

1 Characterisation of space in Great Britain using the 2 Spatial Signatures model

3 Martin Fleischmann^{1, *} and Daniel Arribas-Bel¹

4 ¹Geographic Data Science Lab, Department of Geography and Planning, University of Liverpool, Roxby Building ,
5 74 Bedford St S , Liverpool , L69 7ZT, United Kingdom

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

7 ABSTRACT

8 This is a manuscript template for Data Descriptor submissions to *Scientific Data* (<http://www.nature.com/scientificdata>). The abstract must be no longer than 170 words, and should succinctly describe the study, the assay(s) performed, the resulting data, and the reuse potential, but should not make any claims regarding new scientific findings. No references are allowed in this section.

9 Please note: Abbreviations should be introduced at the first mention in the main text – no abbreviations lists or tables should be
10 included. Structure of the main text is provided below.

11 Background & Summary

12 (700 words maximum) An overview of the study design, the assay(s) performed, and the created data, including any background
13 information needed to put this study in the context of previous work and the literature. The section should also briefly outline
14 the broader goals that motivated the creation of this dataset and the potential reuse value. We also encourage authors to include
15 a figure that provides a schematic overview of the study and assay(s) design. The Background & Summary should not include
16 subheadings. This section and the other main body sections of the manuscript should include citations to the literature as
17 needed.

18 The best paper ever is¹.

19 Methods

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

24 Spatial unit

25 The first major methodological decision needs to be taken on the definition of the spatial unit. As mentioned, it needs to
26 reflect space in a granular manner and we argue that it should fulfil three conditions. First, it should be *indivisible*, meaning
27 that when such a unit would be subdivided into smaller parts, none of them would be enough to capture the nature of spatial
28 signature. Second, it needs to be *internally consistent* - it should always reflect only a single signature type. Last, it should be
29 geographically *exhaustive*, covering entirety of the study area.

30 Spatial units used in literature can be split into three groups. One is using administrative boundaries like city regions, wards
31 or census output areas, that are convenient to obtain and can be easily linked to auxiliary data. However, those rarely reflect
32 the morphological composition of urban space and in some cases may even “obscure morphologic reality” REF Taubenbock
33 2019. At the same time, most of them are divisible and larger units are not always internally consistent. Another group is based
34 on arbitrary uniform grids linked either to spatial indexing method like H3 REF or OS National Grid REF, or to auxiliary data
35 of remote sensing or other origins like a WorldPop grid REF. The issue is that grids cannot be considered internally consistent
36 as they have no relation to the real-life spatial pattern. Finally, urban morphology tends to use morphological elements as street
37 segments REF, blocks REF buildings or plots as a unit of analysis. Some of those could be seen as indivisible and internally
38 consistent but since they are largely based on built-up fabric, they are not exhaustive. When there is no building or street, there
39 is no spatial unit to work with. Plots could be theoretically considered as exhaustive, consistent and indivisible but there is no
40 accepted conceptual definition and unified geometric representation (REF Kropf).

We are, therefore, proposing an application of an alternative spatial unit called *enclosed tessellation cell* (ETC), defined as:
A characterisation of space based on form and function designed to understand urban environments

ETC follows the morphological tradition in a sense that it is based on the physical elements of an environment but overcomes the drawbacks of conventionally used units. Its geometry is generated in three steps illustrated on a Figure . First, a set of features representing physical barriers subdividing space, in our case composed of street network, railways, rivers and a coastline, is combined together, generating a layer of boundaries. These then partition space into smaller enclosed geometries called *enclosures*, which can be very granular or very coarse depending on the geographic context. In dense city centres where a single enclosure represent a single block is a high frequency of small enclosures, while in the countryside, we can observe very few large enclosures as their delimiters are far away from each other. Enclosures are then combined with building footprints, posing as anchors in the space and are subdivided into enclosed tessellation cells using the morphological tessellation algorithm REF, a polygon-based adaptation of Voronoi tessellation. Resulting geometries are indivisible as they contain, at most, a single anchor building, internally consistent due to their granularity and link to morphological elements composing urban fabric, and geographically exhaustive as they cover entire area limited by specified boundaries.

In the case of classification of Great Britain, street networks are extracted from OS Open Roads datasets (REF) representing simplified road centrelines cleaned of road segments under the ground. Railways are retrieved from OS OpenMap - Local ("RailwayTrack" layer) which captures surface railway tracks. Rivers are extracted from OS OpenRivers (REF) representing river network of GB as centrelines, and a coastline is retrieved from OS Strategi® (2016) REF, capturing coastline as a continuous line geometry. Building geometry is extracted, again, from OS OpenMap - Local ("Building" layer) and represents generalised building footprint polygons. 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).

Characterisation of space

Spatial signatures are capturing the character of the built and unbuilt environment based on two components - form and function. Each of them is quantified on the level of individual ETCs using methods appropriate for the specific datasets. While form component is described using urban morphometrics (i.e. quantitative analysis of urban form), function is a composite of a variety of data inputs outlined in detail below.

Form

Morphometric characterisation of urban form is based on the numerical description of four elements capturing the built environment - buildings, streets, ETCs, and enclosures, and reflects their patterns based on six categories of characters - dimensions, shapes, spatial distribution, intensity, connectivity and diversity. Each element is considered across different scales, from the measurements of individual geometries, relations of neighbouring geometries to graph-based analysis of street network. The combination of elements, categories and scales results in a set of 59 individual morphometric characters listed in the table XXX.

However, measuring individual characters is not enough to understand the predominant spatial patterns as for some types of urban form is typical high heterogeneity. That means that using, for example, areas of building footprints would in most cases result in largely discontinuous clusters. We are, therefore, representing each of the morphometric characters using three proxy variables reflecting statistical distributions of measured data within a spatial context of each ETC. Context is defined as 10th order of contiguity based on the mesh composed of contiguous ETCs. Furthermore, each value is weighted by the inverse distance between so called poles of inaccessibility (defined as a centre of a maximum inscribed circle) of each ETC. Three proxy variables then capture the first, the second and the third quartile of the resulting weighted distribution. Such a characterisation is able to capture the contextual tendency of each morphometric character and hence identify contiguous clusters in both homogenous and heterogenous urban tissues.

Function

Characterisation of the function component uses a different approach. While data describing urban form are not generally available in a processed format, hence we have to employ morphometric approaches, different aspects of function are often available as open data products. Therefore, the main goal of our characterisation of ETCs based on function is to develop appropriate transfer methods to link data published as grids or linked to administrative boundaries to enclosed tessellation.

In this work we are using five different transfer methods:

- Areal interpolation
- Building-based dasymetric areal interpolation
- Network-constrained accessibility
- Euclidean accessibility

- Zonal statistics

Areal interpolation is used when the functional data covers the entirety of space in a form of polygon geometry and when there is no assumption that the phenomena it captures is linked directly to the human population, for example land cover data. When there is an assumption of relation to the population, building-based dasymetric areal interpolation is used instead. The main difference is that instead of ETC polygons, building footprint polygons linked to individual ETCs are used as a target of interpolation. That ensures that data like population estimates are linked to ETCs proportionally to their ability to provide accommodation, rather than by their area. Network-constrained accessibility is used when the input data represent points of interest like locations of supermarkets. Points are then snapped to the nearest node on the street network and linked to the ETCs as a number of observations within 15 minutes walking distance (1200m on the street network) and a distance to the nearest point. In some cases, Euclidean (as-crow-flies) accessibility is measured instead to accommodate for phenomena that are often reached outside of drivable network like water bodies. Final method, zonal statistics, is used to transfer data originally stored in a raster format to ETC polygons as a mean value of raster pixels intersecting each polygon geometry. Finally, characters based on interpolation and zonal statistics are expressed using their contextual versions following the method used for form characters to, again, reflect the pattern of measured values. The selection of datasets and the chosen transfer method are listed in the table XXX.

Cluster analysis

When combined, contextual proxies of form and function characters (or characters themselves when they are reflecting the context by definition) compose a dataset describing each ETC by 331 variables (177 for form and 154 for function.) We treat all of them equally (there is no weighting involved), standardize each variable applying Z-score normalization and use them as an input for K-Means cluster analysis.

Due to the nature of the selected K-Means clustering the step preceeding the final analysis is the selection of an optimal number of clusters. We use exploratory clustergram method (REF) reflecting the behaviour of different options, the relationship between clustering solutions regarding the allocation of individual observations to classes and the separation between the clusters within each tested solution. Clustergram is further accompanied by measures of internal validation measures - the Silhouette score diagram, Calinski-Harabasz index and Davies-Bouldin index. The optimal number of classes is selected based on the interpretation of clustergram supported by additional measures aiming at a balance between cluster separation and an appropriate detail of resulting classification.

The results of the top level clustering capture the first layer of a national signature classification. However, since the classified ETCs cover entirety of space from vast natural open spaces to dense city centres, it may result in only a few class representing urban areas. While that is caused by the variable heterogeneity of our dataset in combination with K-Means clustering, the measured characters have the ability to further distinguish sub-classes of already identified clusters. As spatial signatures are focused on urban environment, we further subdivide those clusters covering substantial portion of urban areas using another iteration of K-Means clustering. Resulting classification then provide two hierarchical levels capturing the typology of spatial signatures with a detailed focus on urban development.

Finally, individual spatial signature geometries are generated as a combination of adjacent ETCs belonging to the same signature class.

Data Records

The Data Records section should be used to explain each data record associated with this work, including the repository where this information is stored, and to provide an overview of the data files and their formats. Each external data record should be cited numerically in the text of this section, for example², and included in the main reference list as described below. A data citation should also be placed in the subsection of the Methods containing the data-collection or analytical procedure(s) used to derive the corresponding record. Providing a direct link to the dataset may also be helpful to readers (<https://doi.org/10.6084/m9.figshare.853801>).

Tables should be used to support the data records, and should clearly indicate the samples and subjects (study inputs), their provenance, and the experimental manipulations performed on each (please see 'Tables' below). They should also specify the data output resulting from each data-collection or analytical step, should these form part of the archived record.

Technical Validation

This section presents any experiments or analyses that are needed to support the technical quality of the dataset. This section may be supported by figures and tables, as needed. This is a required section; authors must present information justifying the reliability of their data.

Spatial signatures are unique as a classification method, limiting the potential validation methods to only indirect methods using ancillary datasets capturing conceptually similar aspects of environment. We compare the signatures with three 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 measurable level of association between the two:

- WorldPop settlement patterns of building footprints (2021)
- Classification of Multidimensional Open Data of Urban Morphology (MODUM)(2015)
- Copernicus Urban Atlas (2018)

Validation method

All datasets, spatial signatures as well as those selected as validation contain categorical classification of space linked to their unique geometry. The first task, to make each pair comparable is to transfer data to the same geometry. That can be interpolation of one set of polygon-based data to another (input to ETCs) or converting spatial signatures to the raster representation matching an input raster as the latter is computationally more feasible 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 validation classes. We use contingency tables and a Pearson's χ^2 test to determine whether the frequencies of observed (signature types) and expected (validation types) labels significantly differ in one or more categories. Furthermore, we use Cramér's V statistics to assess the strength of an association (assuming the Pearson's χ^2 test rejected the hypothesis of independence).

WorldPop settlement patterns of building footprints

WorldPop settlement patterns of building footprints aim to derive a typology of morphological patterns based on the gridded approach (spatial unit is a grid of a size 100x100m per cell) and building footprints. Authors measure 6 morphometric characters linked to the grid cells and use them as an input of unsupervised clustering leading to a 6 class typology (Figure XXX). As the classification is dependent on the building footprint data, grid cells that do not contain any information on building-based pattern are treated as missing in the final data product. For the validation of spatial signatures, 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 inclusion of function in spatial signatures, higher granularity of both initial spatial unit and the resulting classification (6 vs 19 classes).

Signature typology is rasterized and linked to the WorldPop grid. The resulting contingency table is shown on Figure XXX. 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.

MODUM

Multidimensional Open Data Urban Morphology (MODUM) classification describes a typology of neighbourhoods derived from 18 indicators capturing built environment as streets, railways or parks, linked to the Census Output Area geometry. The classification identifies 8 types of neighbourhoods illustrated on figure XXX. Compared to the WorldPop classification, MODUM takes into the account more features of built environment than building footprints, which makes it conceptually closer to the spatial signatures. However, it is still focusing predominantly on the form component, although there are some indicators that would be classified as function within the signatures framework (e.g. population). The MODUM method uses a different way of capturing the context compared to the signatures, which leads to some classes being determined predominantly by a single character. For example, *Railway Buzz* type forms a narrow strip around the railway network, which is an effect signatures are trying to avoid. MODUM typology is available only for England and Wales, therefore the validation takes into the account only ETCs covering the same area. 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 on Figure XXX. There is a significant relationship between two typologies, $\chi^2(152, N = 13067584) = 13938867, p < .001$. The strength of association measured as Cramér's V is 0.300, indicating moderate association of a very similar levels we have seen above.

Copernicus Urban Atlas

Copernicus Urban Atlas is the least similar of the validation datasets. It is a high-resolution land use classification of functional urban areas derived primarily from Earth Observation data enriched by other reference data as OpenStreetMap or topographic maps. Its smallest spatial unit in urban areas is 0.25 ha and 1 ha in rural areas, defined primarily by physical barriers. The Urban Atlas classification, which identifies 27 classes, is illustrated on figure XXX. The majority of urban areas is classified as urban fabric further distinguished based on continuity and density resulting in 6 classes of urban fabric. The classification does not consider the type of the pattern or any other aspect. Furthermore, it does not take into account what signatures call the *context* as each spatial unit is classified independently, which in some cases leads to a high heterogeneity of classification

within a small portion of land. Signatures take a different approach, therefore 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. Validation then works with 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 on Figure XXX. 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 weak association.

Summary

Usage Notes

Released dataset is following the widespread standards for geographic data storage and should not pose a challenge for researches wanting to reuse it. Due to the density of signature geometry (resulting from the detailed ETCs), it may be needed to simplify the geometry for smoother interactive experience on weaker machines.

Replication of the analysis optimally requires at least a single computational node with large amount of RAM (100 GB+) due to the size of the input data and detail on which signature characterization is computed. It is also recommended to revisit the state of the development of related software packages, notably momepy, libpysal, tobler and dask-geopandas as they may in the near future offer more efficient drop-in replacements of the custom code used to produce this dataset.

Code availability

The source code used to produce this dataset is openly available in a GitHub repository at https://github.com/urbangrammarai/spatial_signatures and in the form of a website on <https://urbangrammarai.github.io>. Code is organized in a series of Jupyter notebooks and have been executed within the darribas:gds_env Docker container, unless specified otherwise in the individual notebooks. The specific version of the container is listed on top of each notebook.

References

1. Singleton, A. & Arribas-Bel, D. Geographic data science. *Geogr. Analysis* **53**, 61–75 (2021).
2. Hao, Z., AghaKouchak, A., Nakhjiri, N. & Farahmand, A. Global integrated drought monitoring and prediction system (GIDMaPS) data sets. *figshare* <https://doi.org/10.6084/m9.figshare.853801> (2014).
3. Kaufman, D. *et al.* A global database of holocene paleotemperature records. *Sci. Data* **7**, 115, <https://doi.org/10.1038/s41597-020-0445-3> (2020).
4. Figueredo, A. J. & Wolf, P. S. A. Assortative pairing and life history strategy – a cross-cultural study. *Hum. Nat.* **20**, 317–330, <https://doi.org/10.1007/s12110-009-9068-2> (2009).
5. Babichev, S. A., Ries, J. & Lvovsky, A. I. Quantum scissors: teleportation of single-mode optical states by means of a nonlocal single photon. Preprint at <https://arxiv.org/abs/quant-ph/0208066> (2002).
6. Behringer, R. *Manipulating the mouse embryo: a laboratory manual* (Cold Spring Harbor Laboratory Press, New York, 2014).

LaTeX formats citations and references automatically using the bibliography records in your .bib file, which you can edit via the project menu. Use the cite command for an inline citation, e.g.^{3–6}. For data citations of datasets uploaded to e.g. *figshare*, please use the howpublished option in the bib entry to specify the platform and the link, as in the Hao:gidsmaps:2014 example in the sample bibliography file. For journal articles, DOIs should be included for works in press that do not yet have volume or page numbers. For other journal articles, DOIs should be included uniformly for all articles or not at all. We recommend that you encode all DOIs in your bibtex database as full URLs, e.g. <https://doi.org/10.1007/s12110-009-9068-2>.

Acknowledgements

(not compulsory)

Acknowledgements should be brief, and should not include thanks to anonymous referees and editors, or effusive comments. Grant or contribution numbers may be acknowledged.

Author contributions statement

Must include all authors, identified by initials, for example: A.A. conceived the experiment(s), A.A. and B.A. conducted the experiment(s), C.A. and D.A. analysed the results. All authors reviewed the manuscript.

237 **Competing interests**

238 The authors declare no competing interests.

239 **Figures & Tables**

240 Figures, tables, and their legends, should be included at the end of the document. Figures and tables can be referenced in \LaTeX
241 using the ref command, e.g. Figure 1 and Table 1.

242 Authors are encouraged to provide one or more tables that provide basic information on the main ‘inputs’ to the study (e.g.
243 samples, participants, or information sources) and the main data outputs of the study. Tables in the manuscript should generally
244 not be used to present primary data (i.e. measurements). Tables containing primary data should be submitted to an appropriate
245 data repository.

246 Tables may be provided within the \LaTeX document or as separate files (tab-delimited text or Excel files). Legends, where
247 needed, should be included here. Generally, a Data Descriptor should have fewer than ten Tables, but more may be allowed
248 when needed. Tables may be of any size, but only Tables which fit onto a single printed page will be included in the PDF
249 version of the article (up to a maximum of three).

250 Due to typesetting constraints, tables that do not fit onto a single A4 page cannot be included in the PDF version of the
251 article and will be made available in the online version only. Any such tables must be labelled in the text as ‘Online-only’ tables
252 and numbered separately from the main table list e.g. ‘Table 1, Table 2, Online-only Table 1’ etc.



Figure 1. Legend (350 words max). Example legend text.

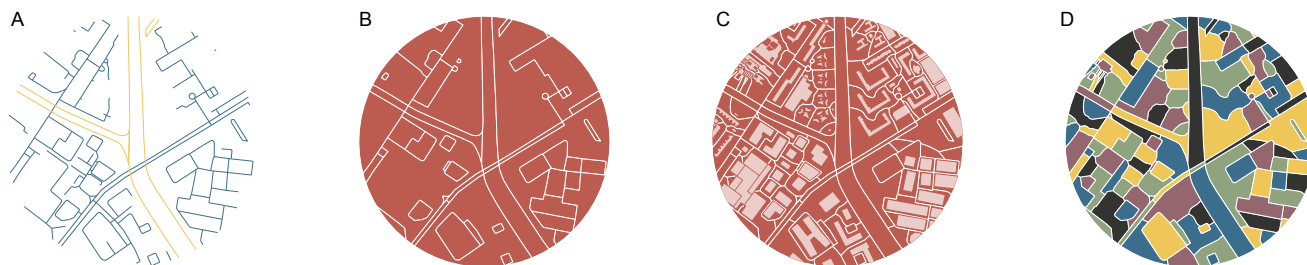


Figure 2. Diagram illustrating the sequential steps leading to the delineation of enclosed tessellation. From a series of enclosing components, where blue are streets and yellow river banks (A), to enclosures (B), incorporation of buildings as anchors (C) to final tessellation cells (D).

Condition	n	p
A	5	0.1
B	10	0.01

Table 1. Legend (350 words max). Example legend text.