correlation  $\alpha_{i,j} \in R^{1 \times 1}$  is also derived by the self-attention operation:

334

$$\alpha_{i,j} = \sum_{i=1}^N \sum_{j=1}^N \text{softmax}\left(\frac{q_i k_j^T}{\sqrt{d}}\right) \quad (8)$$

335

336 And the output features of the mixed-pixel token  $y_i$  that are embedded with the

337 syntax- and context- correlation can be obtained:

338

$$\hat{B} = \mathbf{CDB}(\{y'_i\}_{i=1 \dots N}) = V \cdot \text{SoftMax}(K^T \cdot Q / \sqrt{d}) \quad (9)$$

339

340 where  $\hat{B} \in R^{d \times N}$  represents the ensemble of the mixed-pixel token that is organized by

341 context.

342 Noteworthy is that all the  $\alpha_{i,m,n}$  can be formulated as an attention matrix  $SA \in$

343  $R^{S^2 \times S^2}$ , which is the desired syntax-dependence, and all the  $\alpha_{i,j}$  can be formulated as a

344 matrix  $A \in R^{N \times N}$ , which is the desired context-dependence. Both the syntax and

345 context are visualized in the experimental results (Section 4.1), to demonstrate

346 explicitly the learned correlation rule in the urban scenario. From this point of view, the

347 process of CDB and SDB is consistent with the spatial dependence modelling practice  
 348 of the classical model-driven SPM algorithm, both of which explicitly explore the  
 349 optimal spatial configuration at the sub-pixel-scale. Furthermore, benefitting from the  
 350 global “mutual retrieval” ability of SNS operation, the range of the correlation is no  
 351 longer limited to the eight-neighborhood, but the whole scene, as shown in **Fig. 3**, and  
 352 the Eq. (1) can be reformed to:

$$\hat{B} = \sum_{i=1}^N \sum_{j=1}^N (\text{softmax}(\frac{q_i k_j^T}{\sqrt{d}})) + \sum_{m=1, sp_m \in p_i}^{S^2} \sum_{n=1, sp_n \in p_i}^{S^2} \text{softmax}(\frac{q_m k_n^T}{\sqrt{d}}) \quad (10)$$

354

355 To ensure robust training and accelerate convergence, a layer normalization  
 356 (LayerNorm) operation is adopted to normalize the derived syntax  $SB_i$  and context  $B_i$   
 357 to a Gaussian distribution with zero mean and variance 1 (as shown in **Fig. 6**), given as:

358

$$\text{LayerNorm}(B_i) = g_{ln} \times \frac{\hat{B}_i - \mu(\hat{B}_i)}{\sigma(B_i)} + b_{ln} \quad (11)$$

359

360 where  $g_{ln}$  is the gain parameter, and  $b_{ln}$  is the bias, which are also trainable.

361 A feed forward network (FFN) is then used for each sequential token feature to  
 362 enhance the fitting capability, which adopts the multiple layer perceptron (MLP), as  
 363 was done by Vaswani et al. (2017).

364

$$\text{FFN}(B_i) = W_f \otimes \text{LayerNorm}(\hat{B}_i) + b_f \quad (12)$$

365

366 where  $W_f$  and  $b_f$  represent the  $1 \times 1$  kernel weight and bias, and  $\otimes$  means the  
 367 convolution operation.

368 In summary, the above process can be reorganized by a residual connection  
 369 strategy:

$$370 \quad \hat{B} = Y + \mathbf{CBD}(\mathbf{SBD}(Y)) \quad (13)$$

$$B = \hat{B} + \mathbf{FFN}(\hat{B}) \quad (14)$$

371

372 Note that the size of the output features  $SB$  and  $B$  equal the input sub-pixel token  
 373 and mixed-pixel, that is,  $ch \times S^2$  and  $d \times N$ , which can be used as the input for the next  
 374 layer of SDB and CDB, and enable the SNS layers to be stacked to form an end-to-end  
 375 feed forward architecture.

376

#### 377 2.2.4 Encoder and decoder structure

378 Considering the large scale-variation of urban compositions, a 4-stage multi-scale  
 379 structure is designed in the encoder part to learn the hierarchical representations of the  
 380 syntax and context correlation features (**Fig. 7**), which is analogous to the structure of  
 381 the typical ResNet (He et al., 2016b). Each stage is composed of stacked layers of SBD  
 382 and CBD. The patch merging strategy is designed between each stage to generate multi-  
 383 scale features. After that, the output features of the 4-stage are aggregated in the decoder  
 384 part, and a *Softmax* classifier is adopted to generate the classification result.

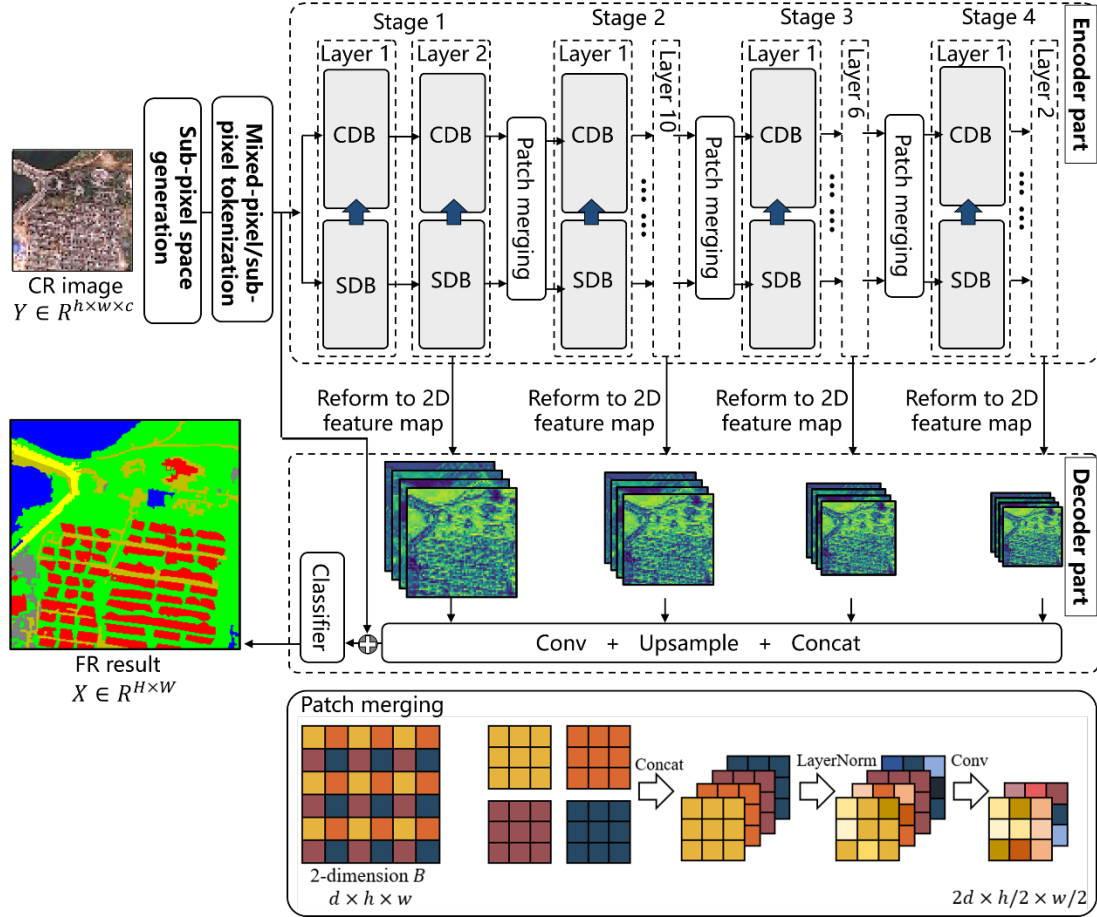


Fig. 7. The network architecture of LECOS.

385

#### 386 2.2.4.1 Patch merging

387 Patch merging is used to create multi-scale features for each stage. Taking the  
 388 output feature of  $B$  from stage-1 as an example, a reshape operation is firstly applied to  
 389 resize  $B$  from 1-dimension  $d \times N$  to 2-dimensions  $d \times h \times w$ . Then, the 2-dimensional  
 390 feature maps are split into four equal parts according to Fig. 7, and are concatenated  
 391 along with the channel dimension, equivalent to a downsampling process. Finally, the  
 392 derived 2-dimensional feature maps are again flattened to a 1-dimensional vector for  
 393 the next stage. In this way, every time the  $B$  pass a stage, the height and width are  
 394 reduced to half, and the channel dimensions are doubled, that is,  $F_1 \in R^{H \times W \times C}$ ,  $F_2 \in$   
 395  $R^{H/2 \times W/2 \times 2C}$ ,  $F_3 \in R^{H/4 \times W/4 \times 8C}$ ,  $F_4 \in R^{H/8 \times W/8 \times 16C}$ . The reason for this design is

396 that the distance between two highly correlated urban compositions is various, which  
 397 cannot be fully captured by a fixed patch size, and an ablation experiment is also  
 398 designed to validate the multi-scale architecture (refer to the supplementary file).

399

#### 400 2.2.4.1 Feature integration

401 In the decoder part, given the output feature  $F1, F2, F3, F4$  from each stage, the  
 402 pixel-shuffle operation (Shi et al., 2016) is used to upsample them to the target size  
 403  $H \times W \times C$ , and they are integrated by concatenation:

404

$$[F2', F3', F4'] = \text{Upsample}([F2, F3, F4]) \quad (15)$$

$$F_{ms} = \text{Concat}([F1, F2', F3', F4']) \quad (16)$$

405

406 where  $F_{ms} \in R^{H \times W \times 4C}$  represents the integrated features embedded with multi-scale  
 407 information.

408 A *SoftMax* classifier is used to classify the feature map  $F_{ms}$  to the FR land cover  
 409 class probability map  $P \in R^{H \times W \times NC}$  for  $NC$  categories:

410

$$F = \sigma \cdot (W_c \otimes F_{ms} + b_c) \quad (17)$$

$$P = \sum_{i=1}^H \sum_{j=1}^W \frac{\exp(F_{i,j,k})}{\sum_{k=1}^{NC} \exp(F_{i,j,k})} \quad (18)$$

411

412 where  $F \in R^{H \times W \times NC}$  is the intermediate feature map, and the  $F_{i,j,k}$  means the feature  
 413 value at  $(i, j)$  location and  $k$ th channel,  $W_c$  and  $b_c$  are weight and bias, and  $\sigma$  is the  
 414 rectified linear unit (ReLU).

415

### 416 2.2.5 Loss function

417 The loss function aims to measure the differences between the prediction  $P$  and  
418 the ground reference probability map  $Z \in R^{H \times W \times NC}$ . The differences are  
419 backpropagated to update the network parameter  $\theta$  by minimizing the loss function  
420 during the training process. Given pairs of  $\{Y, Z\}$ , the optimal network model can be  
421 learned as:

422

$$\{\hat{P}, \theta\} = \arg \min_{\theta} \{Loss(P, Z)\} \quad (19)$$

423

424 Considering the commonly class unbalanced distribution of urban compositions,  
425 the combination of a focal loss function  $L_{focal}$  and dice loss function  $L_{dice}$  is adopted  
426 in our approach (Milletari et al., 2016; Lin et al., 2020).  $L_{focal}$  aims to reduce the weight  
427 of easy samples and enlarge the weight of hard samples, and  $L_{dice}$  aims to maximize  
428 the intersection ratio between the prediction and ground reference, given as:

429

$$L_{focal} = -\alpha(1 - P)^{\gamma} \cdot Z \cdot \log P - (1 - \alpha)P^{\gamma} \cdot (1 - Z) \cdot \log(1 - P) \quad (20)$$

$$L_{dice} = 1 - \frac{2(P \cdot Z)}{P + Z} \quad (21)$$

$$Loss(P, Z) = (1 - \mu) \cdot L_{focal} + \mu \cdot L_{dice} \quad (22)$$

430

431 where  $\alpha \in [0, 1]$  is the balance weight between easy samples and hard samples, and  $\mu \in$   
432  $[0, 1]$  is the balance weight between  $L_{focal}$  and  $L_{dice}$ .

433

### 434 **3. Material/Study area**

#### 435 *3.1 Public dataset*

436 Three public urban scenario datasets were selected to examine the capabilities for  
437 detail reconstruction and correlation excavation of our method (**Fig. 8**). Noteworthy is  
438 that pairs of CR image and FR annotation images are required for SPM network training.  
439 However, due to the difficulty of obtaining matching pairs acquired on exactly the same  
440 date, we instead used the semantic segmentation dataset (WHDL D and LoveDA) to  
441 ensure that the image and its annotation image were matched, and further downsampled  
442 the image by a predefined scale factor to derive a synthetic CR image as the input for  
443 the SPM network. Furthermore, a newly published dataset specifically for SPM tasks  
444 (i.e., the FLAS dataset) was also adopted, which has matched pairs of CR Sentinel-2  
445 image and FR annotated Google images.

446 The WHDL D dataset (Shao et al., 2020) consists of pairs of 2 m Google Earth  
447 images and the corresponding annotation images, located in the urban area of Wuhan,  
448 China. WHDL D contains 4940 image patch pairs with a size of 256×256 pixels and  
449 three bands, densely annotated with six land cover classes (i.e., building, pavement,  
450 vegetation, bare soil and water). After downsampling by four times the original images,  
451 the spatial resolution was converted from 2 m to 8 m and the patch size changed to  
452 64×64 pixels, to serve as the CR image.

453 The LoveDA dataset (Wang et al., 2021) consists of pairs of 0.3 m spaceborne  
454 images and their corresponding annotation images, which are located in the urban and  
455 rural areas of Nanjing, Changzhou and Wuhan in July 2016, with a complex background



456 and rich details. LoveDA contains 5987 image patch pairs with a size of 1024×1024  
457 pixels and three bands, densely labelled with seven classes (i.e., background, road,  
458 water, barren, forest, agriculture and background). The image and annotation of the  
459 LoveDA were firstly resampled to 1.2 m spatial resolution (256×256 pixels), and the  
460 image was further downsampled to 4.8 m resolution (64×64 pixels) to serve as the CR  
461 image. The reason is that 0.3 m and 1.2 m resolution show few scale difference  
462 concerning building, road, pavement, green space in this urban scene, reconstructing  
463 1.2 m image to 0.3 m would not reflect the performance of the SPM algorithm.  
464 Meanwhile, the scale effect becomes significant when it resampled to 4.8 m resolution,  
465 with pixelated boundaries between objects, which can be used to validate the  
466 reconstruction ability of SPM.

467 The FLAS dataset (He et al., 2022b) consists of pairs of 10 m Sentinel-2 images  
468 and the corresponding 1 m annotation images annotated on Google Earth image, located  
469 in Guangdong, China, in summer 2019. The FLAS dataset was re-cropped to 12785  
470 image patch pairs, in which the Sentinel-2 image has a size of 64×64 pixels and seven  
471 bands, and the annotation image has a size of 320×320 pixels. Training, validation and  
472 test sets were split randomly at the ratio of 7:1:2 for the above three datasets  
473 (Considering that the self-attention based model requires large amount of training  
474 samples, we also examined the effect of various ratio on performance, such as 6:1:3,  
475 the results of which can be referred to supplementary file).

476 Noteworthy is that although the reconstruction scale for the synthetic dataset  
477 (WHDL and LoveDA) is set to 4 in this study, it was sufficiently discussed with an

478 ablation analysis in the supplementary file, which is referred to (Wu et al., 2018).

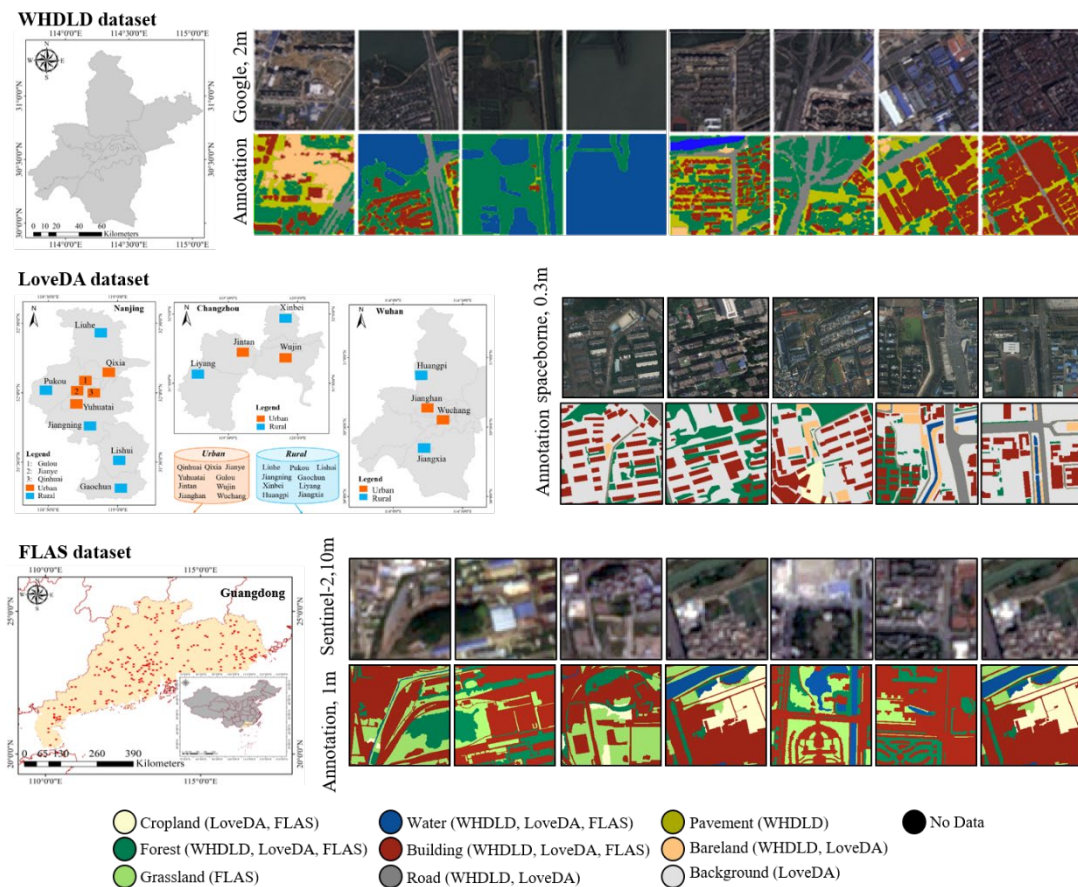
479

480 **Table 1**

481 Public datasets for validation

Dataset	CR image	FR label	Scale	CR size	FR size	Number	Purpose
WHDL D	Downsampled Google image, 8.0 m/pixel	Annotated Google image, 2.0 m/pixel	4	64×64	256×256	4940	Validation on synthetic dataset
LoveDA	Downsampled spaceborne image, 4.8 m /pixel	Annotated Spaceborne image, 1.2 m/pixel	4	64×64	256×256	5987	Validation on synthetic dataset
FLAS	Sentinel-2 image, 10 m /pixel	Annotated Google image, 2.0 m/pixel	5	64×64	320×320	3461	Validation on real dataset

482



**Fig. 8.** Publicly available urban scenario dataset material for this study. (The dataset figures are cited from Shao et al., 2018, Wang et al., 2021 and He et al., 2022b).

483

484 **3.2 Study area**

485        Instead of validation on the public dataset, we further inspected performance on  
486 the downtown areas of the several provincial metropolises in China (**Fig. 9**). These  
487 study areas range from north to south, from the temperate monsoon climate to the  
488 subtropical monsoon climate, and from a first-tier city to a second-tier city, which have  
489 diverse scenes that can inspect the generalization of the LECOS.

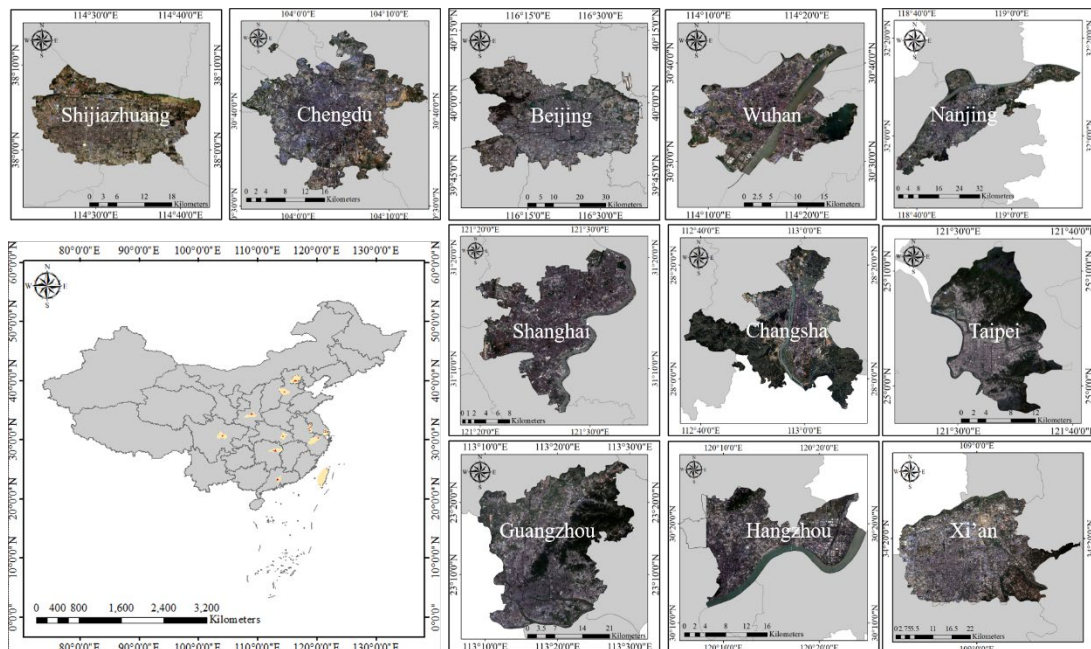
490        Besides, the derived results from the LECOS are compared with the publicly  
491 prevailing land cover products, including the 10-m Finer Resolution Observation and  
492 Monitoring-Global Land Cover (FROMGLC30) product that was produced by  
493 semantic segmentation network and post-processing on Sentinel-2 images (Gong et al.,  
494 2019), the 30 -m GlobeLand30 product that was produced by classical machine learning  
495 and manual post-processing on Landsat images (Chen et al., 2014), the 10 -m Esri Land  
496 Cover product that was produced by semantic segmentation network on Sentinel-2  
497 images (Karra et al. 2021), and the 10-m WorldCover product that was produced by  
498 deep learning method on Sentinel-1 and Sentinel-2 image conducted by European  
499 Space Agency (ESA) (Zanaga et al., 2021).

500        Specifically, we used the Sentinel-2A MSI (Multi Spectral Instrument) of these  
501 study areas, which can be accessed at [https://developers.google.com/earth-](https://developers.google.com/earth-engine/datasets/catalog/COPERNICUS_S2)  
502 [engine/datasets/catalog/COPERNICUS\\_S2](https://developers.google.com/earth-engine/datasets/catalog/COPERNICUS_S2). 10 m Red, Green, Blue, NIR, 20 m SWIR-  
503 1, SWIR-2, and 60 m Coastal bands were selected, which were uniformly resampled to  
504 a 10 m spatial resolution. Cloud area was initially subtracted by a QA60 bitmask band  
505 with cloud mask information, and all the data in the growing period of April-September  
506 2020 were composited by the median-mosaic operation to obtain complete coverage.

507 The pre-trained model trained on the FLAS dataset above was used for 2 m resolution  
508 mapping of these study area.

509 Presently, Sentinel-2 supplies the best available time-series images that have  
510 moderate spatial resolution and a large coverage. However, these images are  
511 insufficient to distinguish the intricate spatial heterogeneity and dynamics of urban  
512 patterns. Therefore, this research attempts to reconstruct the Sentinel-2 images to very  
513 fine spatial resolution to reveal finer urban spatial details.

514



**Fig. 9.** Study area of the selected cities in China.

515

## 516 **4. Results**

### 517 *4.1 Experiment setting*

#### 518 *4.1.1 Implementation details*

519 Firstly, we validated the effectiveness of our proposed method using public dataset  
520 material (Section 4.2.1). Considering the advantage of the explicit correlation modeling  
521 design of our method, we visualized the underlying attention matrix to inspect the

522 learned syntax correlation patterns among the sub-pixels and the context correlation  
523 patterns among the mixed-pixels (Section 4.2.2 and 4.2.3). Besides, we further derived  
524 a statistic based on the attention score of each correlation to quantify the connection  
525 strength between different urban compositions, which can reflect the characteristic,  
526 accessibility and rationality of the urban spatial patterns (Section 4.2.4).

527 The LECOS model was then applied to the large-scale study areas to examine its  
528 practicability (Section 4.3). Specifically, the LECOS model that was well trained on  
529 the FLAS datasets was adopted to reconstruct the 10 m resolution Sentinel-2 images of  
530 the five metropolises, and generate the 2 m resolution fine-grained urban maps. Besides,  
531 based on the quantification strategy of the correlation patterns in Section 4.2.4, we also  
532 evaluated the urban spatial pattern of the five cities according to the correlation statistic,  
533 which was tested by spatial overlap analysis with the 100 m local climate zone map  
534 (Section 4.3).

535 For the parameter settings, the minibatch size was set to 32, and the Adam  
536 optimizer was used with the beta-1 parameter set to 0.9 and the weight decay set to  
537 0.0005, which is recommended in (Liang et al., 2021). The learning rate was initially  
538 set to 0.001 and reduced on plateau with the reduction factor set to 0.8 and the patience  
539 set to 10, which is also suggested for training self-attention network (Liang et al., 2021).  
540 The training epoch was set to 200, the validation performed at each epoch, and we chose  
541 the best model from the validation set for evaluation using the test set. All the  
542 experiments were implemented in PyTorch 1.9.0 and Python 3.6.13 on a server with  
543 three NVIDIA GeForce RTX3090 graphic processing unit (GPU) accelerators (with 24-

544 GB GPU memory).

545

#### 546 *4.1.2 Comparison algorithm*

547       The classical model-driven SPM algorithms like Spatial Attraction based Sub-  
548 pixel Mapping (SASM) (Mertens et al., 2006), the Adaptive MAP model and a winner-  
549 takes-all Class Determination strategy for Sub-pixel Mapping (AMCDSM) (Zhong et  
550 al. 2015), and the state-of-the-art data-driven SPM algorithms like SPMCNN-ESPCN  
551 (He et al., 2021a) are used as benchmark methods for comparison. A comparable  
552 experimental strategy is adopted to compare the model-driven and data-driven methods  
553 (He et al., 2021a). Further, a pixel-level classification methods are also provided for  
554 comparison, including the traditional support vector machine (SVM), and state-of-the-  
555 art data-driven semantic segmentation networks including UNet++ (Ronneberger et al.,  
556 2015), DeepLabV3+ (Chen et al., 2018), Swin Transformer (Liu et al., 2021). Pixel-  
557 level classification results were then downscaled to the target resolution to make them  
558 comparable to the SPM result.

559

#### 560 *4.1.3 Evaluation metric*

561       For quantitative assessment of the performance of the compared methods, the  
562 overall accuracy (*OA*), *Kappa* coefficient, the mean value of the intersection-over-union  
563 (*mIoU*) and the F1-score (*mF1-score*) for each class were used as metrics. In addition,  
564 the producer’s accuracy (*PA*) was selected to examine the class-wise accuracy.

565       The specific evaluation procedure was as follows: Firstly, our method was

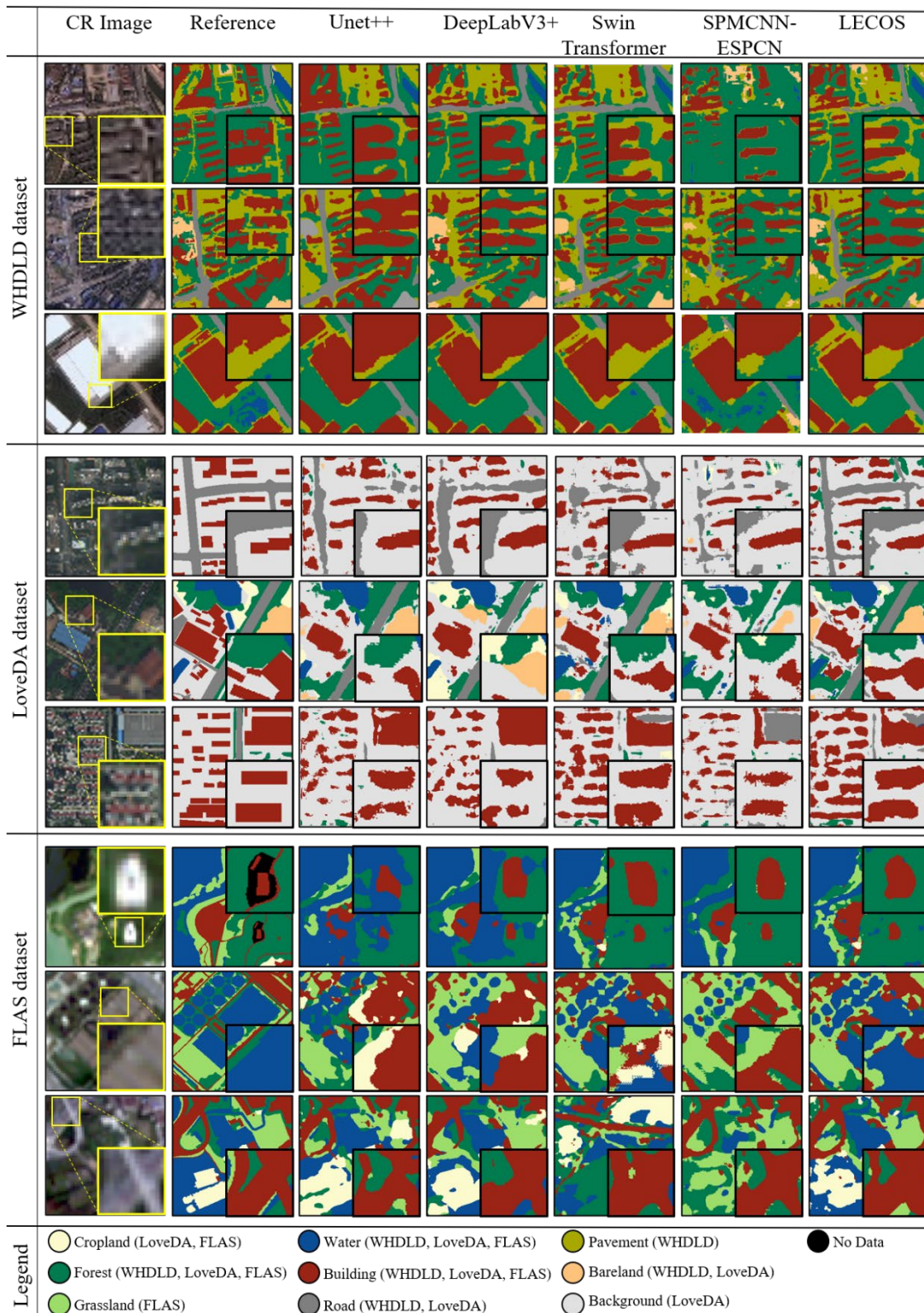
566 qualitatively and quantitatively evaluated by the test set of each public dataset, and the  
567 benchmark SPM algorithms were also applied for comparison (Section 4.2). Then, we  
568 deployed the LECOS model on several metropolises in China to inspect its performance  
569 in revealing very fine spatial resolution urban distributions across large-areas, and the  
570 state-of-the-art land cover products, i.e., 10 m WorldCover, 10 m FROM-GLC10, 10 m  
571 Esri Land Cover and 30 m GlobeLand30 were used for comparison (Section 4.3), to  
572 inspect that whether the 2 m resolution results from LECOS can significantly improve  
573 details of urban pattern compared to the existing products. Besides, we also conducted  
574 an ablation analysis to evaluate the contribution of SBD, CBD, the network structure  
575 (including the number of layers, the 4-stage structure), and the number of training  
576 sample patches, by removing the corresponding design and inspecting the performance  
577 change (see in the supplementary file).

578

## 579 *4.2 Results on public datasets*

### 580 *4.2.1 Visual and quantitative examination*

581 Three patches in the test set were selected from each public dataset for visual  
582 examination (**Fig. 10**), and quantitative assessment of the comparison algorithms is  
583 provided in **Tables 2-4**.



**Fig. 10.** Visual examination of the performance of each algorithm on three public datasets.

584 For visual inspection relating to the WHDLD dataset, LECOS can reconstruct the  
585 detailed compositions and smooth edges well for all the land cover classes in the urban  
586 scene, and the outlines of the buildings are closer to the ground reference. Besides, with



587 the aid of the learned building-pavement correlation rule, LECOS can distinguish  
588 pavements nearby buildings in the residential district, while the other convolutional-  
589 based models that utilize only local information confuse the pavement with roads, and  
590 the traditional self-attention based models cannot well reconstruct the spatial  
591 morphological shape of the objects, as can be seen in **Fig. 10**.

592 For the LoveDA dataset, due to the lack of consideration of the correlation between  
593 road and green parcels, and road and road, that essentially exists in urban scenes, the  
594 roads are broken and the green parcels along the roads are almost lost in the results of  
595 the comparison method. In contrast, LECOS can reconstruct the continuous roads and  
596 green belts better because the context-dependency between urban components is  
597 explicitly learned in the network. Besides, LECOS can better recover the morphological  
598 shapes and edge contours of the buildings, because the syntax-dependency between  
599 sub-pixels is also modelled explicitly.

600 For the FLAS dataset, LECOS is able to precisely identify the pond, grass and  
601 buildings, because the context dependency between pond and grass, and grass and  
602 building are learned. However, the convolutional-based method is prone to misclassify  
603 the pond. Noteworthy is that despite the erroneous identification, the traditional SPM  
604 method (i.e., SPMCNN-ESPCN) is able to reconstruct a more complete morphological  
605 shape of the pond than the other CNN-based method, demonstrating the effectiveness  
606 of sub-pixel-scale manipulation. It is also reflected in our method that the shape of the  
607 pond is more complete, and the impervious road is more continuous, benefiting from  
608 the learned syntax-dependence.

609

610 **Table 2**

611 Quantitative assessment on WHDL D dataset

	Model-driven			Data-driven				Model- and Data-driven
	SVM	SASM	AMCDSM	UNet++	DeepLabV3+	Swin Transformer	SPMCNN-ESPCN	LECOS
Building	0.5144	0.1941	0.1906	0.6208	0.6492	0.7145	0.5857	<b>0.7372</b>
Road	0.0000	0.1846	0.1843	0.5815	<b>0.6596</b>	0.5731	0.1985	0.6079
Pavement	0.4760	0.2508	0.2513	0.5729	<b>0.5888</b>	0.5381	0.4824	0.5874
Vegetation	0.8704	0.6879	0.7007	<b>0.9088</b>	0.8766	0.8825	0.9042	0.8987
Bare Soil	0.2544	0.4255	0.4269	0.3652	0.2762	<b>0.5203</b>	0.3499	0.4879
Water	0.4612	0.0907	<u>0.0874</u>	<u>0.9456</u>	<b>0.9654</b>	0.9494	0.8589	0.9562
<i>OA</i> (%)	61.69	39.08	39.50	81.61	81.90	81.77	76.02	<b>82.66</b>
<i>Kappa</i>	0.4280	0.1205	0.1239	0.7375	0.7355	0.7446	0.4691	<b>0.7562</b>
<i>mIoU</i>	0.3127	0.1445	0.1529	0.5609	0.5549	0.5664	0.4554	<b>0.5842</b>

612

613 **Table 3**

614 Quantitative assessment on LoveDA dataset

	Model-driven			Data-driven				Model- and Data-driven
	SVM	SASM	AMCDSM	UNet++	DeepLabV3+	Swin Transformer	SPMCNN-ESPCN	LECOS
Background	0.7210	0.0504	0.0503	0.6947	0.7361	0.7144	<b>0.7731</b>	0.7262
Building	0.3965	0.0112	0.0112	0.5986	0.5661	<b>0.6056</b>	0.3845	0.6043
Road	0.0042	0.0159	0.0157	0.4268	0.4301	0.4903	0.2305	<b>0.5180</b>
Water	0.0454	0.2484	0.2518	0.7343	0.7006	0.7419	0.5788	<b>0.7549</b>
Barren	0.0987	0.1787	0.1767	0.5105	0.3329	0.3412	0.1792	<b>0.5143</b>
Forest	0.4359	0.7315	0.7448	0.4771	0.4544	0.5305	0.4199	0.5844
Agriculture	0.1793	0.3717	0.3730	0.5939	0.5813	0.5503	0.3107	<b>0.6817</b>
<i>OA</i> (%)	44.14	15.44	15.57	61.75	61.26	62.33	53.13	<b>65.90</b>
<i>Kappa</i>	0.1741	0.0527	0.0536	0.4789	0.4601	0.4842	0.3284	<b>0.5377</b>
<i>mIoU</i>	0.1706	0.0317	0.0323	0.4355	0.4212	0.4301	0.2877	<b>0.4882</b>

615

616 **Table 4**

617 Quantitative assessment on FLAS dataset

	Model-driven			Data-driven				Model- and Data-driven
	SVM	SASM	AMCDSM	UNet++	DeepLabV3+	Swin Transformer	SPMCNN-ESPCN	LECOS
Cropland	0.0000	0.5558	0.6069	0.6394	0.6762	0.6025	0.0000	<b>0.7050</b>
Tree	0.7042	0.6593	0.7139	0.7397	0.7093	0.6349	<b>0.7475</b>	0.6739
Grass	0.4188	0.1297	0.1308	0.4472	0.4737	0.4888	<b>0.6051</b>	0.4123
Water	0.8336	0.8022	0.8270	0.7127	0.7376	<b>0.8542</b>	0.6150	0.7376
Impervious	0.8953	0.7369	0.7876	0.8646	0.8617	0.8791	0.7933	<b>0.9098</b>
<i>OA</i> (%)	74.11	62.49	66.88	77.02	77.18	77.25	69.40	<b>78.80</b>
<i>Kappa</i>	0.5852	0.3071	0.3479	0.6485	0.6527	0.6462	0.5479	<b>0.6687</b>
<i>mIoU</i>	0.4432	0.2251	0.1791	0.5373	0.5411	0.5535	0.3981	<b>0.5495</b>

618

619 In terms of quantitative assessment (**Table 2-4**), LECOS achieves the best *OA* and620 *Kappa* with 82.66% and 0.7562 on the WHDL D dataset, which surpasses the

621 benchmark SPM method (SPMCNN-ESPCN) by 6.64% and 0.2851, and outperforms

622 the benchmark semantic segmentation method (DeepLabv3+) by 0.76% and 0.0187.

623 Besides, LECOS also performs well in terms of  $PA$  accuracy on each urban component  
624 class that can be hard to identify and reconstruct, such as building, road and barren. The  
625 accuracies of these are increased by 0.0057, 0.0879 and 0.0038, respectively, compared  
626 to the benchmark on the LoveDA dataset. For the FLAS dataset, our method again  
627 achieves the highest accuracy, with  $OA$  and  $Kappa$  accuracy of 78.80% and 0.6687, and  
628 outperforms the benchmark by 1.62% and 0.0160, respectively. Noteworthy is that the  
629 model-driven methods have a low accuracy on the WHDL D and LoveDA dataset. The  
630 reason is that the RGB image has insufficient spectral information to derive robust  
631 endmember spectrum and abundance images, and the large uncertainty of the  
632 abundance images impacts the subsequent SPM process.

633

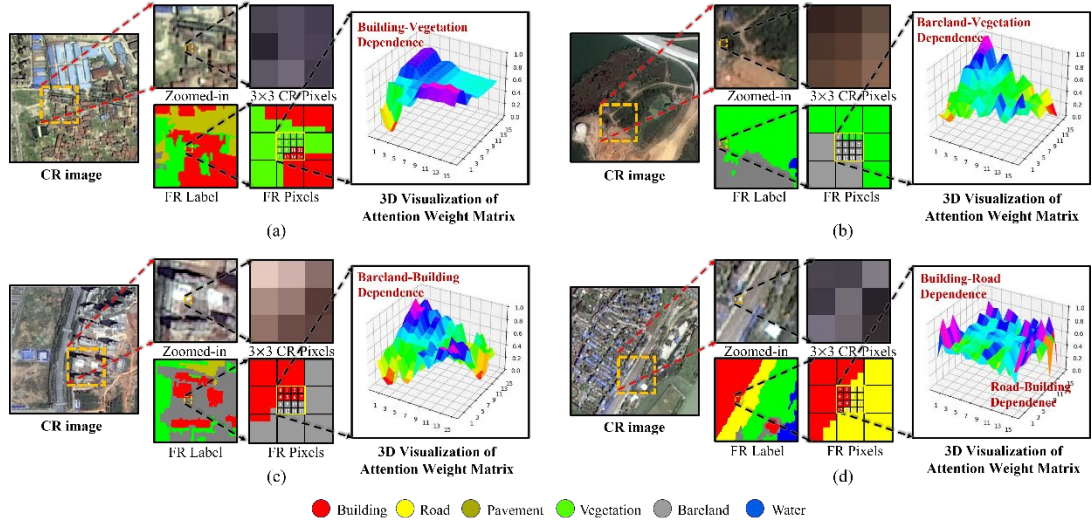
#### 634 *4.2.2 Explicability of the sub-pixel-level syntax-dependence*

635 A reliable syntax that indicates the sub-pixel location of each urban composition  
636 is the key to a successful detail reconstruction. Thus, we visualize and inspect the  
637 learned syntax in this section. Specifically, we selected four representative scenes in the  
638 WHDL D dataset (correlation patterns of different classes are similar in a given dataset),  
639 and focused the evaluation on the zoomed-in area in the CR image within  $3 \times 3$  mixed-  
640 pixels. Our method was used to reconstruct the center mixed-pixel to 16 sub-pixels.  
641 Then, the intermediate variable (i.e., the syntax-dependence represented by the  
642 attention score matrix  $SA \in R^{S^2 \times S^2}$  in stage-1 layer 1, which is defined in Eq. (4)) was  
643 visualized in 3-dimensional space (**Fig. 11**). The size of the matrix  $SA$  is  $16 \times 16$ , and  
644 the z-dimension represent the correlation strength between the two sub-pixels.

645 Taking **Fig. 11a** as an example, a high attention score occurs between building  
646 sub-pixel and vegetation sub-pixel (as the blue and purple area in the 3-D visualization),  
647 which means vegetation and building have a close correlation and often appear together  
648 in the urban scene, and their distributions are mutually attracted when predicting the  
649 sub-pixel locations. Therefore, with the guide of the syntax-dependence indication, the  
650 detailed edge and corners of the building can be well restored.

651 Taking **Fig. 11b** as an example, with rapid urban land development and utilization,  
652 most of the bareland in the urban scene is usually derived from deforestation. Therefore,  
653 bareland is often accompanied by vegetation, and the bareland sub-pixel exhibits a  
654 strong attention score to vegetation sub-pixel. However, the attention score is not  
655 symmetrical (i.e., the vegetation seldom pays attention to the bareland), since the  
656 vegetation (e.g., greenbelt, horticulture and arboriculture) is not necessarily  
657 accompanied by bareland.

658 Similarly, **Fig. 11c** demonstrates a construction site, where the bareland sub-pixel  
659 pays more attention to the building sub-pixel, and in **Fig. 11d** with residential areas near  
660 the main street, the building sub-pixel and the road sub-pixel mutually attract each other  
661 with high attention score, which meets the needs for convenient settlements.



**Fig. 11.** Visualization of the learned syntax-dependence.

662

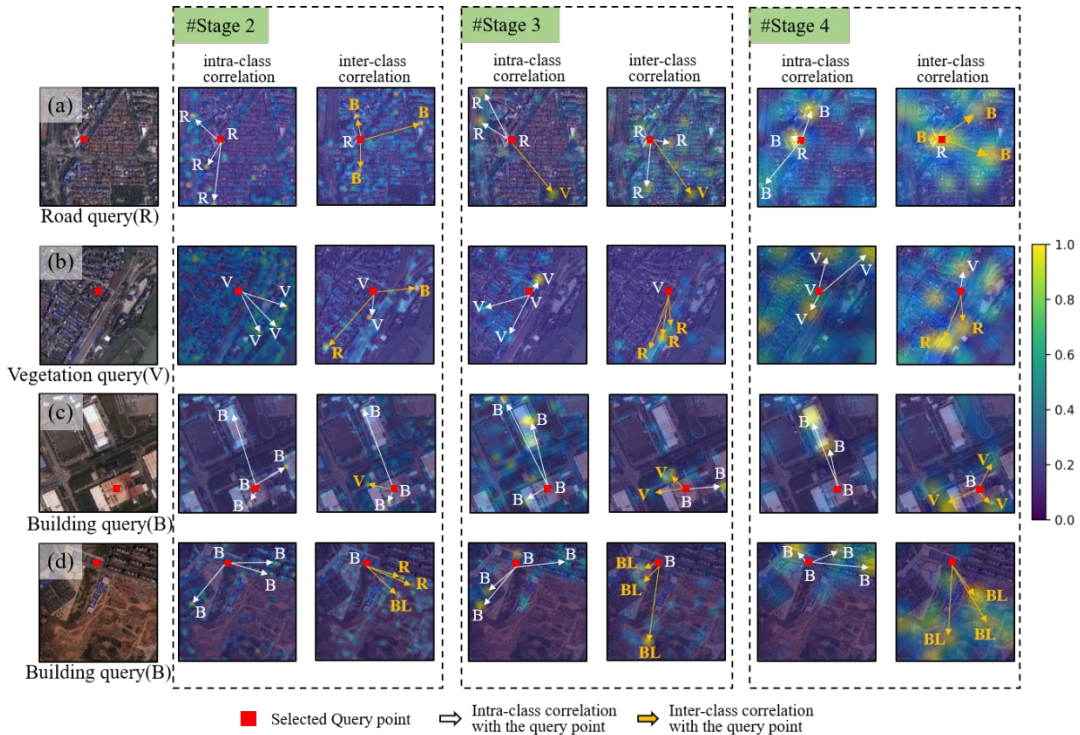
#### 663 4.2.3 Explicability of the pixel-level context dependence

664 Four scenes in the WHDL D dataset are demonstrated for illustration. Taking **Fig.**  
665 **12a** as an example, we selected the road pixel  $i$  (denoted as red point) as the *Query*, and  
666 used the *Key* of the residual pixels to generate the context attention score matrix  
667  $\{\alpha_{i,j}\}_{j=1\dots N}$  (defined in Eq. (9)) across different stages throughout the LECOS network,  
668 which represents the correlation strength of each pixel to the selected road pixel. In  
669 each stage, we randomly selected two attention matrices  $\{\alpha_{i,j}\}_{j=1\dots N}$  at different layers  
670 and different heads, one is representative for inter-class correlation and one is for intra-  
671 class correlation. Noteworthy is that stage-1 is not visualized, since the context  
672 information is not obvious in the shallow layer. Generally, we highlight three highest  
673 attention scores in  $\{\alpha_{i,j}\}_{j=1\dots N}$  with directional arrows in each attention score map, to  
674 demonstrate the correlation patterns between the red point and the residual pixels. For  
675 discrimination, the heat points and arrows with white color notations (i.e., B for  
676 building, R for road, V for vegetation, BL for bareland) indicate the intra-class

677 correlation with the red point, while the heat points with orange notation represent the  
 678 intra-class correlation.

679 Two phenomena are obvious: 1) we find that the attention tends to focus on the  
 680 local small land covers in stage-2 or 3 (e.g., road usually pays more attention to the  
 681 elongated road and scattered vegetation in **Fig. 12a**), while high attentions are spatially  
 682 aggregated in stage-4 and tend to focus on the dominant land covers in the scene (e.g.,  
 683 the road pays more attention to buildings), indicating that the context gradually  
 684 converges towards global perception; 2) Taking stage-4 as an example, the roads are  
 685 highly correlated with buildings (**Fig. 12a**), the vegetation pays more attention to road  
 686 and itself in residential areas (**Fig. 12b**), the buildings have a large correlation with the  
 687 surrounding vegetation (**Fig. 12c**), and the buildings are highly correlated with  
 688 themselves and bareland in construction sites (**Fig. 12d**).

689



**Fig. 12.** Visualization of the context dependence across different stage.

690

#### 691 *4.2.4 Quantitative analysis of the connection strength*

692 We construct a statistic to quantify the correlation strength for each correlation  
693 pattern in the above analysis. Specifically, the output 1-dimensional sequential feature  
694  $B \in R^{(4 \times 4 \times C) \times (W/4 \times H/4)}$  in stage-1 is reshaped to a 2-dimensional feature with size  
695  $H/4 \times W/4 \times 16C$ , and the FR labels are resampled from  $H \times W$  to  $H/4 \times W/4$  to  
696 match the size of the reshaped feature  $B$ , which is used to provide category information  
697 for each pixel of  $B$ . Subsequently, the context-dependence represented by the attention  
698 score matrix  $A \in R^{(W/4 \times H/4) \times (W/4 \times H/4)}$  in stage-1, which indicates the correlation  
699 between each of the two pixels in  $B$ , is used for the attention score statistic. For example,  
700 when calculating the building-road correlation pattern (i.e., B-R), the attention score  
701 between all the building pixels and the road pixels in the matrix  $A$ , indicated by the FR  
702 labels, is summarized. In this way, the correlation patterns can be quantified.

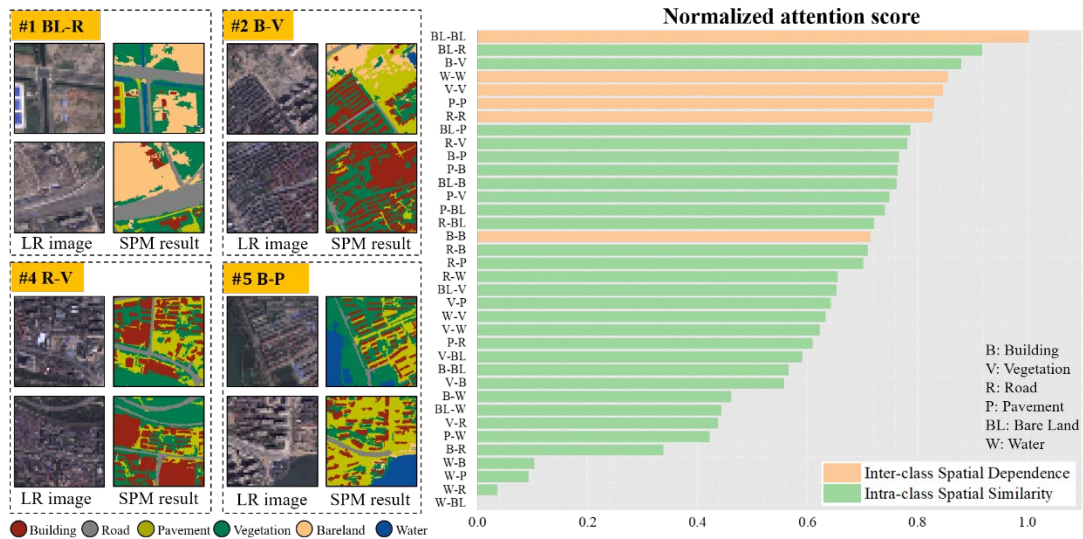
703 We construct a statistic for the WHDL D dataset. WHDL D has six classes such that  
704 there are, in total, 36 correlation patterns between each of the two categories, which are  
705 ranked according to total attention score (**Fig. 13**), and which reflect the most  
706 contributing and important correlation patterns in the WHDL D urban scenario.

707 From the ranking list, we find that the intra-class correlation is usually stronger  
708 than the inter-class correlation, which is reasonable due to spatial autocorrelation theory  
709 where the attraction within the same class is expected to be greater than between  
710 different classes. As for inter-class correlations, we visualize the top four patterns in  
711 **Fig. 13** (i.e., BL-R, B-V, R-V, B-P), for examining the reliability of the quantification.

712 We find that: 1) for the BL-R correlation pattern, the bareland in the urban scene is  
713 usually represented as a construction site. Thus, a frequent transportation of the  
714 construction materials is highly necessary for the construction project, which can  
715 explain the frequent occurrence of the BL-R pattern; 2) The high attention score of the  
716 B-V and R-V correlation pattern demonstrates that the building is highly correlated with  
717 vegetation, and the road is also tightly correlated with vegetation, indicating the  
718 presence of a considerable number of residential scenarios in the WHDL dataset, and  
719 also indicating that the residential horticulture are well organized; 3) The B-P  
720 correlation pattern also ranks highly, which means that buildings are also highly  
721 correlated with pavements, which meets the convenient transportation requirements of  
722 people living in the residential building.

723 From the above qualitative and quantitative analysis of the syntax- and context-  
724 dependence, we find that they are not only explicable, but can also be used as an  
725 indicator to understand the complex urban spatial patterns (e.g., estimating the human  
726 exposure to greenspace in urban scenes by the B-V and R-V attention score, and  
727 estimating the transportation convenience by the B-R and B-P attention score, etc.).





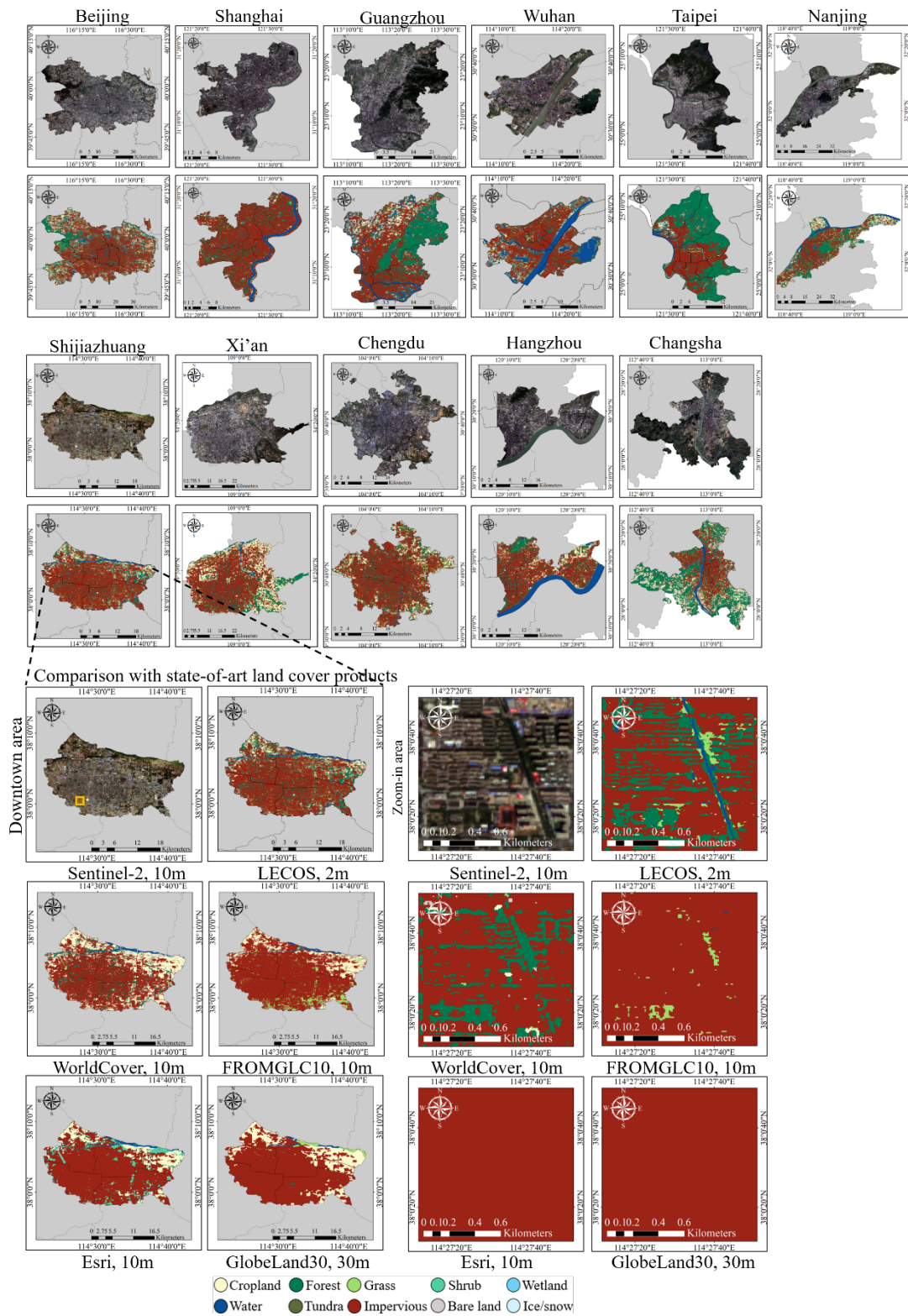
**Fig. 13.** Visualization of the normalized attention score for demonstrating the strength of each correlation pattern.

728

### 729 4.3 Results on five metropolises

#### 730 4.3.1 Visual and quantitative examination

731 The 2 m spatial resolution land cover maps in the downtown area of the five cities  
 732 are displayed visually in **Fig. 14** (taking Shijiazhuang as an example), as well as the  
 733 comparison with the public land cover products. A zoomed-in area was selected to  
 734 inspect the performance in terms of detail reconstruction. Results for other cities are  
 735 given in Appendix II in the supplementary file.



**Fig. 14.** Comparison of our method with state-of-the-art public land cover products.

736

737

It is obvious that the detail distribution of the urban compositions such as buildings,

738

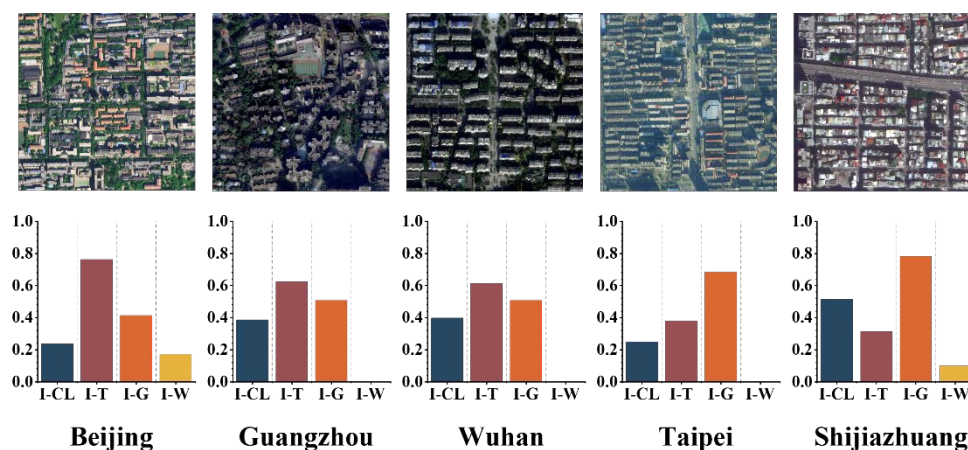
urban greenspace and grassland is usually missing in the Esri Land Cover and

739 GlobeLand30 products. Although detailed compositions can be reconstructed by  
 740 WorldCover and FROMGLC10, the morphological shape is vague and incomplete with  
 741 jagged boundaries. In contrast, our method can restore a fine-grained structure into the  
 742 urban compositions, especially for small-sized or elongated urban greenspace that is  
 743 usually mixed with impervious surfaces in the original Sentinel-2 images and difficult  
 744 to detect.

745

#### 746 4.3.2 Connection strength of each correlation pattern

747 The normalized attention score statistic of each correlation pattern in the several  
 748 cities is provided in **Fig. 15**. Specifically, impervious surface was selected as the *Query*,  
 749 and the other classes were used as the *Key*. Thus, four correlation patterns (i.e., I-CL,  
 750 I-T, I-G, I-W; I for impervious, CL for cropland, T for Tree, G for Grass and W for  
 751 water) were inspected.



**Fig. 15.** The attention score of each correlation patterns in the five cities.

752

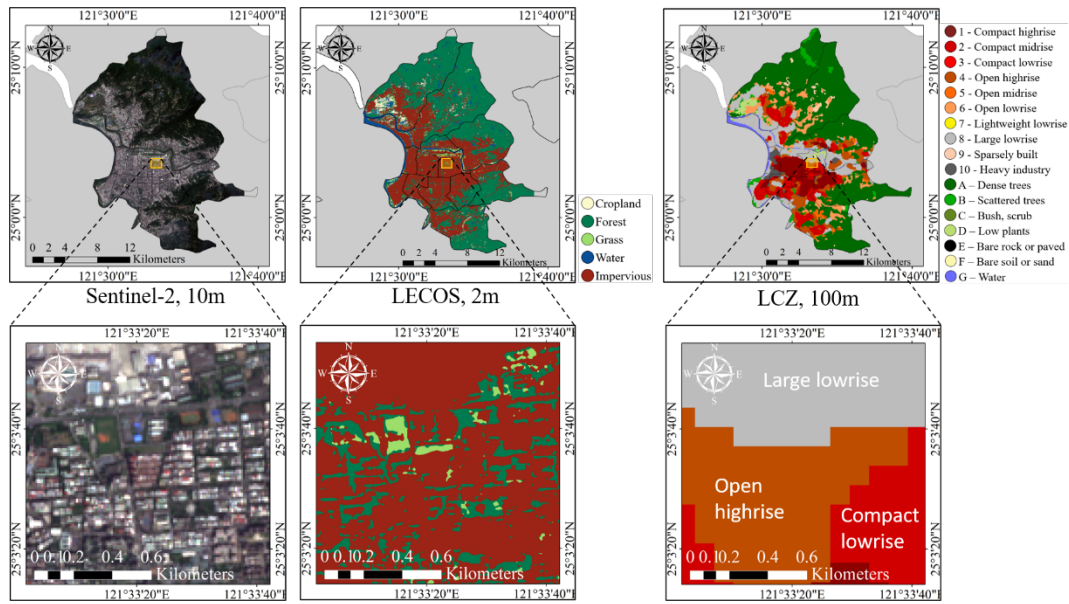
753 Focusing on the I-T score representing human exposure to greenspace, it is  
 754 obvious that Beijing achieves the highest score, which means that the urban greenspace

755 is cross-interweaved with impervious surfaces with highly accessibility. Meanwhile,  
756 although Taipei has a larger greenspace area ratio, there is less greenspace interspersed  
757 within the impervious surface buildings or roads, resulting in a weaker connection  
758 strength and lower exposure to urban greenspace. In this way, we find that the attention  
759 score statistic can precisely reflect the urban patterns.

760

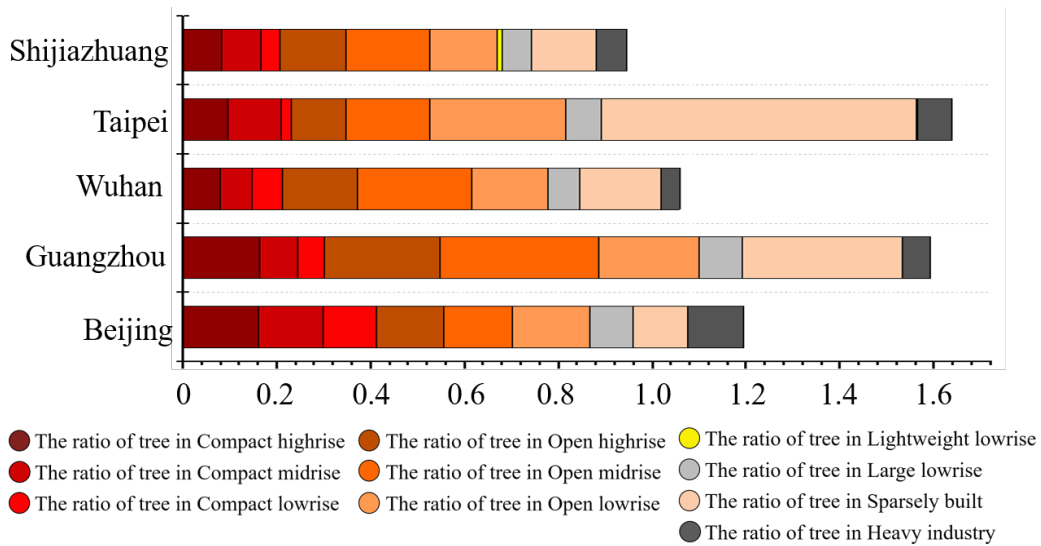
#### 761 *4.3.3 Overlay with Local Climate Zone (LCZ) dataset*

762 From the above analysis, Beijing achieves the highest I-T score, which means that  
763 more trees are planted within the buildings group. To validate this conclusion and  
764 further inspect the effectiveness of the I-T attention score statistic in terms of spatial  
765 accessibility, we further consider a spatial overlay analysis between our result and the  
766 Local Climate Zone (LCZ) dataset. The LCZ dataset divides impervious surfaces into  
767 10 finer classes, in terms of the density (compact, open, sparse), height (high rise, mid  
768 rise, low rise), and material (heavy and lightweight), as shown in **Fig. 16**. We measure  
769 the area of the urban trees that fall into the area of each impervious class, which is  
770 normalized by dividing the total area of that impervious class. A higher tree ratio in  
771 denser classes implies greater spatial accessibility and greenspace exposure.  
772 Specifically, the 100 m spatial resolution GLOBAL LCZ MAP (Demuzere et al., 2022)  
773 was used. A zoomed-in area in Taipei is provided in **Fig. 16** to spatially demonstrate the  
774 overlay of the two products, and the urban tree ratio statistic for each impervious class  
775 in the five cities is provided in **Fig. 17**.



**Fig. 16.** The ratio of trees distributed in each impervious-related category region.

776



**Fig. 17.** The ratio of trees distributed in each impervious-related category region.

777

778 Beijing has the highest greenspace ratio amongst the three compact impervious  
 779 classes, which is consistent with the I-T score analysis. Specifically, the proportion of  
 780 greenspace in compact areas is significantly higher than that in other cities, indicating  
 781 that green exposure and accessibility are greater, which can serve more people and, for  
 782 example, help to mitigate the urban heat island effect. On the other hand, in Taipei city,  
 783 a large greenspace ratio is configured in the sparse built area, while the proportion of

784 greenspace in the compact area is small, indicating that the configuration of green  
785 vegetation may be sub-optimal for some residents, and resulting in lower greenspace  
786 accessibility.

787

## 788 **5. Discussion**

789       The experimental results demonstrated the superiority of our method, which is due  
790 to integrating a data-driven learning procedure with model-driven spatial correlation  
791 characterization, which can achieve more fine-grained land cover identification than  
792 using either of them alone. Other findings are 1) from a comparison between pixel-level  
793 classification and the SPM method, we find that SPM can reconstruct more complete  
794 morphological shapes and details owing to its sub-pixel manipulation ability; while the  
795 pixel-level classification method commonly achieves higher pixel-level accuracies due  
796 to its strong class identification ability. Our method combines the advantages of both  
797 (i.e., context for class identification and syntax for sub-pixel location inference) and,  
798 thus, achieves greater identification accuracy and detail reconstruction performance; 2)  
799 from the visualization of the correlation, the large spatial correlation between two  
800 objects usually has long-distance due to the heterogeneity of urban scenes, which  
801 demonstrates the necessity of global “mutual retrieval”. Besides, datasets with different  
802 urban scenes reveal different intensities of the correlation patterns, which demonstrate  
803 the explicability of our method; 3) from the ablation experiment, the self-attention in  
804 self-attention operation (SNS) can learn more spatial detail features than the typical  
805 self-attention.

806 From the experimental results in the selected cities, we found that 1) the 2 m spatial  
807 resolution of our products can effectively reveal a heterogeneous distribution of the  
808 urban compositions, compared to the prevailing 10 m or 30 m resolution land cover  
809 products, which are more suitable for urban pattern recognition; 2) the traditional green  
810 coverage indices cannot reflect human accessibility to greenspace (Chen et al., 2022).  
811 For example, Taipei has larger ratio of greenspace than other cities, however, it is  
812 distributed in the mountainous area and far away from the residential area, leading to a  
813 relatively small correlation between impervious and tree reflected by our learned  
814 correlation; 3) an overlay analysis with the local climate zone product also confirms  
815 that more greenspace ratio is distributed in compact impervious areas in Beijing than in  
816 Taipei, which might be expected to serve more people; 4) the correlation derived from  
817 our method is able to consider the global pattern of each composition, which can reflect  
818 the cross-interweaved degree of both. Such inferences may find future application in  
819 social and environmental analyses of cities, for example, for understanding urban  
820 fairness and measuring resource distribution.

821 There exist several limitations of the proposed method. First, the fixed size of  
822 tokenization may adversely affect the learning of the spatial correlation, since the  
823 distance between two correlated urban components may vary. Second, the design of  
824 interpolating the CR image at the beginning of the network may increase the  
825 computational burden, since the size of the feature maps may be increased throughout  
826 the network. Third, we make a statistic of the computational efficiency of each  
827 algorithm in terms of parameter number (Para) and floating point of per second (FLOPs)

828 (refer to the supplementary file), and find that although the number of parameters in  
829 LECOS model is relatively larger than that in other networks, the computational  
830 complexity is at the same order of magnitude. Fourth, when evaluating the correlation  
831 score of each city, it can only reflect a relative comparison that, e.g., impervious is more  
832 correlated with tree than with cropland, or the correlation score of one city is higher  
833 than another city, however, it is difficult to give a real correlation value for each city to  
834 test our generated correlation score due to the lack of full city ground truth.

835

## 836 **6. Conclusion**

837 Very fine spatial resolution observation of the urban pattern dynamic is the  
838 foundation to reveal urban spatial heterogeneity in cities, and inform judgements and  
839 planning with respect to issues such as urban fairness and access to greenspace.  
840 However, the current large-area land use/land cover products are usually limited by a  
841 moderate spatial resolution (i.e., 10 m or 30 m), which makes them unsuitable for  
842 mapping intricate urban components.

843 The contribution of this research has two folds including algorithm innovation and  
844 scientific findings of urban spatial pattern. In aspect of algorithm, the spatial resolution  
845 of land use/land cover products were increased, specifically for urban scenarios, by  
846 developing a new sub-pixel mapping algorithm, which addresses two challenges in the  
847 sub-pixel mapping community: 1) local spatial autocorrelation is usually limited to a  
848 fixed window, which ignores global contextual correlation that can inform inference at  
849 sub-pixel locations; 2) urban scenarios are complex, which means that the “more



850 proximate, more similar” assumptions of classical spatial autocorrelation are unable to  
851 support inference of a heterogeneously distributed sub-pixel pattern. Therefore, how to  
852 learn the various correlations amongst urban components that can help sub-pixel  
853 location reasoning remains an open question. To this end, we designed an end-to-end  
854 contextual spatial correlation learnable network architecture (LECOS) to address  
855 directly the drawbacks of the fixed autocorrelation, providing two innovations: 1) a  
856 “mutually retrieve” mechanism was designed in an end-to-end network to learn  
857 correlation patterns adaptively; 2) a “self-attention in self-attention” operation was  
858 designed to explicitly infer the classes of the sub-pixels.

859 Validation experiments were conducted on three challenging urban datasets, and  
860 we found that the LECOS was better able to reconstruct the outlines and edges of urban  
861 buildings, roads, trees, etc., with an average *OA* of 82.66%, which significantly  
862 outperformed the benchmark. From the ablation analysis, we found that the designed  
863 “self-attention in self-attention” operation added greatly to the restoration of spatial  
864 details, with an increase in accuracy of 2.14%. From visualization of the learned  
865 correlations in LECOS, we found that datasets with different urban scenes reveal  
866 different types and intensities of correlation patterns, and all of them were consistent  
867 with the *in-situ* reference, demonstrating the explicability of our method.

868 In aspect of scientific findings, several typical metropolises were selected to  
869 examine the practical applicability of the proposed method, and we found that our  
870 method not only revealed heterogeneous urban patterns with very fine spatial resolution  
871 compared to the prevailing 10 m or 30 m resolution land cover products, but also the

872 statistic of the intermediate attention score was able to reflect the human exposure  
873 patterns that are indicative of urban fairness. For example, Beijing was found to have  
874 greater urban greenspace accessibility, relative to Taipei. Such inferences may be useful  
875 in supporting planning decisions.

876 In summary, this research provides the first explicable sub-pixel mapping method  
877 for reconstructing very fine spatial resolution urban spatial patterns. As such, it  
878 represents the first attempt to combine model-driven spatial correlation theory with  
879 data-driven learning practice, which makes it possible to characterize the spatial  
880 heterogeneity of urban spatial patterns. We hope that LECOS will be used in future to  
881 support further sustainable urban development research.

882  
883

#### 884 **Acknowledgements**

885 This work was supported by the National Key R&D Program of China Grant No.  
886 2022YFB3903402, in part by the National Natural Science Foundation of China under  
887 Grants 42222106, 61976234, 42201340, Natural Science Foundation of Guangdong  
888 Province Grant No. 2020A1515110708.

889

#### 890 **Data availability statement**

891 The executable code that supports the findings of this research are available from the  
892 corresponding author upon reasonable request.

893

#### 894 **Reference**

895 Adams, J. B., Sabol, D. E., Kapos, V., Filho, R. A., Roberts, D. A., Smith, M. O., et al.

896 (1995). Classification of multispectral images based on fractions of endmembers:  
897 Application to land-cover change in the Brazilian Amazon. *Remote Sens. Environ.*,  
898 52, 137–154. [https://doi.org/10.1016/0034-4257\(94\)00098-8](https://doi.org/10.1016/0034-4257(94)00098-8)

899 Ardila, J. P., Tolpekin, V. A., Bijker, W., Stein, A., 2011. Markov-random-field-based  
900 super-resolution mapping for identification of urban trees in VHR images. *ISPRS*  
901 *Journal of Photogrammetry and Remote Sensing* 66(6): 762-775. [https://doi.org/](https://doi.org/10.1016/j.isprsjprs.2011.08.002)  
902 [10.1016/j.isprsjprs.2011.08.002](https://doi.org/10.1016/j.isprsjprs.2011.08.002).

903 Atkinson, P. M. 1997. Mapping sub-pixel boundaries from remotely sensed images.  
904 *Innovations in GIS IV*. London, U.K.; Taylor and Francis, 166–180.  
905 [https://www.taylorfrancis.com/chapters/edit/10.1201/9781482272956-](https://www.taylorfrancis.com/chapters/edit/10.1201/9781482272956-25/mapping-sub-pixel-boundaries-remotely-sensed-images-peter-atkinson)  
906 [25/mapping-sub-pixel-boundaries-remotely-sensed-images-peter-atkinson](https://www.taylorfrancis.com/chapters/edit/10.1201/9781482272956-25/mapping-sub-pixel-boundaries-remotely-sensed-images-peter-atkinson).

907 Atkinson, P. M. 2005. Sub-pixel target mapping from soft-classified, remotely sensed  
908 imagery. *Photogramm. Eng. Remote Sensing* 71(7): 839-846. [https://doi.org/10.](https://doi.org/10.14358/PERS.71.7.839)  
909 [14358/PERS.71.7.839](https://doi.org/10.14358/PERS.71.7.839).

910 Chen, B., Wu, S., Song, Y., Webster, C., Xu, B., Gong, P. 2022. Contrasting inequality  
911 in human exposure to greenspace between cities of Global North and Global South.  
912 *Nature commun.*, 13(1), 1-9. <https://doi.org/10.1038/s41467-022-32258-4>.

913 Chen, J., Ban, Y., and Li, S., 2014. China: Open Access to Earth land-cover Map.  
914 *Nature* 514 (7523): 434. <https://doi.org/10.1038/514434c>.

915 Demuzere, M., Kittner, J., Martilli, A., Mills, G., Moede, C., Stewart, I. D., van Vliet,  
916 J., and Bechtel, B., 2022. A global map of Local Climate Zones to support earth  
917 system modelling and urban scale environmental science, *Earth Syst. Sci. Data*

918 14(8) 3835-3873. <https://doi.org/10.5194/essd-14-3835-2022>.

919 Deng, C., Zhu, Z., 2020. Continuous subpixel monitoring of urban impervious surface  
920 using Landsat time series. *Remote Sens. Environ.*, 238, 110929. <https://doi.org/10.1016/j.rse.2018.10.011>.

921

922 Fisher, Y., and Koltun V., 2015. Multi-scale context aggregation by dilated  
923 convolutions. arXiv preprint arXiv:1511.07122.

924 Foley, J. A. et al., 2005. Global consequences of land use. *Science* 309, 570–574.  
925 <https://doi.org/10.1126/science.1111772>.

926 Gong P., et al., 2019. Stable classification with limited sample: transferring a 30-m  
927 resolution sample set collected in 2015 to mapping 10-m resolution global land  
928 cover in 2017. *Sci. Bull.*, 64, 370-373. <https://doi.org/10.1016/j.scib.2019.03.002>

929 He, D., Zhong, Y., Feng, R., Zhang, L., 2016a. Spatial-Temporal Sub-Pixel Mapping  
930 Based on Swarm Intelligence Theory. *Remote Sens.*, 8(11), 30. [https://doi.org](https://doi.org/10.3390/rs8110894)  
931 [/10.3390/rs8110894](https://doi.org/10.3390/rs8110894).

932 He, D., Zhong, Y., Wang, X., Zhang, L., 2021a. Deep convolutional neural network  
933 framework for subpixel mapping. *IEEE Trans. Geosci. Remote Sens.*, 59(11),  
934 9518–9539. <https://doi.org/10.1109/TGRS.2020.3032475>.

935 He, D., Shi, Q., Liu, X., Zhong, Y., Zhang, X., 2021b. Deep Subpixel Mapping Based  
936 on Semantic Information Modulated Network for Urban Land Use Mapping. *IEEE*  
937 *Trans. Geosci. Remote Sens.*, 59(12), 10628-10646. [https://doi.org/10.1109/TGRS.](https://doi.org/10.1109/TGRS.2021.3050824)  
938 [2021.3050824](https://doi.org/10.1109/TGRS.2021.3050824).

939 He, D., Shi, Q., Liu, X., Zhong, Y., He, D., Zhong, Y., Feng, R., Zhang, L., 2022a.

940       Generating 2m fine-scale urban tree cover product over 34 metropolises in China  
941       based on deep context-aware sub-pixel mapping network. *Int. J. Appl. Earth Obser.*  
942       *Geoinf.* 106, 102667. <https://doi.org/10.1016/j.jag.2021.102667>.

943   He, D., Shi, Q., Liu, X., Zhong, Y., Xia, G., Zhang, L., 2022b. Generating annual high  
944       resolution land cover products for 28 metropolises in China based on a deep super-  
945       resolution mapping network using Landsat imagery. *GIScience & Remote Sensing*,  
946       59(1), 2036-2067. <https://doi.org/10.1080/15481603.2022.2142727>.

947   He, K., Zhang, X., Ren, S., Sun, J., 2016b. Deep Residual Learning for Image  
948       Recognition.” *IEEE Conference on Computer Vision and Pattern Recognition*  
949       (CVPR), Las Vegas, NV, USA, 770–778. <https://doi.org/10.1109/CVPR.2016.90>.

950   IPCC Climate Change 2013: The Physical Science Basis (eds Stocker, T. F. et al.)  
951       (Cambridge Univ. Press, 2013).

952   Karra, K., Kontgis, C., Statman-Weil, Z., Mazzariello, J. C., Mathis, M., Brumby S. P.,  
953       2021. Global Land use/land Cover with Sentinel 2 and Deep Learning. *IEEE*  
954       International Geoscience and Remote Sensing Symposium (IGARSS), Brussels,  
955       Belgium, 4704–4707. <https://doi.org/10.1109/IGARSS47720.2021.9553499>

956   Kasetkasem, T., Arora, M. K., Varshney, P. K., 2005. Super-resolution land cover  
957       mapping using a Markov random field based approach. *Remote Sens. Environ.*  
958       96(3-4): 302-314. <https://doi.org/10.1016/j.rse.2005.02.006>.

959   Li, X., Ling, F., Foody, G. M., Ge, Y., Zhang, Y., Du, Y., 2017. Generating a series of  
960       fine spatial and temporal resolution land cover maps by fusing coarse spatial  
961       resolution remotely sensed images and fine spatial resolution land cover maps.

962 Remote Sens. Environ. 196: 293-311. <https://doi.org/10.1016/j.rse.2017.05.011>.

963 Liang, J., Cao, J., Sun, G., Zhang, K., Van Gool, L., & Timofte, R. (2021). SwinIR:  
964 Image restoration using swin transformer. In Proceedings of the IEEE/CVF  
965 international conference on computer vision (pp. 1833-1844). [https://arxiv.org/  
966 abs/2108.10257](https://arxiv.org/abs/2108.10257).

967 Lin, T.-Y., Goyal, P., Girshick, R., He, K., Dollar, P., 2020. Focal loss for dense object  
968 detection. IEEE Trans. Pattern Anal. Mach. Intell. 42 (2), 318–327. [https://doi.org/  
969 10.1109/TPAMI.2018.2858826](https://doi.org/10.1109/TPAMI.2018.2858826).

970 Liu, X., Pei, F., Wen, Y., Li, X., Wang, S., et al., 2019. Global urban expansion offsets  
971 climate-driven increases in terrestrial net primary productivity. Nature  
972 Communication 10: 5558. <https://doi.org/10.1038/s41467-019-13462-1>.

973 Liu, X., Huang, Y., Xu, X., Li, X., Li, X., et al., 2020. High-spatiotemporal-resolution  
974 mapping of global urban change from 1985 to 2015. Nature Sustainability 3: 564–  
975 570. <https://doi.org/10.1038/s41893-020-0521-x>.

976 Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., ... Guo, B., 2021. Swin transformer:  
977 Hierarchical vision transformer using shifted windows. In Proceedings of the  
978 IEEE/CVF international conference on computer vision, 10012-10022.

979 Mertens, K. C., Baets, B. D., Verbeke, L. P. C., Wulf, R. D., 2006. A sub-pixel mapping  
980 algorithm based on sub-pixel/pixel spatial attraction models. Int. J. Remote Sens.  
981 27 (15): 3293-3310. <https://doi.org/10.1080/01431160500497127>.

982 Milletari, F., Navab, N., Ahmadi, S.-A., 2016. V-net: Fully convolutional neural  
983 networks for volumetric medical image segmentation. IEEE International

984 Conference on 3D Vision (3DV), 565-571. <https://doi.org/10.1109/3DV.2016.79>

985 Poudyal, A., 2013. Spatial Statistics and Super Resolution Mapping for Precision  
986 Agriculture Using VHR Satellite Imagery [M]. Netherland: University of Twente.

987 Rashed, T., Weeks, J. R., Roberts, D., Rogan, J., Powell, R., 2003. Measuring the  
988 physical composition of urban morphology using multiple endmember spectral  
989 mixture models. *Photogramm. Eng. Remote Sensing*, 69, 1011–1020.  
990 <https://doi.org/10.14358/PERS.69.9.1011>

991 Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for  
992 biomedical image segmentation. In *International Conference on Medical image  
993 computing and computer-assisted intervention*, Springer, Cham, 234-241.

994 Shao, Z., Yang, K., Zhou, W., 2018. Performance Evaluation of Single-Label and  
995 Multi-Label Remote Sensing Image Retrieval Using a Dense Labeling Dataset.  
996 *Remote Sens.*, 10(6), 964. <https://doi.org/10.3390/rs10060964>

997 Small, C., 2003. High spatial resolution spectral mixture analysis of urban reflectance.  
998 *Remote Sens. Environ.*, 88, 170– 186. <https://doi.org/10.1016/j.rse.2003.04.008>

999 Song, M., Zhong, Y., Ma, A., Feng, R., 2019. Multiobjective Sparse Subpixel Mapping  
1000 for Remote Sensing Imagery. *IEEE Trans. Geosci. Remote Sens.* 57(7), 4490-  
1001 4508. <https://doi.org/10.1109/TGRS.2019.2891354>.

1002 Stokes, E. C., Seto, K., 2018. Characterizing and measuring urban landscapes for  
1003 sustainability. *Environ. Res. Lett.* 14, 045002. <https://doi.org/10.1088/1748-9326/aafab8>

1004

1005 Su, Y., 2019. Integrating a scale-invariant feature of fractal geometry into the Hopfield

1006 neural network for super-resolution mapping. *Int. J. Remote Sens.* 40(23): 8933-  
1007 8954. <https://doi.org/10.1080/01431161.2019.1624865>.

1008 Verdonck, M. L., Okujeni, A., van der Linden, S., Demuzere, M., De Wulf, R., Van  
1009 Coillie, F., 2017. Influence of neighbourhood information on ‘Local Climate  
1010 Zone’ mapping in heterogeneous cities. *Int. J. Appl. Earth. Obs. Geoinf.*, 62, 102-  
1011 113. <https://doi.org/10.1016/j.jag.2017.05.017>

1012 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., et al., 2017. Attention is  
1013 all you need. *Advances in neural information processing systems*, 30.

1014 Wang, Q., Shi, W., Atkinson, P.M., Zhao, Y., 2015. Downscaling modis images with  
1015 area-to-point regression kriging. *Remote Sens. Environ.* 166, 191–204.  
1016 <https://doi.org/10.1016/j.rse.2015.06.003>.

1017 Wang, J., Zheng, Z., Ma, A., Lu, X., Zhong, Y., 2021. A Remote Sensing Land-Cover  
1018 Dataset for Domain Adaptive Semantic Segmentation. *Proceedings of the Neural  
1019 Information Processing Systems Track on Datasets and Benchmarks*. [https://  
1020 doi.org/ 10.48550/arXiv.2110.08733](https://doi.org/10.48550/arXiv.2110.08733)

1021 Wang, X., Ross G., Abhinav G., and He, K., 2018. Non-local neural networks. In  
1022 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,  
1023 7794-7803.

1024 Wu, C., Murray, A. T., 2003. Estimating impervious surface distribution by spectral  
1025 mixture analysis. *Remote Sens. Environ.*, 84, 493–505. [https://doi.org/10.  
1026 1016/S0034-4257\(02\)00136-0](https://doi.org/10.1016/S0034-4257(02)00136-0)

1027 Wu, S., Ren, J., Chen, Z., Jin, W., Liu, X., Li, H., . . . Guo, W., 2018. Influence of



1028 reconstruction scale, spatial resolution and pixel spatial relationships on the sub-  
1029 pixel mapping accuracy of a double-calculated spatial attraction model. Remote  
1030 Sens. Environ., 210, 345-361. [https://doi.org/https://doi.org/10.1016/j.rse.2018.](https://doi.org/https://doi.org/10.1016/j.rse.2018.03.015)  
1031 [03.015](https://doi.org/https://doi.org/10.1016/j.rse.2018.03.015)

1032 Xia, H., Chen, Y., Song, C., Li, J., Quan, J., Zhou, G., 2022. Analysis of surface urban  
1033 heat islands based on local climate zones via spatiotemporally enhanced land  
1034 surface temperature. Remote Sens. Environ., 273, 112972. [https://doi.org/10.](https://doi.org/10.1016/j.rse.2022.112972)  
1035 [1016/j.rse.2022.112972](https://doi.org/10.1016/j.rse.2022.112972)

1036 Xiao, J., Moody, A., 2005. A comparison of methods for estimating fractional green  
1037 vegetation cover within a desert-to-upland transition zone in central New Mexico,  
1038 USA. Remote Sens. Environ., 98(2-3), 237-250. [https://doi.org/10.1016/](https://doi.org/10.1016/j.rse.2005.07.011)  
1039 [j.rse.2005.07.011](https://doi.org/10.1016/j.rse.2005.07.011)

1040 Xu, X., Tong, X., Plaza, A., Li, J., Zhong, Y., Xie, H., et al., 2018. A New Spectral-  
1041 Spatial Sub-Pixel Mapping Model for Remotely Sensed Hyperspectral Imagery.  
1042 IEEE Trans. Geosci. Remote Sens. 56(11): 6763-6778. [https://doi.org/10.1109/](https://doi.org/10.1109/TGRS.2018.2842748)  
1043 [TGRS.2018.2842748](https://doi.org/10.1109/TGRS.2018.2842748).

1044 Zanaga, D., Van De Kerchove, R., De Keersmaecker, W., Souverijns, N., Brockmann,  
1045 et al., 2021. ESA WorldCover 10 m 2020 v100. [https://doi.org/10.5281/zenodo.](https://doi.org/10.5281/zenodo.5571936)  
1046 [5571936](https://doi.org/10.5281/zenodo.5571936).

1047 Zhang, G., Ghamisi, P., Zhu, X., 2019. Fusion of heterogeneous earth observation data  
1048 for the classification of local climate zones. IEEE Trans. Geosci. Remote Sens. 57  
1049 (10), 7623–7642. <https://doi.org/10.1109/TGRS.2019.2914967>

1050 Zhong, Y., Wu, Y., Xu, X., Zhang, L., 2015. An Adaptive Subpixel Mapping Method  
1051 Based on MAP Model and Class Determination Strategy for Hyperspectral  
1052 Remote Sensing Imagery. IEEE Trans. Geosci. Remote Sens. 53 (3): 1411-1426.  
1053 <https://doi.org/10.1109/TGRS.2014.2340734>.  
1054