

# Geographical Characterisation of British Urban Form and Function using the Spatial Signatures Framework

Martin Fleischmann<sup>1,\*</sup> and Daniel Arribas-Bel<sup>1</sup>

<sup>1</sup>Geographic Data Science Lab, Department of Geography and Planning, University of Liverpool, Roxby Building , 74 Bedford St S , Liverpool , L69 7ZT, United Kingdom

\*corresponding author(s): Martin Fleischmann (m.fleischmann@liverpool.ac.uk)

## ABSTRACT

The spatial arrangement of the building blocks that make up cities matters to understand the rules directing their dynamics. Our study outlines the development of the national open-source classification of space according to its form and function into a single typology. We create a bespoke granular spatial unit, the enclosed tessellation, and measure characters capturing its form and function within a relevant spatial context. Using K-Means clustering of individual enclosed tessellation cells, we generate a classification of space for the whole of Great Britain. Contiguous enclosed tessellation cells belonging to the same class are merged forming spatial signature geometries and their typology. We identify 16 distinct types of spatial signatures stretching from wild countryside, through various kinds of suburbia to types denoting urban centres according to their regional importance. The open data product presented here has the potential to serve as boundary delineation for other researchers interested in urban environments and policymakers looking for a unique perspective on cities and their structure.

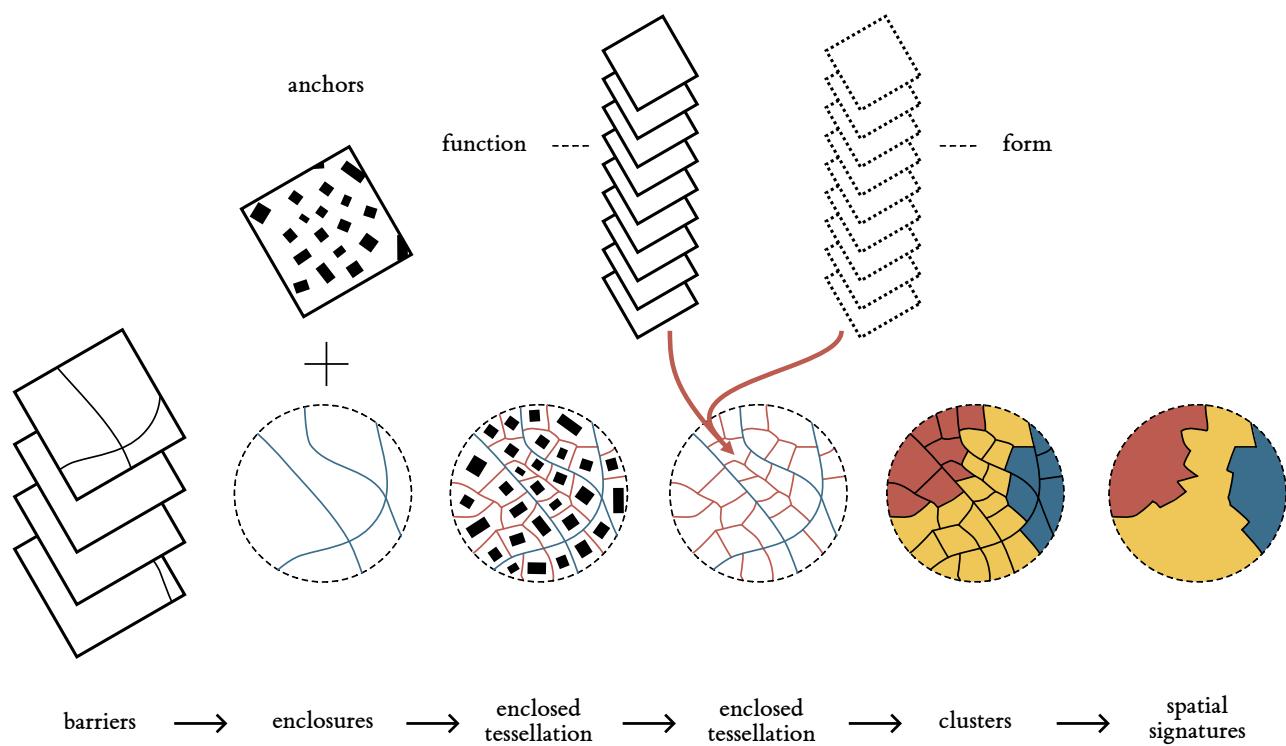
## Background & Summary

How the building blocks that make up cities are spatially arranged is worth quantifying and understanding. By "building blocks", we mean both the activities and agents that inhabit cities, as well as the (infra)structure that supports them. The former can be conceptualised as *urban function*, while the latter falls under the study of *urban form*. Understanding urban form and function is important for two main reasons. First, the combination of both *encodes* rich information about the history, character and evolution of cities. For example, the shape and properties of the street network encode the technology of the time (e.g., automobile); while the degree of mix in land uses can reflect cultural values. Second, the spatial pattern of urban form and function also acts as a frame that *influences* a variety of outcomes, from economic productivity to socio-economic cohesion to environmental sustainability.

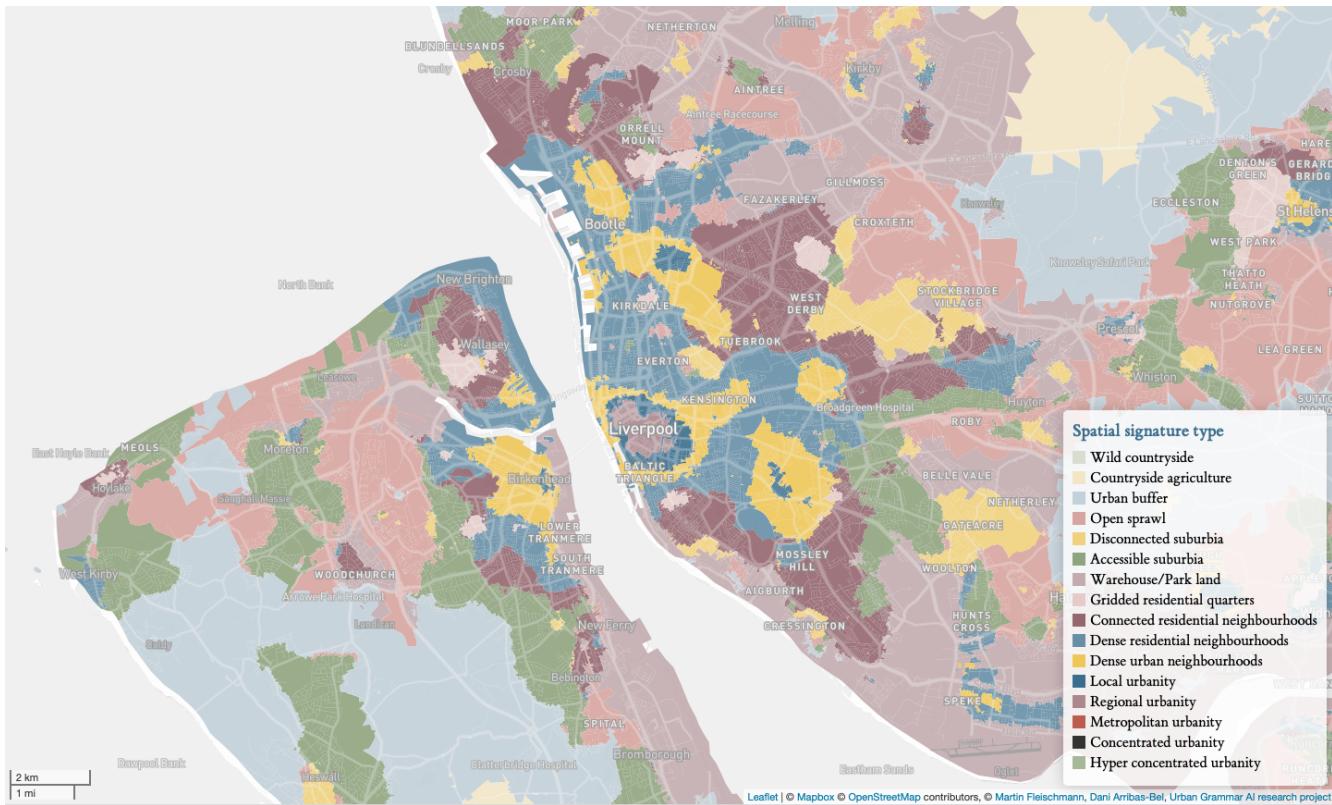
In this paper, we use the Spatial Signatures framework<sup>1,2</sup>, which develops a "characterisation of space based on form and function designed to understand urban environments"<sup>1</sup>. Spatial signatures are theory-informed, data-driven computable classes that describe the form and function of a consistent patch of geography. Figure 1 presents an overview of the development of a spatial signature classification. We build a series of enclosures that we combine with building footprints to further subdivide geographical space into what we call enclosed tessellation cells (ETCs). We then attach form and function characters to each of these subdivisions, and use those to group them into consistent and differentiated classes we call signatures. Each phase is expanded in detail in the next section.

We introduce an open data product (ODP<sup>3</sup>) containing a classification of spatial signatures for Great Britain (illustrated in a figure 2). In doing so, we provide an analysis-ready layer that brings together urban form and function consistently, in detail, and at national scale. To the best of our knowledge, this is the first dataset capturing urban form and function published both with a degree of detail and scale as ours. Our results are based on the analysis of more than 14 million of ETCs, to each of which we attach more than 300 characters capturing a wide range of aspects relating to urban form and function. We provide access to both granular geographical boundaries of the delineated spatial signatures as well as measurements for each character at the signature level. The ODP also includes a web map that allows exploration without any technical requirement other than a web browser, and we have open sourced all the code, including details on the computational backend. The uniqueness of our ODP makes it challenging to set up a technical validation as a comparison with existing datasets. Nevertheless, we relate our signatures to a few well-established data products that capture each a subset of the form and function dimensions we consider. Our results are encouraging in that they show broad agreement in expected areas, but also highlight aspects that can only be discovered when considering form and function in tandem.

The approach and outputs presented bring several benefits to a range of stakeholders interested in cities. This spatial signatures ODP provides insight generated from detailed, comprehensive and computationally intensive data analysis and



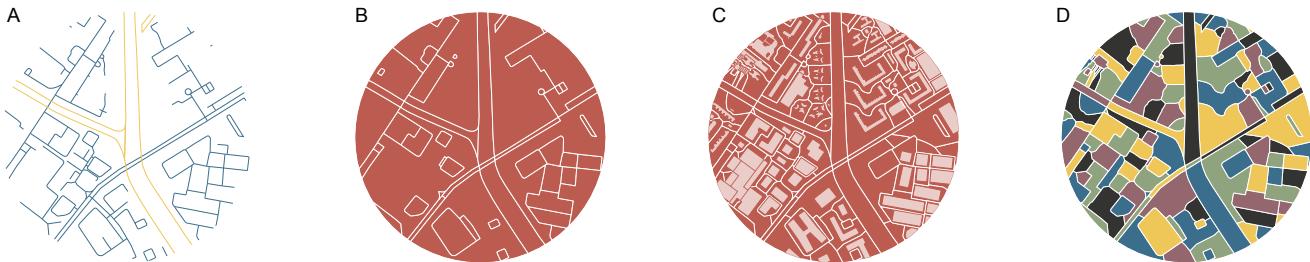
**Figure 1.** Diagram illustrating the sequential steps leading to the delineation of spatial signatures. From a series of enclosing components, to enclosures, enclosed tessellation (ET), the addition of form and function characters to ET cells, and the development of spatial signatures.



**Figure 2.** Illustration of a classification of spatial signatures in Liverpool and Birkenhead area, in the north west of England.

39 presents it in a way that is easy to access, work with and integrate into larger projects. Together with the importance of form  
 40 and function discussed above, we anticipate the output will be relevant to both academic researchers as well as policymakers  
 41 and practitioners. As a framework, the spatial signatures provide a flexible yet generalisable way to understand, characterise  
 42 and quantify urban form and function. One way to understand our results is as an application to Great Britain of a more general  
 43 approach to quantitatively characterise the spatial dimension of cities. As such, our conceptual approach can be applied in many  
 44 more local contexts and regions beyond Great Britain. It is true that Great Britain currently represents an unusual case in that it  
 45 is specially “data dense”, with a large variety of open data that may not be readily available in other parts of the world. However,  
 46 given form and function reinforce each other, spatial signatures are designed to be robust to variations in the specific data  
 47 sources used, and two different classifications do not need to be based on exactly the same data to be useful. At the same time,  
 48 we note that the combination of volunteered geographic information (e.g., OpenStreetMap) and technologies such as modern  
 49 satellites and artificial intelligence are filling many of these gaps very rapidly, and we anticipate near-future developments that  
 50 will make the implementation of classifications such as the one presented here possible in almost any (urban) area of the planet.  
 51 In this sense, our ODP (data, code, and methodology) can be a useful illustration for researchers and practitioners who, even if  
 52 not specifically interested in the British use case, would like to implement a similar approach on their own.

53 As illustration of potential applications, we provide two. The spatial signatures may be used to delineate types of (origin  
 54 and destination) locations in mobility analysis, that could unveil patterns of commuting or migration in situations like the  
 55 COVID-19 pandemic. A second application may focus directly on supporting policy on inequalities. For example the spatial  
 56 signatures can underpin analysis on equality of access to services and amenities within the UKs Levelling Up agenda<sup>4</sup>, using  
 57 them to target areas based on their signature type, since they will share key structural components. It is important to note we do  
 58 not expect signatures to focus on a single aspect of urban environment as, for example, Local Climate Zones<sup>5</sup> do with climate,  
 59 but instead on a wider range of uses due to their inclusion of both form and function and a data driven nature reflecting the  
 60 specific place rather than abstract conceptual classes. In this respect, we hope the present paper serves not only to document our  
 61 own work but to inspire future efforts aimed at urban form and function.



**Figure 3.** Diagram illustrating the sequential steps leading to the delineation of enclosed tessellation. From a series of enclosing components, where blue are streets and yellow river banks (A), to enclosures (B), incorporation of buildings as anchors (C) to final tessellation cells (D).

## 62 Methods

63 The method of identification of spatial signatures consists of three top-level steps. First, we delineate a spatial unit of analysis  
 64 that reflects the structure of urban phenomena on a very granular level. Then we characterise each of them according to form  
 65 and function, capturing the nature of each unit and its spatial context. Finally, we use cluster analysis to derive a typology of  
 66 our spatial units that, once combined into contiguous areas, forms a typology of spatial signatures.

### 67 Spatial unit

68 The first major methodological decision relates to the definition of the spatial unit. An ideal candidate needs to reflect space in  
 69 a granular manner, and we argue it should fulfil three conditions. First, it should be *indivisible*, meaning that any subdivision  
 70 would result in a unit that is incapable of capturing the nature of urban form and function. Second, it needs to be *internally*  
 71 *consistent* - it should always reflect only a single signature type. Last, it should be geographically *exhaustive*, covering the  
 72 entirety of the study area.

73 Spatial units used in literature can be split into three groups. One is using administrative boundaries like city regions<sup>6</sup>,  
 74 wards or census output areas<sup>7</sup>, that are convenient to obtain and can be easily linked to auxiliary data. However, those rarely  
 75 reflect the morphological composition of urban space and, in some cases, may even “obscure morphologic reality”<sup>8</sup>. At the  
 76 same time, most of them are divisible, and larger units are not always internally consistent. Another group is based on arbitrary  
 77 uniform grids linked either to spatial indexing methods like H3<sup>9</sup> or Ordnance Survey National Grid, or to ancillary data of  
 78 remote sensing or other origins like a WorldPop grid<sup>10</sup>. Grids however cannot be considered internally consistent as they do not  
 79 consider the underlying structure of the landscape. Finally, urban morphology studies tend to use morphological elements as  
 80 street segments<sup>11</sup>, blocks<sup>12</sup>, buildings<sup>13</sup> or plots<sup>14</sup> as units of analysis. Some of those could be seen as indivisible and internally  
 81 consistent, but since they are largely based on built-up fabric, they are not exhaustive. For example, in areas without any  
 82 building or street, there is no spatial unit to work with. Plots could be theoretically considered as exhaustive, consistent and  
 83 indivisible, but there is no accepted conceptual definition and unified geometric representation<sup>15</sup>.

84 We are, therefore, proposing an application of an alternative spatial unit called *enclosed tessellation cell* (ETC), defined  
 85 as "the portion of space that results from growing a morphological tessellation within an enclosure delineated by a series of  
 86 natural or built barriers identified from the literature on urban form, function and perception"<sup>1</sup>. ETCs follow the morphological  
 87 tradition in that it is based on the physical elements of an environment but overcome the drawbacks of conventionally used  
 88 units. Its geometry is generated in the three steps illustrated in Figure 3. First, a set of features representing physical barriers  
 89 subdividing space, in our case composed of the street network, railways, rivers and a coastline, is combined, generating a layer  
 90 of boundaries (3 A). These then partition space into smaller enclosed geometries called *enclosures* (3 B), which can be very  
 91 granular or very coarse depending on the geographic context. In dense city centres where a single enclosure represents a single  
 92 block is a high frequency of small enclosures. At the same time, in the countryside, this approach leads to very few large  
 93 enclosures as their delimiters are far away from each other. Enclosures are then combined with building footprints (3 B), which  
 94 act as anchors in space and potentially subdivide enclosures into enclosed tessellation cells using the morphological tessellation  
 95 algorithm<sup>16</sup> (3 D), a polygon-based adaptation of Voronoi tessellation. The resulting geometries are indivisible as they contain,  
 96 at most, a single anchor building, internally consistent due to their granularity and link to morphological elements composing  
 97 urban fabric, and geographically exhaustive as they cover an entire area limited by specified boundaries.

98 In our ODP for Great Britain, street networks are extracted from OS Open Roads datasets<sup>17</sup> representing simplified road  
 99 centrelines cleaned of underground road segments. Railways are retrieved from OS OpenMap - Local<sup>18</sup> ("RailwayTrack"  
 100 layer) which captures surface railway tracks. Rivers are extracted from OS OpenRivers<sup>19</sup> representing river network of GB  
 101 as centrelines, and a coastline is retrieved from OS Strategi®<sup>20</sup>, capturing coastline as a continuous line geometry. Building

102 geometry is extracted, again, from OS OpenMap - Local ("Building" layer) and represents generalised building footprint  
103 polygons.<sup>1</sup>

## 104 **Characterisation of space**

105 Spatial signatures capture the character of the built and unbuilt environment based on two components - form and function.  
106 Each of them is quantified at the level of individual ETCs using methods appropriate for each specific dataset. While form is  
107 described using urban morphometrics (i.e. quantitative analysis of urban form)<sup>21</sup>, function is a composite of a variety of data  
108 inputs. We outline each component with a bit more detail below.

### 109 **Form**

110 Morphometric characterisation of urban form is based on the numerical description of four elements capturing the built  
111 environment - buildings, streets, ETCs, and enclosures - and reflects their patterns based on six categories of characters:  
112 dimensions, shapes, spatial distribution, intensity, connectivity and diversity<sup>22</sup>. Each element is considered across different  
113 scales, from the measurement of individual geometries, to relations of neighbouring geometries, to a graph-based analysis of  
114 the street network. The combination of elements, categories and scales results in a set of 59 individual morphometric characters  
115 listed in the table 1. The selection builds on the principles outlined by<sup>21</sup> and later explored by<sup>23</sup>, both following the rules  
116 derived by<sup>24</sup>. The gist is to include as many characters present in literature as is feasible, while minimising potential collinearity  
117 and limiting redundancy of information.

118 However, measuring individual characters is not enough to understand the predominant spatial patterns. For some types  
119 of urban environment, high heterogeneity is not uncommon. This means that using, for example, areas of building footprints  
120 would, in most cases, result in largely discontinuous clusters that do not capture the pattern within an area. Therefore, we  
121 represent each of the morphometric characters using three summary variables reflecting statistical distributions of measured  
122 data within a spatial context of each ETC. Context is defined as tenth order of contiguity computed across the mesh composed  
123 of contiguous ETCs as illustrated in figure 4. Furthermore, each value is weighted by the inverse distance between so-called  
124 poles of inaccessibility (defined as a centre of a maximum inscribed circle) of each ETC. Three proxy variables then capture the  
125 first, the second and the third quartile of the resulting weighted distribution. Such a characterisation can capture the contextual  
126 tendency of each morphometric character and hence identify contiguous clusters in both homogenous and heterogeneous urban  
127 tissues. These contextual values are then used as an input for cluster analysis while the original non-contextualised versions are  
128 left out, making the final form component composed of 177 contextual characters.

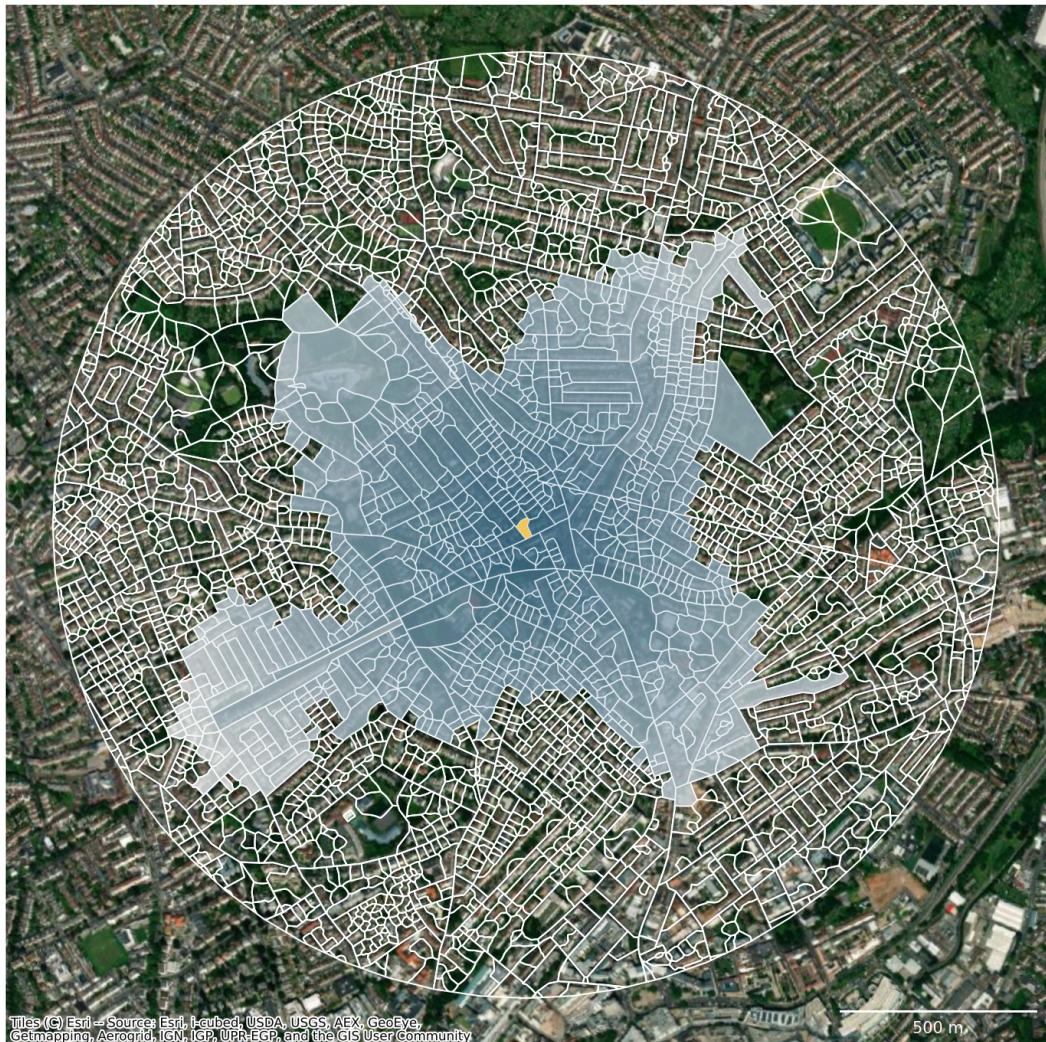
### 129 **Function**

130 Characterisation of the function component uses a different approach. While data describing urban form are not generally  
131 available in a processed format, forcing us to employ morphometric approaches, different aspects of function are often available  
132 as open data products. We guide the compilation of functional characters following three main principles: first, we identify from  
133 the literature on urban function key areas to be represented; second, we translate those abstract areas into measurable features;  
134 and third, we select open data available in for Great Britain that allows for the redistribution of derivative products. With a  
135 list of function characters selected, the main goal of our characterisation of ETCs based on function is to develop appropriate  
136 transfer methods to link data published as grids or linked to administrative boundaries to ETCs.

137 In this work, we are using five different transfer methods: Areal interpolation, Building-based dasymetric areal interpolation<sup>37</sup>  
138 using building footprint area, Network-constrained accessibility, Euclidean accessibility, and Zonal statistics. Areal  
139 interpolation is used when the functional data covers the entirety of space in the form of polygon geometry and when there is no  
140 assumption that the phenomena it captures are linked directly to the human population, such as land cover data. When there is  
141 an assumption of relation to the population, building-based dasymetric areal interpolation is used instead. The main difference  
142 is that instead of ETC polygons, building footprint polygons linked to individual ETCs are used as a target of interpolation.  
143 That ensures that data like population estimates are linked to ETCs proportionally to their ability to house population rather  
144 than by their area. Network-constrained accessibility is used when the input data represents points of interest like locations of  
145 supermarkets. Points are then snapped to the nearest node on the street network and linked to the ETCs through the count of  
146 observations accessible from the cell within 15 minutes of walk (1200m on the street network) and a distance to the nearest  
147 point. In some cases, Euclidean (as-crow-flies) accessibility is measured instead to accommodate for phenomena that are often  
148 outside the reach of a drivable network like water bodies. Zonal statistics are used to transfer data originally stored in a raster  
149 format to ETCs as the mean value of raster pixels intersecting each polygon geometry. Finally, characters based on interpolation  
150 and zonal statistics are expressed using their contextual versions following the method used for form characters to, again, reflect  
151 the contextual pattern of measured values. As in the case of morphometric characters, only contextual versions are then used in  
152 the cluster analysis. The selection of datasets and the chosen transfer method are listed in the table 2.

<sup>1</sup>Note that the dataset does not distinguish between individual buildings when they are adjacent (e.g. perimeter block composed of multiple buildings is represented by a single polygon).

character	category	reference
area of building	dimension	25
perimeter of building	dimension	26
courtyard area of building	dimension	27
circular compactness of building	shape	21
corners of building	shape	28
squareness of building	shape	28
equivalent rectangular index of building	shape	29
elongation of building	shape	28
centroid - corner distance deviation of building	shape	23
centroid - corner mean distance of building	dimension	27
orientation of building	distribution	27
street alignment of building	distribution	27
cell alignment of building	distribution	23
longest axis length of ETC	dimension	23
area of ETC	dimension	13
circular compactness of ETC	shape	23
equivalent rectangular index of ETC	shape	23
orientation of ETC	distribution	23
covered area ratio of ETC	intensity	30
length of street segment	dimension	12
width of street profile	dimension	11
openness of street profile	distribution	11
width deviation of street profile	diversity	11
linearity of street segment	shape	11
area covered by edge-attached ETCs	dimension	23
buildings per meter of street segment	intensity	23
area covered by node-attached ETCs	dimension	23
alignment of neighbouring buildings	distribution	31
mean distance between neighbouring buildings	distribution	31
perimeter-weighted neighbours of ETC	distribution	23
area covered by neighbouring cells	dimension	23
reached ETCs by neighbouring segments	intensity	23
reached area by neighbouring segments	dimension	23
node degree of junction	distribution	32
mean distance to neighbouring nodes of street n...	dimension	23
mean inter-building distance	distribution	33
weighted reached enclosures of ETC	intensity	23
reached ETCs by tessellation contiguity	intensity	23
reached area by tessellation contiguity	dimension	23
area of enclosure	dimension	21
perimeter of enclosure	dimension	12
circular compactness of enclosure	shape	27
equivalent rectangular index of enclosure	shape	29
compactness-weighted axis of enclosure	shape	34
orientation of enclosure	distribution	12
perimeter-weighted neighbours of enclosure	distribution	23
area-weighted ETCs of enclosure	intensity	23
local meshedness of street network	connectivity	34
mean segment length within 3 steps	dimension	23
local cul-de-sac length of street network	dimension	23
reached area by local street network	dimension	23
reached ETCs by local street network	intensity	23
local node density of street network	intensity	23
local proportion of cul-de-sacs of street network	connectivity	35
local proportion of 3-way intersections of stre...	connectivity	32
local proportion of 4-way intersections of stre...	connectivity	32
local degree weighted node density of street ne...	intensity	21
local closeness of street network	connectivity	36



**Figure 4.** Illustration of a definition of spatial context used to capture the distribution of values around each ET cell. For the yellow ET cell in the middle, we propose to define a neighbourhood of 10 topological steps on the tessellation and weight the importance of each cell within such an area by inverse distance between poles of inaccessibility of each cell.

character	data
Population	Population estimates
Night lights	Night Lights
Workplace population [Agriculture, energy and water]	Workplace population
Workplace population [Manufacturing]	Workplace population
Workplace population [Construction]	Workplace population
Workplace population [Distribution, hotels and restaurants]	Workplace population
Workplace population [Transport and communication]	Workplace population
Workplace population [Financial, real estate, professional and administrative activities]	Workplace population
Workplace population [Public administration, education and health]	Workplace population
Workplace population [Other]	Corine land cover
Land cover [Airports]	Corine land cover
Land cover [Non-irrigated arable land]	Corine land cover
Land cover [Industrial or commercial units]	Corine land cover
Land cover [Salt marshes]	Corine land cover
Land cover [Estuaries]	Corine land cover
Land cover [Sport and leisure facilities]	Corine land cover
Land cover [Green urban areas]	Corine land cover
Land cover [Discontinuous urban fabric]	Corine land cover
Land cover [Pastures]	Corine land cover
Land cover [Broad-leaved forest]	Corine land cover
Land cover [Mineral extraction sites]	Corine land cover
Land cover [Port areas]	Corine land cover
Land cover [Road and rail networks and associated land]	Corine land cover
Land cover [Water bodies]	Corine land cover
Land cover [Land principally occupied by agriculture, with significant areas of natural vegetation]	Corine land cover
Land cover [Mixed forest]	Corine land cover
Land cover [Peat bogs]	Corine land cover
Land cover [Natural grasslands]	Corine land cover
Land cover [Moors and heathland]	Corine land cover
Land cover [Transitional woodland-shrub]	Corine land cover
Land cover [Continuous urban fabric]	Corine land cover
Land cover [Intertidal flats]	Corine land cover
Land cover [Sea and ocean]	Corine land cover
Land cover [Coniferous forest]	Corine land cover
Land cover [Construction sites]	Corine land cover
Land cover [Sparsely vegetated areas]	Corine land cover
Land cover [Bare rocks]	Corine land cover
Land cover [Inland marshes]	Corine land cover
Land cover [Dump sites]	Corine land cover
Land cover [Fruit trees and berry plantations]	Corine land cover
Land cover [Complex cultivation patterns]	Corine land cover
Land cover [Beaches, dunes, sands]	Corine land cover
Land cover [Water courses]	Corine land cover
Land cover [Burnt areas]	Corine land cover
Land cover [Agro-forestry areas]	Corine land cover
Land cover [Coastal lagoons]	Corine land cover
NDVI	NDVI
Supermarkets [distance to nearest]	Retail POIs (supermarkets)
Supermarkets [counts within 1200m]	Retail POIs (supermarkets)
Listed buildings [distance to nearest]	Listed Buildings
Listed buildings [counts within 1200m]	Listed Buildings
FHRS points [distance to nearest]	Food Hygiene Rating Scheme Rating
FHRS points [counts within 1200m]	Food Hygiene Rating Scheme Rating
Cultural venues [distance to nearest]	Culture (theatres, cinemas)
Cultural venues [counts within 1200m]	Culture (theatres, cinemas)
Water bodies [distance to nearest]	Water bodies
Retail centres [distance to nearest]	Retail centres

**Table 2.** Functional characters used to describe the function component of spatial signatures. For details of the

153 **Cluster analysis**

154 When combined, contextual summaries of form and function characters (or characters themselves when they are reflecting  
155 the context by definition) compose a dataset describing each ETC by 331 variables (177 contextual characters representing  
156 59 initial characters for form and 154 for function composed of 144 contextual characters representing 48 characters that do  
157 not capture context by design and 10 accessibility-based characters that do). Assigning equal weight to each variable, we  
158 standardize them applying Z-score normalization, and use them as input for K-Means cluster analysis. Although collinearity is  
159 likely to be present between several of them, we do not view this as a problem: we select each character not from a purely  
160 statistical point of view (i.e., which ones will be more effective at segmenting the dataset), but instead from a conceptual one.  
161 Each variable has been identified by the literature on urban form and function as a relevant aspect that contributes to collectively  
162 characterising these two more abstract concepts. We thus see this situation as a way of adding robustness to the measurement of  
163 more conceptual notions which are ultimately our aim. We opt for K-Means because we consider it strikes a compromise in the  
164 trade-off between performance and scalability. K-Means is widely used in the literature on unsupervised learning, and in much  
165 of that concerning the clustering of geographic entities<sup>38</sup>. To select the algorithm, we experimented with a random subset of our  
166 dataset, comparing K-Means with alternatives such as Gaussian Mixture Models (GMM) or Self-Organising Maps (SOM).  
167 We found results from the latter two were not notably better in terms of cluster compactness and qualitative examination of  
168 the geographic clusters, but were significantly slower in computation runtime, posing serious challenges to be run at scale.  
169 Although K-Means does not consider space explicitly, our approach incorporates information about the geographic context of  
170 each observation through the operation described above and illustrated in Figure 4. We prefer this over a spatially-constrained  
171 algorithm (e.g., SKATER<sup>39</sup>) that restricts the clustering only among spatially contiguous observations because we are not  
172 interested in areas that are spatially contiguous unless they are sufficiently similar to each other on the attribute space. Our  
173 contextual approach is more similar to spatially-encouraged algorithms such as the GeoSOM<sup>40</sup> or spatially-encouraged spectral  
174 clustering<sup>41</sup> that incorporate geographic proximity when clustering but do not restrict. Our choice in this case was led by its  
175 scalability over other such algorithms. Nevertheless, we consider this a fruitful avenue for future research.

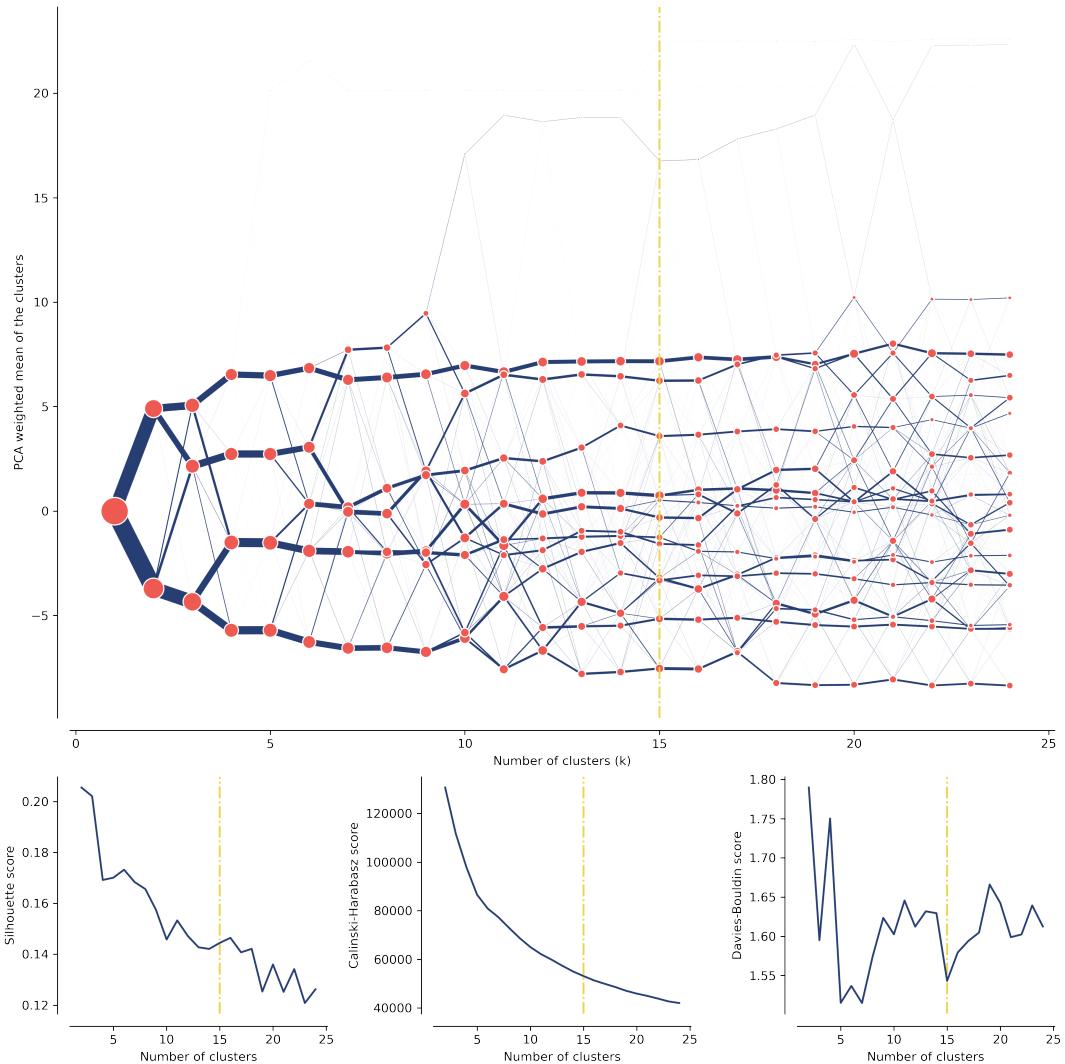
176 Due to the nature of the selected K-Means clustering, the step preceding the final analysis is the selection of an optimal  
177 number of clusters. We use the clustergram exploratory method<sup>42</sup>, reflecting the behaviour of different options, the relationship  
178 between clustering solutions regarding the allocation of individual observations to classes, and the separation between the  
179 clusters within each tested solution (figure 5). Clustergram is further accompanied by measures of internal validation measures  
180 - the Silhouette score diagram, Calinski-Harabasz index<sup>43</sup> and Davies-Bouldin index<sup>44</sup>. The optimal number of classes is  
181 selected based on the interpretation of clustergram supported by additional measures aiming at a balance between cluster  
182 separation and an appropriate detail of resulting classification. We use mini batch K-Means with a batch size of 1,000,000 and  
183 100 initialisations to create the clustergram and test number of clusters between 2 and 25. The results indicate 10 clusters as an  
184 optimal solution. The final clustering solution is generated using mini batch K-Means with a batch size of 1,000,000 and 1,000  
185 initialisations to ensure the stability of the outcome.

186 The results of the clustering capture the first group of a national signature classification composed of ten clusters. However,  
187 since the classified ETCs cover the entirety of space, from vast natural open spaces to dense city centres, it may result in only a  
188 few classes representing urban areas. While that is caused by the variable heterogeneity of our dataset in combination with  
189 K-Means clustering, the measured characters have the ability to further distinguish classes of already identified clusters. As  
190 spatial signatures are focused on the urban environment, we further subdivide those clusters covering a substantial portion of  
191 urban areas using another iteration of K-Means clustering (one class into nine and another into three clusters). Both subdivisions  
192 were created using standard K-Means (single batch) using 1,000 initialisations. The resulting classification then provide  
193 classification capturing the typology of spatial signatures with a detailed focus on urban development.

194 Finally, individual spatial signature geometries are generated as a combination of adjacent ETCs belonging to the same  
195 signature class. To describe each geometry and each signature type, we measure mean values of the original, non-contextualised  
196 characters, and release it as additional descriptive tables. The resulting numerical profile of each signature type is available as  
197 table 3. Table 4 contains pen portraits derived from these numerical profiles.

198 **Data Records**

199 The data product described in this article is available through the Consumer Data Research Centre Open Data repository  
200 available at <https://data.cdrc.ac.uk/dataset/spatial-signatures-great-britain> under the Open Government Licence v3.0 license and  
201 archived at <https://doi.org/10.6084/m9.figshare.16691575.v1>. The dataset stored in the repository contains a GeoPackage with a  
202 signature geometry (OSGB36 / British National Grid (EPSG : 27700) CRS) and related signature type, plain-text pen portraits  
203 describing individual signature types, a series of CSV files describing individual signatures and signature types, and a CSV  
204 files linking signature types to the Output Area and Lower Super Output Area geometry. An online interactive map of spatial  
205 signatures for the whole of Great Britain is available on the project website (<https://urbangrammarai.xyz/great-britain>).



**Figure 5.** Clustergram and relevant metrics of a goodness of fit (Silhouette score, Calinski-Harabasz score, Davies-Bouldin score) for tested numbers of clusters. The clustergram suggest two potential solutions, the very conservative option of 4 clusters and 10 clusters selected as an optimal result (indicated by a vertical yellow line).

type	Accessible suburbia	Connected res.
area of building	176.95	
perimeter of building	53.90	
courtyard area of building	0.48	
circular compactness of building	0.53	
corners of building	4.25	
squareness of building	0.78	
equivalent rectangular index of building	0.99	
elongation of building	0.64	
centroid - corner mean distance of building	9.60	
centroid - corner distance deviation of building	0.36	
orientation of building	19.56	
longest axis length of ETC	50.84	
area of ETC	1147.25	
circular compactness of ETC	0.47	
equivalent rectangular index of ETC	0.97	
orientation of ETC	20.40	
covered area ratio of ETC	0.19	
cell alignment of building	7.38	
alignment of neighbouring buildings	5.31	
mean distance between neighbouring buildings	17.82	
perimeter-weighted neighbours of ETC	0.04	
area covered by neighbouring cells	8620.11	
weighted reached enclosures of ETC	0.00	
mean inter-building distance	21.97	
width of street profile	28.38	
width deviation of street profile	3.30	
openness of street profile	0.42	
length of street segment	187.61	
linearity of street segment	0.93	
mean segment length within 3 steps	2327.31	
node degree of junction	2.87	
local meshedness of street network	0.08	
local proportion of 3-way intersections of street network	0.74	
local proportion of 4-way intersections of street network	0.07	
local proportion of cul-de-sacs of street network	0.19	
local closeness of street network	0.00	
local cul-de-sac length of street network	228.58	
square clustering of street network	0.03	
mean distance to neighbouring nodes of street network	132.49	
local node density of street network	0.02	
local degree weighted node density of street network	0.03	
street alignment of building	8.73	
area covered by node-attached ETCs	22426.36	
area covered by edge-attached ETCs	36496.96	
buildings per meter of street segment	0.11	
reached ETCs by neighbouring segments	49.09	
reached area by neighbouring segments	113290.06	
reached ETCs by local street network	166.98	
reached area by local street network	451276.21	
reached ETCs by tessellation contiguity	36.80	
reached area by tessellation contiguity	60511.46	
area of enclosure	242778.35	
perimeter of enclosure	2046.29	
circular compactness of enclosure	0.40	
equivalent rectangular index of enclosure	0.85	
compactness-weighted axis of enclosure	515.77	
orientation of enclosure	19 <sup>14</sup> / <sub>21</sub>	
perimeter-weighted neighbours of enclosure	0.01	
area-weighted ETCs of enclosure	36.32	
Perimeter of building	1.51	

	0
Wild countryside	In “Wild countryside”, human influence is the least intensive. This signature covers large open spaces.
Countryside agriculture	“Countryside agriculture” features much of the English countryside and displays a high degree of agricultural activity.
Urban buffer	“Urban buffer” can be characterised as a green belt around cities. This signature includes mostly agricultural land.
Open sprawl	“Open sprawl” represents the transition between countryside and urbanised land. It is located in the outskirts of towns and cities.
Disconnected suburbia	“Disconnected suburbia” includes residential developments in the outskirts of cities or even towns.
Accessible suburbia	“Accessible suburbia” covers residential development on the urban periphery with a relatively legal and planned layout.
Warehouse/Park land	“Warehouse/Park land” covers predominantly industrial areas and other work-related developments.
Gridded residential quarters	“Gridded residential quarters” are areas with street networks forming a well-connected grid-like pattern.
Connected residential neighbourhoods	“Connected residential neighbourhoods” are relatively dense urban areas, both in terms of population density and connectivity.
Dense residential neighbourhoods	A “dense residential neighbourhood” is an abundant signature often covering large parts of cities.
Dense urban neighbourhoods	“Dense urban neighbourhoods” are areas of inner-city with high population and built-up density.
Local urbanity	“Local urbanity” reflects town centres, outer parts of city centres or even district centres. In all cases, it is a highly developed area.
Regional urbanity	“Regional urbanity” captures centres of mid-size cities with regional importance such as Liverpool.
Metropolitan urbanity	Signature type “Metropolitan urbanity” captures the centre of the largest cities in Great Britain such as London.
Concentrated urbanity	“Concentrated urbanity” is a signature type found in the city centre of London and nowhere else in the country.
Hyper concentrated urbanity	The epitome of urbanity in the British context. “Hyper concentrated urbanity” is a signature type found in the central business districts of major cities like London and Manchester.

**Table 4.** Interpretative pen portraits characterising each signature type based on its numerical profile.

	relative importance
covered area ratio of ETC (Q1)	0.036944
covered area ratio of ETC (Q2)	0.031717
perimeter-weighted neighbours of ETC (Q2)	0.023476
mean inter-building distance (Q2)	0.016662
area of ETC (Q3)	0.016005
area covered by node-attached ETCs (Q3)	0.014813
longest axis length of ETC (Q2)	0.014501
weighted reached enclosures of ETC (Q1)	0.014115
reached area by neighbouring segments (Q3)	0.014000
reached area by neighbouring segments (Q1)	0.013904

**Table 5.** Relative importance of top 10 most important characters in predicting spatial signature types using the Random Forest model.

## 206 Composition and comparison

### 207 Character importance

208 The characters used in the cluster analysis have each different importance in distinguish between signature types. Those  
 209 characters which spatial distribution most closely matches the distribution of signatures can be seen as more important than  
 210 those that are seemingly random or mostly invariant (as some of the land cover classes are). Unpacking the importance of  
 211 individual characters from K-Means clustering cannot be done directly, but a useful method is to train a supervised model, in  
 212 our case Random Forest, designed to predict individual signature types from input data. Such a model then provides a feature  
 213 importance - a relative measure of a strength of each character in distinguishing between the types. The results of this approach  
 214 are shown in a table 5. As you can see, form-based characters dominate the top 10 characters but it is worth noting that these  
 215 top 10 characters together bear only 0.196 of the overall importance.

216 A similar exercise can be done on a level of individual clusters, with a binary Random Forest model trained to distinguish  
 217 that particular class from the other. Resulting relative importance of top 10 characters for each signature type is presented in a  
 218 table 6. While it is clear that form-based characters still dominate the prediction, the more urban signature types are, the higher  
 219 the importance of function seems to be. Complete tables with all characters are available as online tables 1 and 2.

### 220 Comparison

221 Spatial signatures are unique as a classification method, limiting the potential validation. Therefore, we rather present a  
 222 comparison of signatures and ancillary datasets capturing conceptually similar aspects of the environment. We compare the

Wild countryside name	Country-side agriculture rel. importance	Country-side agriculture name	Gridded res. rel. importa
longest axis length of ETC (Q1)	0.196609	covered area ratio of ETC (Q1)	0.154
covered area ratio of ETC (Q2)	0.151118	covered area ratio of ETC (Q2)	0.144
covered area ratio of ETC (Q1)	0.145754	mean inter-building distance (Q2)	0.078
area of ETC (Q2)	0.096485	area of ETC (Q2)	0.072
perimeter-weighted neighbours of ETC (Q3)	0.075078	area covered by node-attached ETCs (Q2)	0.066
reached area by neighbouring segments (Q1)	0.048869	mean distance to neighbouring nodes of street n...	0.066
reached area by tessellation contiguity (Q1)	0.018289	reached area by neighbouring segments (Q1)	0.062
area of ETC (Q3)	0.015991	Land cover [Discontinuous urban fabric] (Q2)	0.055
mean distance between neighbouring buildings (Q2)	0.015013	perimeter-weighted neighbours of ETC (Q2)	0.021
mean inter-building distance (Q2)	0.010559	longest axis length of ETC (Q2)	0.020

**Table 6.** Relative importance of top 10 most important characters for each signature type in predicting using the Random Forest model.

signatures with four of such datasets, each focusing on a different classification perspective, but all related to our classification to a degree when we can assume there will be a measurable level of association between the two:

- WorldPop settlement patterns of building footprints (2021)<sup>10</sup>
- Classification of Multidimensional Open Data of Urban Morphology (MODUM) (2015)<sup>7</sup>
- Copernicus Urban Atlas (2018)<sup>45</sup>
- Local Climate Zones (2019)<sup>46</sup>

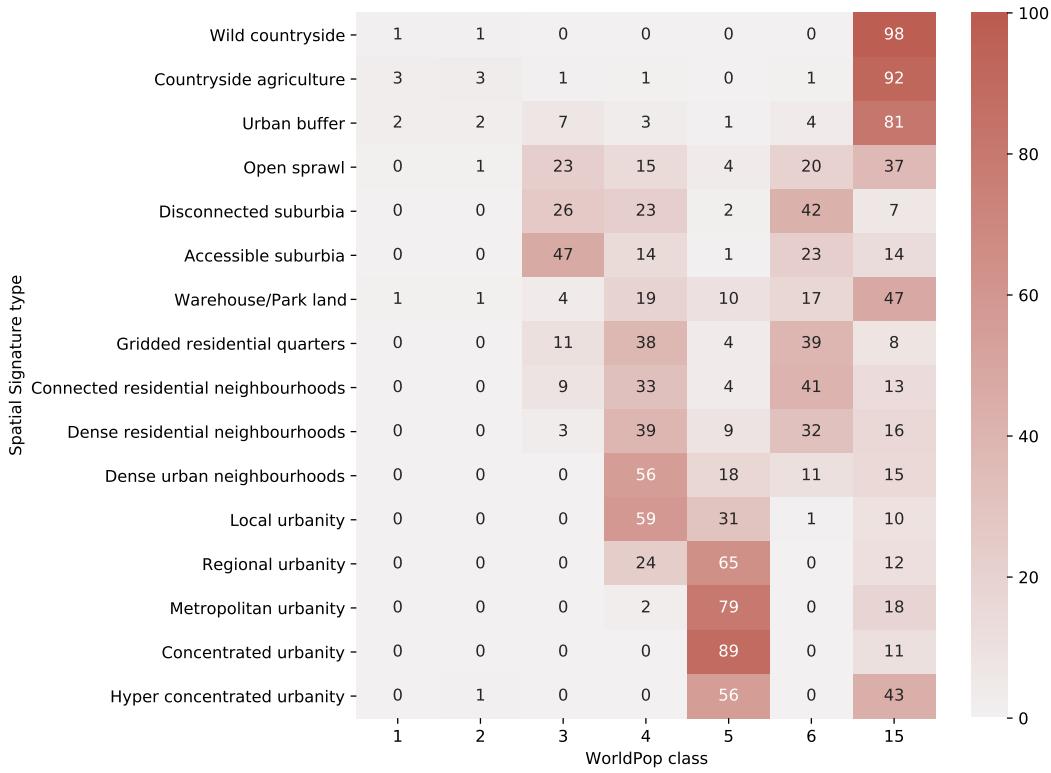
## Comparison approach

All datasets, spatial signatures and those selected for a comparison contain a categorical classification of space linked to their unique geometry. The first requirement to be able to compare data products is to transfer their information to the same geometry. We take two approaches for this step, depending on the dataset we are comparing the signatures with: an interpolation of one set of polygon-based data to another (input to ETCs); or the conversion of spatial signatures to the raster representation matching an input raster, which is computationally more efficient when one of the layers is already a raster. The second step is a statistical comparison of two sets of classification labels, one representing spatial signature typology and the other comparison classes. We use contingency tables and Pearson's  $\chi^2$  test to determine whether the frequencies of observed (signature types) and expected (comparison types) labels significantly differ in one or more categories. Furthermore, we use Cramér's V statistics<sup>47</sup> to assess the strength of the association.

## WorldPop settlement patterns of building footprints

WorldPop settlement patterns of building footprints dataset aims to derive a typology of morphological patterns based on a gridded approach with cells of 100x100m, and building footprints. Authors measure six morphometric characters linked to the grid cells and use them as input for an unsupervised clustering algorithm leading to a six-class typology. As the classification is dependent on building footprints, grid cells that do not contain any information on the building-based pattern are treated as missing in the final data product. For the comparison, this *missing* category is treated as a single class. It is assumed that the top-level large scale patterns detected by the WorldPop method and spatial signatures will provide similar results. However, there will be differences caused by the inclusion of function in spatial signatures, higher granularity of both initial spatial units and the resulting classification (6 vs 19 classes).

Signature typology is rasterized and linked to the WorldPop grid. The resulting contingency table is shown in Figure 6. There is a significant relationship between two typologies,  $\chi^2(114, N = 22993921) = 13341832, p < .001$ . The strength of association measured as Cramér's V is 0.311, indicating moderate association. The contingency table shows that WorldPop classes tend to be linked to groups of signature types of a similarly degree of urbanity. A WorldPop class 15 is "undefined" due to the lack of building footprints in the area, therefore overlapping a large portion of signatures. The difference between classifications is likely driven by two main aspects - one is the different number of classes. We can see that WorldPop classes tend to cluster within a limited number of signature types and vice versa. The only exception is allocation of signature types into classes 4 and 6, which seems to heavily overlap. That is possibly caused by the second aspect - inclusion of function. Both



**Figure 6.** Contingency table showing frequencies (in %) of WorldPop classes within signature types.

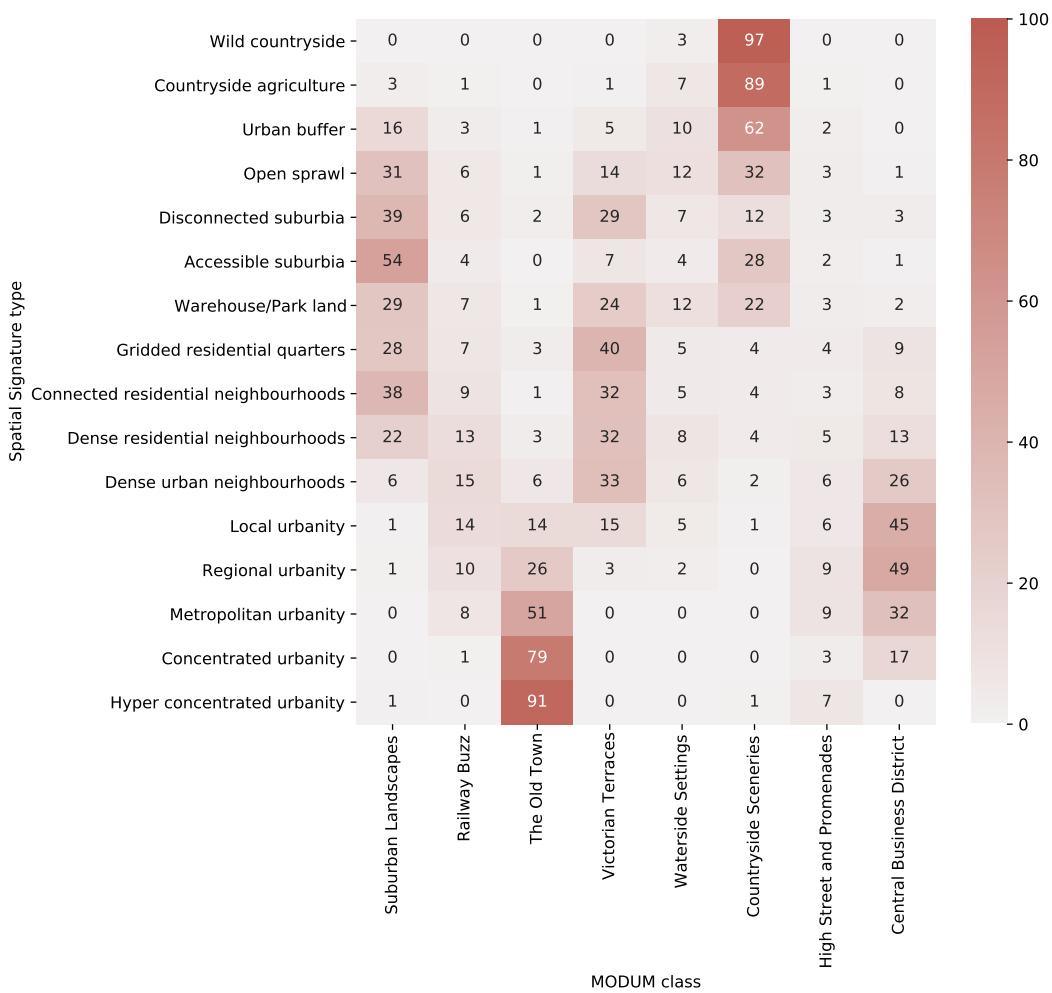
256 classes 4 and 6 tend to be outside of city centres but still within urban areas. While it is the footprint-based form that is driving  
 257 the difference between them, signatures in the same area are often distinguished by function and varies access to amenities and  
 258 services.

## 259 MODUM

260 Multidimensional Open Data Urban Morphology (MODUM) classification describes a typology of neighbourhoods derived  
 261 from 18 indicators capturing built environment as streets, railways or parks, linked to the Census Output Area geometry. The  
 262 classification identifies 8 types of neighbourhoods. Compared to the WorldPop classification, MODUM takes into account  
 263 more features of the built environment than building footprints, which makes it conceptually closer to the spatial signatures.  
 264 However, it is still focusing predominantly on the form component, although there are some indicators that would be classified  
 265 as function within the signatures framework (e.g. population). The MODUM method uses a different way of capturing context  
 266 compared to the signatures, which leads to some classes being determined predominantly by a single character. For example,  
 267 the *Railway Buzz* type forms a narrow strip around the railway network, which is an effect signatures avoid. MODUM typology  
 268 is available only for England and Wales. Therefore the comparison takes into account only ETCs covering the same area.  
 269 The classification is linked to the ETC geometry is based on the proportion (the type covering the largest portion of ETC is  
 270 assigned). The resulting contingency table is shown in Figure 7. There is a significant relationship between two typologies,  
 271  $\chi^2(152, N = 13067584) = 13938867, p < .001$ . The strength of association measured as Cramér's V is 0.300, indicating  
 272 moderate association of very similar levels we have seen above. The contingency table indicates similar relationships, where a  
 273 single MODUM class overlaps a group of signature types. However, the groups tend to be well defined and formed based on  
 274 the similarity of types. Signature types are minimally present in MODUM classes driven by a single character (*Railway Buzz*,  
 275 *Waterside Settings*, *High Street and Promenades*), suggesting the more balanced weight of characters.

## 276 Copernicus Urban Atlas

277 Copernicus Urban Atlas is the least similar of the comparison datasets. It is a high-resolution land use classification of functional  
 278 urban areas derived primarily from Earth Observation data enriched by other reference data as OpenStreetMap or topographic  
 279 maps. Its smallest spatial unit in urban areas is 0.25 ha and 1 ha in rural areas, defined primarily by physical barriers. It  
 280 identifies 27 predefined classes using the supervised method. The majority of urban areas is classified as urban fabric further  
 281 distinguished based on continuity and density resulting in six classes of the urban fabric. The classification does not consider



**Figure 7.** Contingency table showing frequencies (in %) of MODUM classes within signature types.

the type of the pattern or any other aspect. Furthermore, it does not take into account what signatures call *context* as each spatial unit is classified independently, which in some cases leads to the high heterogeneity of classification within a small portion of land. Signatures take a different approach. Consequently, it is expected that the similarity between the two will be limited. Urban Atlas is available only for functional urban areas (FUA), leaving rural areas unclassified. Comparison then applies to FUAs only. The classification is linked to the ETC geometry based on the proportion (the type covering the largest portion of ETC is assigned). The resulting contingency table is shown in Figure 8. There is a significant relationship between two typologies,  $\chi^2(450, N = 8396642) = 5229900, p < .001$ . The strength of association measured as Cramér's V is 0.186, indicating a weak association. The contingency table shows the difference in the aim of spatial signatures and that of Urban Atlas with a majority of signatures being linked to a few of Urban Atlas classes. Within relevant classes, we see a tendency of signature types to cluster within Urban Atlas classes based on the level of urbanity, albeit not as strong as in the previous two cases. The main reason behind such a large difference are the aims of both classifications. While the Copernicus Urban Atlas attempts to capture land cover, resulting in a large number of non-urban classes, spatial signatures are aimed at urban environment with 13 out of 16 classes covering primarily urbanised areas.

## Local Climate Zones

Local climate zones (LCZ) are conceptual classes originally designed to support study of urban climate as temperature. It consists of 17 classes of which 10 can be classified as urban and 7 and natural ones. In the context of Great Britain, the dataset used in this study does not contain 2 of them, *Lightweight low-rise* and *Compact highrise* as they are not present in the British landscape. The datasets produced by<sup>46</sup> released LCZs in a 100 meters grid based on the 2016 data. As the LCZs are remotely sensed in this case, authors report overall average accuracy of 80 %. As a conceptual classification aimed to cover all possible types of primarily urban climate zones globally, LCZs may not be optimal when looking into a single country with specific history of urban development. This is further indicated by classes that are missing. It is therefore likely that large parts of British cities will fall into only a few of LCZ classes, while being represented by a much larger number of signature types.

Signature typology is rasterized and linked to the LCZ grid. The resulting contingency table is shown in Figure 9. There is a significant relationship between two typologies,  $\chi^2(225, N = 16203338) = 18467242, p < .001$ . The strength of association measured as Cramér's V is 0.276, indicating a modest to weak association, close to values we've seen in first two cases. As expected, urban signature types are clustered primarily within *Compact midrise* and *Open lowrise* LCZs, while non-urban signatures mostly fall into the *Low plants* LCZ.

The difference between signatures and LCZs can be accounted to two aspects. One, as we've seen before is the inclusion of function in spatial signatures, differentiating e.g. LCZ's *Open lowrise* into many signature types. The other is data-driven nature of signatures compared to conceptual LCZs, where differences in signature types are below the resolution capability of simple matrix composed of density and compactness levels. On the other, it is encouraging to see that most of signature types fall predominantly in a single LCZ class, suggesting that while both classifications are built differently, they are able to capture similar large-scale patterns in cities.

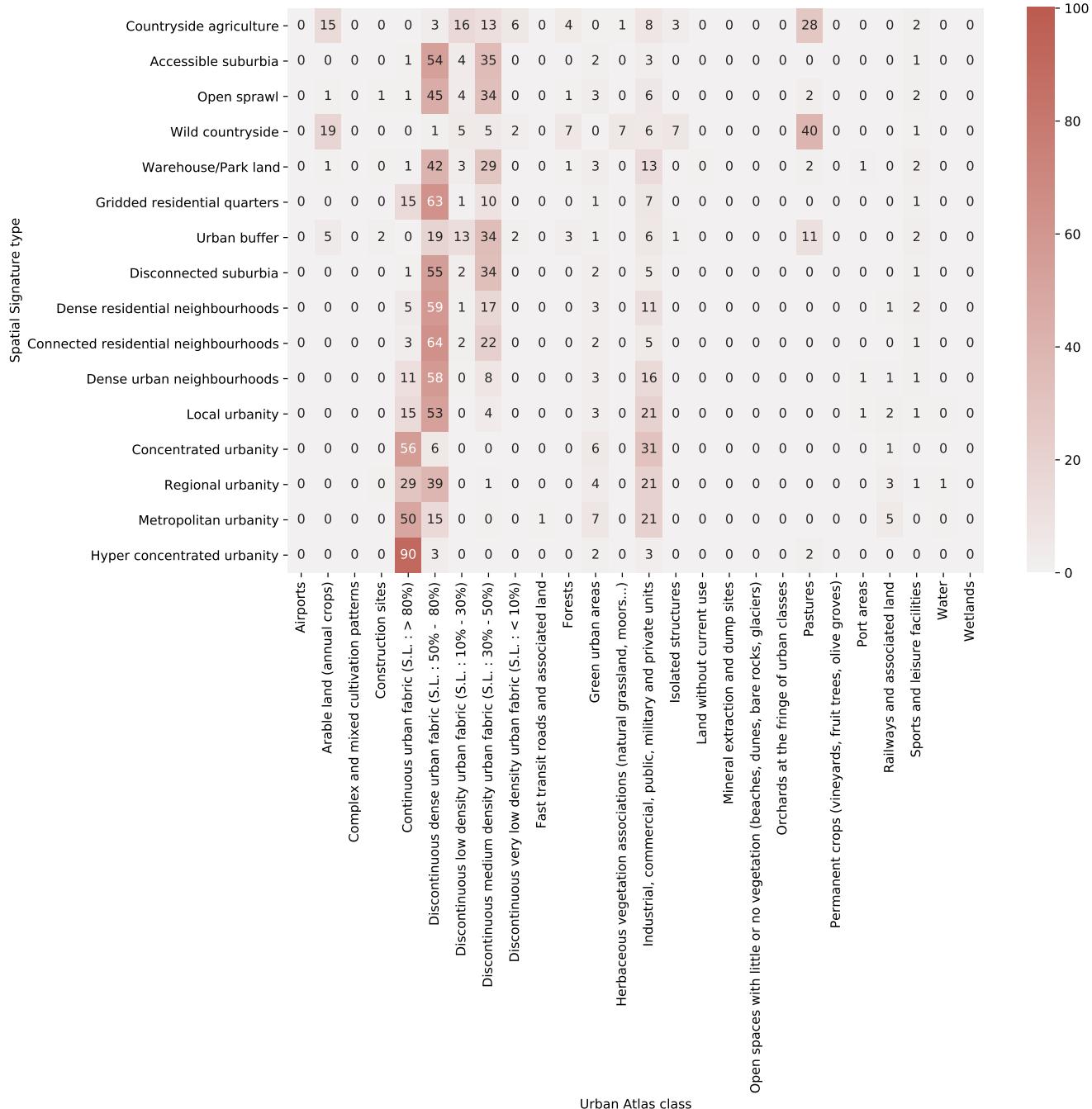
## Summary

None of the comparisons shows more than a moderate association, but since none of the comparison datasets is aiming to capture the same conceptualization of space as spatial signatures do, such a result is expected. The moderate association with both WorldPop settlements patterns and MODUM is reassuring as both are conceptually closer to signatures than the Urban Atlas (especially in their unsupervised design). Urban Atlas, though very different in its aims and methods, still shows a measurable association, which we interpret as sign that the key structural aspects forming cities are captured by both. The comparison exercise suggests that general patterns forming cities are shared among signatures and existing typologies. Signature types tend to form groups when we look at their relation to comparison classes and it is not uncommon that a single signature type is present in multiple groups linked to different classes. However, all these groups tend to be formed based on the similarity and illustrate the granularity of the presented classification compared to existing datasets, allowing us to distinguish, for example, five types of signature types forming town an city centres.

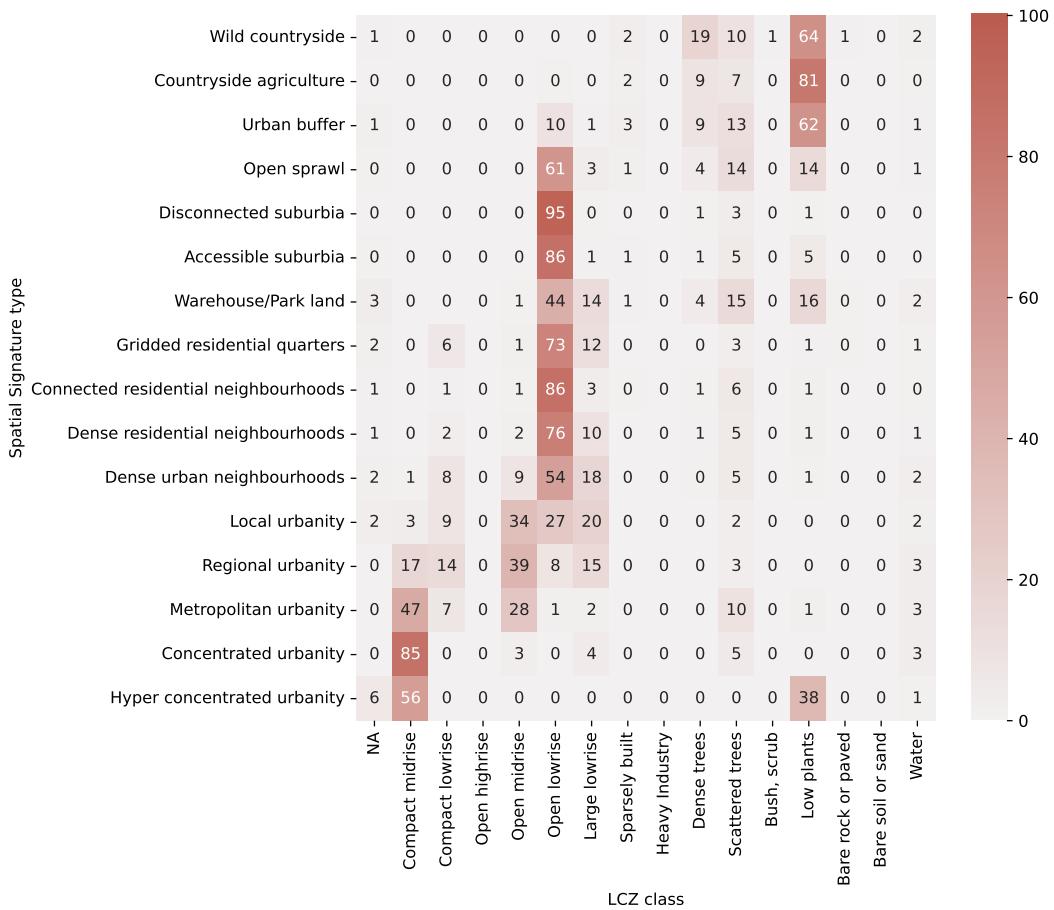
## Usage Notes

The released data product follows widespread standards for geographic data storage and should be easy to integrate with other data and methods by researchers wanting to reuse it. However, due to the density of signature geometry (resulting from the detailed ETCs), it may be needed to simplify the geometry for a smoother interactive experience on machines with limited resources.

Replication of the analysis optimally requires at least a single computational node with a large amount of RAM (+100GB) due to the size of the input data and detail on which signature characterization is computed. It is also recommended to revisit the



**Figure 8.** Contingency table showing frequencies (in %) of Urban Atlas classes within signature types.



**Figure 9.** Contingency table showing frequencies (in %) of Local Climate Zones within signature types.

333 state of the development of related software packages, notably `momepy`<sup>48</sup>, `libpsyal`<sup>49</sup>, `tobler`<sup>37</sup> and `dask-geopandas`  
334 as they may soon offer more efficient drop-in replacements of the custom code used to produce this dataset.

## 335 **Code availability**

336 The source code used to produce this dataset is openly available in a GitHub repository at  
337 [https://github.com/urbangrammarai/spatial\\_signatures](https://github.com/urbangrammarai/spatial_signatures) and in the form of a website on <https://urbangrammarai.xyz>. Code is  
338 organized in a series of Jupyter notebooks and have been executed within the `darribas:gds_dev:6.1`<sup>50</sup> Docker container,  
339 unless specified otherwise in the individual notebooks.

## 340 **References**

- 341 1. Arribas-Bel, D. & Fleischmann, M. Spatial Signatures - Understanding (urban) spaces through form and function (2021).  
342 Mimeo.
- 343 2. Fleischmann, M. & Arribas-Bel, D. Classifying urban form at national scale - The British morphosignatures (2021).  
344 Proceedings of XXVIII International Seminar on Urban Form.
- 345 3. Arribas-Bel, D., Green, M., Rowe, F. & Singleton, A. Open Data Products: a framework for creating valuable analysis-ready  
346 data. *J. Geogr. Syst.* (*forthcoming*).
- 347 4. HM Government. Levelling Up the United Kingdom. <https://www.gov.uk/government/publications/levelling-up-the-united-kingdom> (2022). London: The Stationery Office.
- 348 5. Stewart, I. D. & Oke, T. R. Local climate zones for urban temperature studies. *Bull. Am. Meteorol. Soc.* **93**, 1879–1900  
349 (2012).
- 350 6. Angel, S., Arango Franco, S., Liu, Y. & Blei, A. M. The shape compactness of urban footprints. *Prog. Plan.* **139**, 100429,  
351 [10/gg638j](https://doi.org/10/gg638j) (2020).
- 352 7. Alexiou, A., Singleton, A. & Longley, P. A. A Classification of Multidimensional Open Data for Urban Morphology. *Built  
353 Environ.* **42**, 382–395, [10/gddwsn](https://doi.org/10/gddwsn) (2016).
- 354 8. Taubenböck, H. *et al.* A new ranking of the world's largest cities—do administrative units obscure morphological realities?  
355 *Remote. Sens. Environ.* **232**, 111353 (2019).
- 356 9. Brodsky, I. H3: Uber's hexagonal hierarchical spatial index. Available from Uber Eng. website: [https://eng.uber.com/h3/\[22 June 2019\]](https://eng.uber.com/h3/[22 June 2019]) (2018).
- 357 10. Jochem, W. C. & Tatem, A. J. Tools for mapping multi-scale settlement patterns of building footprints: An introduction to  
358 the R package `foot`. *PLoS one* **16**, e0247535, [10/gh7sjr](https://doi.org/10/gh7sjr) (2021).
- 359 11. Araldi, A. & Fusco, G. From the street to the metropolitan region: Pedestrian perspective in urban fabric analysis:. *Environ.  
360 Plan. B: Urban Anal. City Sci.* **46**, 1243–1263, [10.1177/2399808319832612](https://doi.org/10.1177/2399808319832612) (2019).
- 361 12. Gil, J., Montenegro, N., Beirão, J. N. & Duarte, J. P. On the Discovery of Urban Typologies: Data Mining the Multi-  
362 dimensional Character of Neighbourhoods. *Urban Morphol.* **16**, 27–40 (2012).
- 363 13. Hamaina, R., Leduc, T. & Moreau, G. Towards Urban Fabrics Characterization Based on Buildings Footprints. In *Bridging  
364 the Geographic Information Sciences*, vol. 2, 327–346, [10.1007/978-3-642-29063-3\\_18](https://doi.org/10.1007/978-3-642-29063-3_18) (Springer, Berlin, Heidelberg,  
365 Berlin, Heidelberg, 2012).
- 366 14. Bobkova, E., Berghauser Pont, M. & Marcus, L. Towards analytical typologies of plot systems: Quantitative profile of five  
367 European cities. *Environ. Plan. B: Urban Anal. City Sci.* [239980831988090](https://doi.org/10/ggbgsm), [10/ggbgsm](https://doi.org/10/ggbgsm) (2019).
- 368 15. Kropf, K. Plots, property and behaviour. *Urban Morphol.* **22**, 5–14 (2018).
- 369 16. Fleischmann, M., Feliciotti, A., Romice, O. & Porta, S. Morphological tessellation as a way of partitioning space:  
370 Improving consistency in urban morphology at the plot scale. *Comput. Environ. Urban Syst.* **80**, 101441, [10.1016/j.compenvurbsys.2019.101441](https://doi.org/10.1016/j.compenvurbsys.2019.101441) (2020).
- 371 17. Ordnance Survey. OS Open Roads (2020).
- 372 18. Ordnance Survey. OS OpenMap - Local (2020).
- 373 19. Ordnance Survey. OS Open Rivers (2020).
- 374 20. Ordnance Survey. Strategi (2016).

- 378 21. Dibble, J. *et al.* On the origin of spaces: Morphometric foundations of urban form evolution. *Environ. Plan. B: Urban*
- 379 Anal. City Sci. **46**, 707–730 (2019).
- 380 22. Fleischmann, M., Romice, O. & Porta, S. Measuring urban form: Overcoming terminological inconsistencies for a
- 381 quantitative and comprehensive morphologic analysis of cities. *Environ. Plan. B: Urban Anal. City Sci.* 239980832091044,
- 382 [10/ggngw6](https://doi.org/10/ggngw6) (2020).
- 383 23. Fleischmann, M., Feliciotti, A., Romice, O. & Porta, S. Methodological foundation of a numerical taxonomy of urban
- 384 form. *Environ. Plan. B: Urban Anal. City Sci.* 23998083211059835 (2021).
- 385 24. Sneath, P. H., Sokal, R. R. *et al.* *Numerical taxonomy. The principles and practice of numerical classification.* (Freeman,
- 386 San Francisco, 1973).
- 387 25. Hallowell, G. D. & Baran, P. K. Suburban change: A time series approach to measuring form and spatial configuration.
- 388 *The J. Space Syntax* **4**, 74–91 (2013).
- 389 26. Vanderhaegen, S. & Canters, F. Mapping urban form and function at city block level using spatial metrics. *Landsc. Urban*
- 390 *Plan.* **167**, 399–409, [10.1016/j.landurbplan.2017.05.023](https://doi.org/10.1016/j.landurbplan.2017.05.023) (2017).
- 391 27. Schirmer, P. M. & Axhausen, K. W. A multiscale classification of urban morphology. *J. Transp. Land Use* **9**, 101–130,
- 392 [10.5198/jtlu.2015.667](https://doi.org/10.5198/jtlu.2015.667) (2015).
- 393 28. Steiniger, S., Lange, T., Burghardt, D. & Weibel, R. An Approach for the Classification of Urban Building Structures
- 394 Based on Discriminant Analysis Techniques. *Transactions GIS* **12**, 31–59, [10.1111/j.1467-9671.2008.01085.x](https://doi.org/10.1111/j.1467-9671.2008.01085.x) (2008).
- 395 29. Basaraner, M. & Cetinkaya, S. Performance of shape indices and classification schemes for characterising perceptual shape
- 396 complexity of building footprints in GIS. *Int. J. Geogr. Inf. Sci.* **31**, 1952–1977, [10.1080/13658816.2017.1346257](https://doi.org/10.1080/13658816.2017.1346257) (2017).
- 397 30. Hamaina, R., Leduc, T. & Moreau, G. A New Method to Characterize Density Adapted to a Coarse City Model.
- 398 In *OpenStreetMap in GIScience*, 249–263, [10.1007/978-3-642-31833-7\\_16](https://doi.org/10.1007/978-3-642-31833-7_16) (Springer International Publishing, Berlin,
- 399 Heidelberg, 2013).
- 400 31. Hijazi, I. *et al.* Measuring the homogeneity of urban fabric using 2D geometry data. *Environ. Plan. B: Plan. Des.* 1–25,
- 401 [10.1177/0265813516659070](https://doi.org/10.1177/0265813516659070) (2016).
- 402 32. Boeing, G. A multi-scale analysis of 27,000 urban street networks: Every US city, town, urbanized area, and Zillow
- 403 neighborhood. *Environ. Plan. B: Urban Anal. City Sci.* **219**, 239980831878459, [10.1177/239980831878459](https://doi.org/10.1177/239980831878459) (2018).
- 404 33. Caruso, G., Hilal, M. & Thomas, I. Measuring urban forms from inter-building distances: Combining MST graphs with a
- 405 Local Index of Spatial Association. *Landsc. Urban Plan.* **163**, 80–89, [10.1016/j.landurbplan.2017.03.003](https://doi.org/10.1016/j.landurbplan.2017.03.003) (2017).
- 406 34. Feliciotti, A. *RESILIENCE AND URBAN DESIGN: A SYSTEMS APPROACH TO THE STUDY OF RESILIENCE IN*
- 407 *URBAN FORM*. Ph.D. thesis, University of Strathclyde, Glasgow (2018).
- 408 35. Lowry, J. H. & Lowry, M. B. Comparing spatial metrics that quantify urban form. *Comput. Environ. Urban Syst.* **44**,
- 409 59–67, [10.1016/j.compenvurbsys.2013.11.005](https://doi.org/10.1016/j.compenvurbsys.2013.11.005) (2014).
- 410 36. Porta, S., Crucitti, P. & Latora, V. The network analysis of urban streets: A primal approach. *Environ. Plan. B: Plan. Des.*
- 411 **33**, 705–725, [10.1068/b32045](https://doi.org/10.1068/b32045) (2006).
- 412 37. eli knaap *et al.* pysal/tobler: Release v0.8.2, [10.5281/zenodo.5047613](https://doi.org/10.5281/zenodo.5047613) (2021).
- 413 38. Webber, R. & Burrows, R. *The predictive postcode: the geodemographic classification of British society* (Sage, 2018).
- 414 39. Lage, J. P., Assunção, R. M. & Reis, E. A. A minimal spanning tree algorithm applied to spatial cluster analysis. *Electron.*
- 415 *Notes Discret. Math.* **7**, 162–165 (2001).
- 416 40. Baçao, F., Lobo, V. & Painho, M. The self-organizing map, the geo-som, and relevant variants for geosciences. *Comput. &*
- 417 *geosciences* **31**, 155–163 (2005).
- 418 41. Wolf, L. J. Spatially–encouraged spectral clustering: a technique for blending map typologies and regionalization. *Int. J.*
- 419 *Geogr. Inf. Sci.* **35**, 2356–2373 (2021).
- 420 42. Schonlau, M. The clustergram: A graph for visualizing hierarchical and nonhierarchical cluster analyses. *The Stata J.* **2**,
- 421 391–402, [10.ghh97z](https://doi.org/10.ghh97z) (2002).
- 422 43. Caliński, T. & Harabasz, J. A dendrite method for cluster analysis. *Commun. Stat.* **3**, 1–27, [10.1080/03610927408827101](https://doi.org/10.1080/03610927408827101)
- 423 (1974). <https://www.tandfonline.com/doi/pdf/10.1080/03610927408827101>.
- 424 44. Davies, D. L. & Bouldin, D. W. A cluster separation measure. *IEEE transactions on pattern analysis machine intelligence*
- 425 224–227 (1979).

- 426 45. European environment agency (EEA). Urban Atlas (2018).
- 427 46. Demuzere, M., Bechtel, B., Middel, A. & Mills, G. Mapping europe into local climate zones. *PloS one* **14**, e0214474  
428 (2019).
- 429 47. Cramér, H. *Mathematical Methods of Statistics (PMS-9), Volume 9* (Princeton university press, 2016).
- 430 48. Fleischmann, M. momepy: Urban morphology measuring toolkit. *J. Open Source Softw.* **4**, 1807, [10.21105/joss.01807](https://doi.org/10.21105/joss.01807)  
431 (2019).
- 432 49. Rey, S. J. *et al.* The pysal ecosystem: Philosophy and implementation. *Geogr. Analysis* (2021).
- 433 50. Arribas-Bel, D. gds\_env: A containerised platform for geographic data science, [10.5281/zenodo.4642516](https://doi.org/10.5281/zenodo.4642516).

## 434 Acknowledgements

435 (not compulsory)

436 M.F. and D.A. kindly acknowledge funding by the UK's Economic and Social Research Council through the project  
437 "Learning an urban grammar from satellite data through AI", project reference ES/T005238/1.

## 438 Author contributions statement

439 M.F. and D.A. designed the method, M.F. conducted the experiments, M.F. and D.A. analysed the results. M.F. and D.A. wrote  
440 and reviewed the manuscript.

## 441 Competing interests

442 The authors declare no competing interests.

## 443 Figures & Tables