

# UrbanDataLayer: A Unified Data Pipeline for Urban Science

Yiheng Wang<sup>1</sup>, Tianyu Wang<sup>1</sup>, Yuying Zhang<sup>1</sup>  
Hongji Zhang<sup>1</sup>, Haoyu Zheng<sup>1</sup>, Guanjie Zheng<sup>\*1</sup>, Linghe Kong<sup>\*1</sup>  
<sup>1</sup> Shanghai Jiao Tong University, Shanghai, China  
{yhwang0828, wty500, shjtzyy01, zhanghongji, langanzheng, gjzheng, linghe.kong}@sjtu.edu.cn

## 1 Detailed Introduction of Urban Data Layer Components

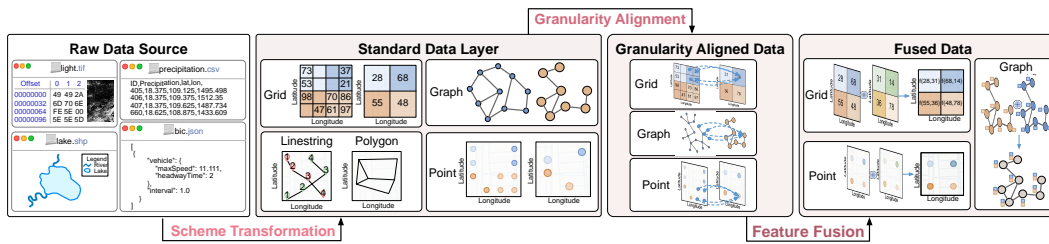


Figure 1: The components of UrbanDataLayer. The words in red are the data processing steps.

The overall procedure of UDL data processing is shown in Fig. 1. It can be viewed as four stages of data wrappers (in black bold characters) divided by three data processing steps (in red characters).

### 1.1 Four Data Wrappers

The four data wrappers contain a series of data intermediate states from raw data to fused data which can be directly used by the models.

- (1) **Raw data source.** The raw data refers to unprocessed source data obtained directly from sensors or geographically relevant data obtained from the network such as raster data. These data is usually stored in various forms of structured files such as JSON, CSV and several image files. Raw data presents different structures and must undergo data pre-processing operations before being fed into the model.
- (2) **Standardized data layer.** The standardized data layer is a uniformly defined data structure and consists of different types of urban data. It contains both data spatio-temporal information and the ready-to-use data itself. We define five types of data layer, including grid data, graph data, point data, linestring data and polygon data.
  - **Grid data:** Grid data represents the city data at a specific range of latitude and longitude and at a certain granularity. The concerned area is partitioned into  $I \times J$  grids according to the geographical coordinates. Given urban data  $D = \{d_{i,j}, 0 \leq i \leq I - 1, 0 \leq j \leq J - 1\}$ ,  $d_{i,j}$  indicates the feature value within the latitude and longitude range to which the grid  $(i, j)$  belongs.

\*Corresponding Author.

Table 1: Summary of data formats. The common methods of saving, loading and getting data are omitted.

Type	Properties	Methods
Graph layer <sup>1</sup>	<u>name</u> , <sup>2</sup> lon, lat, directed, year, data	construct_node, construct_edge, construct_graph
Grid layer	<u>name</u> , start_lat, end_lat, step_lat, start_lon, end_lon, step_lon, year, data	construct_grid, get_value, get_region, get_value_by_grid, print_info, save_grid
Point layer	<u>name</u> , feature_name, <sup>3</sup> year, data	add_points, delete_points, get_value, get_value_by_range, to_gpd
Linestring layer	<u>name</u> , year, data	add_linestrings, delete_linestrings
Polygon layer	<u>name</u> , year, data	add_polygons, delete_polygons

<sup>1</sup> The methods of NetworkX [1] are also available for the GraphLayer object.

<sup>2</sup> The underline is the key of the node features.

<sup>3</sup> The feature\_name is the list of point features' types besides the main feature.

20  $i(i = 0, 1, \dots, I - 1)$  and  $j(j = 0, 1, \dots, J - 1)$  denote the indice of latitude (row) and longitude  
21 (column), respectively. At this stage, each grid data layer stores a single feature.

22 • **Graph data:** Graph data constructs the data as a graph  $G = (V, E, X)$ .  $V$  is the set of nodes,  $E$   
23 is the set of edges where each edge can be represented as  $e_{ij} = (v_i, v_j)$  and  $X$  is the set of node  
24 features. The properties of the edges are optional. At this stage, each graph data layer stores  
25 individual node characteristics along with the latitude and longitude of node.

26 • **Point data:** Point data is defined as  $P = \{p_i\}$ , where each point  $p_i = (x_i, y_i, \mathbf{v}_i)$  where  $(x_i, y_i)$   
27 stands for its geographical coordinate location and  $\mathbf{v}_i$  stands for its features' value. It can also be  
28 denoted as  $\mathbf{X}_i \in \mathbb{R}^C$  where  $C$  is the number of features. Each point layer may contain multiple  
29 features at this stage.

30 • **Linestring data:** Linestring data is usually applied to the representation of trajectories in city.  
31 Data in linestring layer is  $L = \{l_i\}$ , where  $l_i = [(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)]$  describes the  
32 lines between the points. The elements in the list can also be the type of point data. Each layer  
33 can be consist of multiple linestrings.

34 • **Polygon data:** Polygon data  $D_{plg} = d_i$  uses the points of a boundary to denote a linearly enclosed  
35 area and multiple areas cannot cross each other where  $d_i = [(x_0, y_0), (x_1, y_1), \dots, (x_m, y_m)]$ .  
36 The holes bounded by linear rings in the polygon region can also be stated. A polygon data  
37 layer can contain multiple polygons.

38 The data formats of UDL vary between different types of data layers and also have some common  
39 properties. Each UDL class basically comprises the layer name (the meaning of the urban data,  
40 e.g., nightlight, population, etc.), the year (year of the recorded data) and the data itself. Due to  
41 the urban nature of UDL, all data layers have geographic information identified by latitude and  
42 longitude (can also be replaced with other unique geographic identifiers).

43 For the defined five types of UDL, data operations like constructing data, modifying data and  
44 querying data by coordinates are provided. Subtle differences exist between various classes of  
45 UDL and a uniform interfaces to save and load data is accessible. The exclusive properties and  
46 methods of each class are described briefly in Table 1. More details can be found in the document  
47 of UDL<sup>1</sup>.

48 (3) **Granularity aligned data.** Granularity aligned data is processed from UDL layer data with  
49 different granularity through offered granularity transformation methods, specifying required  
50 spatial granularity. Data after granularity transformation differs from the original standardized  
51 layer data only in spatial granularity by some aggregation method, such as averaging, summing,  
52 and so on. Multi-source data of the same layer type can be consistent in both spatial coverage and  
53 granularity after this process.

54 (4) **Fused data.** Fused data is several aligned layer data fused from multiple urban data sources.  
55 Urban data is heterogeneous and the relationships between different domains also cannot be  
56 ignored [2]. Like in forecasting traffic flow, external factors including both weather conditions,  
57 temperature and wind speed were considered in [3]. Fused data is processed from granularity  
58 aligned data with the same granularity and then multi-source data can be integrated. The greater  
59 the variety of data types involved in a given task, the higher the probability of encountering a

<sup>1</sup><https://urbandatalayer-doc.readthedocs.io/en/latest/>

60 data scarcity issue [4]. The obtained fused data is also in UDL data format, where values are  
 61 concatenated or aggregated in a specified way.

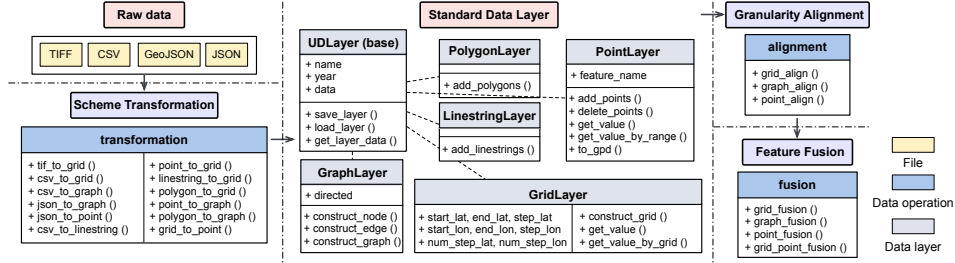


Figure 2: The design and structure of UDL interface. The raw data file format can be directly transformed into UDL through *Scheme Transformation* as some layers receiving files to initialize.

## 62 1.2 Three Data Processing Steps

63 The UrbanDataLayer builds the data layers and user-friendly APIs that make it easier to process and  
 64 reuse city data in urban research, provided with both codes and some city data. The overall structure  
 65 and processing flow are shown in Fig. 2. And the processing operations of UrbanDataLayer between  
 66 four data wrappers are as follows.

67 (1) **Scheme Transformation.** Scheme transformation consists of two types: from raw data source to  
 68 standard data layer and between data layers. For the type from raw data, methods automatically  
 69 transforming raw data into UDL layer data are provided. Structured data handling is considered,  
 70 such as CSV files, TIFF files and JSON files, etc. For instance, `tif_to_grid()` interface  
 71 transforms a TIFF file to a grid layer data with customized scope and granularity. We design  
 72 this pipeline to address the challenge for researchers to go through the complex data processing  
 73 process repeatedly when using similar urban datasets and we attempt to define it as a unified  
 74 data transformation paradigm. All raw files need to contain geographical location information.  
 75 After this, all data is converted to harmonized standard layer data. The mutual conversion  
 76 within data layers are also supplied. Here we provide eight transformations between UDL data  
 77 layers as illustrated in Fig. 3. For instance, when transforming polygons to graphs such as  
 78 `polygon_to_graph()`, vertices and line segments become nodes and edges, respectively. It  
 should be noted that such a basic transformation may be lossy (marked in red arrows).

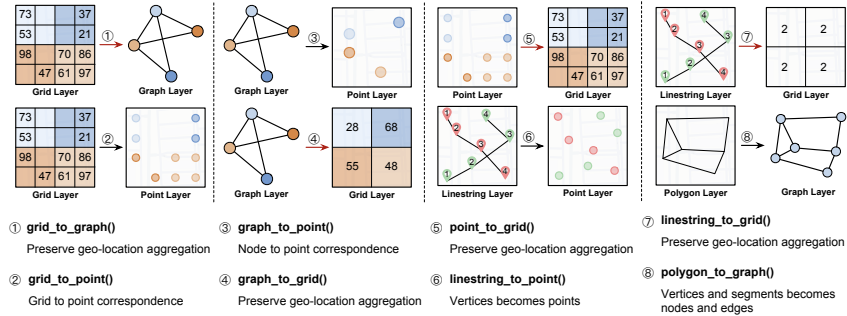


Figure 3: Transformation within layers. Red arrows indicate that there is intra-area aggregation during the transformation process, which may lose some precision.

79

80 (2) **Granularity Alignment.** Granularity alignment is a transformation among different granularities  
 81 of data in the same data layer according to a specified aggregation, mainly considering the  
 82 transformation from fine-grained data to coarse-grained data. Identical data sources may be  
 83 processed at multiple granularities in different scenarios. Instantiated by a specific task, when

Table 2: Data used in cases and corresponding UrbanDataLayer types.

Dataset	Region(s)	Type	Year <sup>2</sup>
Population	Shanghai, New York State	Grid	2020
Nightlight	Shanghai, New York State	Grid	2016
Built-up surface	Shanghai, New York State	Grid	2020
Roadnet intersection	Shanghai, New York State	Grid	2022, 2023
PM <sub>2.5</sub>	Shanghai, New York State	Grid	2016, 2019
SMOD	Shanghai, New York State	Grid	2015
POI	Shanghai, New York City	Point	2022
Boundary	Shanghai, New York City	Polygon	2021, 2023
El Nino	Equatorial Pacific	Graph <sup>1</sup>	1989 - 1998

<sup>1</sup> Its original data type is Point.

<sup>2</sup> The two years indicate the year of data for two different regions, respectively.

84 predicting PM<sub>2.5</sub> concentrations for cities, it may be necessary to try different divisions of the  
 85 city in order to find the most suitable grid units for prediction by `grid_align()`. This type  
 86 of interface is designed mainly for grid, graph and point layers. We provide several optional  
 87 aggregation methods and users can also define the aggregation function by themselves. Also,  
 88 granularity alignment for multiple data source is a imperative operation before data fusion when  
 89 faced with multi-sources data of varying granularities.

90 (3) **Feature Fusion.** Feature fusion mainly occurs mainly among identical UDL layer type with the  
 91 same granularity. This interface can be accessed after the granularity alignment. Considering  
 92 the solutions of urban tasks usually combine data fusion into the model and training process,  
 93 especially learning a representation of the original features from diverse datasets through the  
 94 utilization of deep neural networks (DNN) [5, 2]. We provide aggregation and concatenation of  
 95 cross-domain data which are most commonly used and easy to handle in the subsequent operations.  
 96 We also provide a hybrid data fusion between point and grid `grid_point_fusion()` which is  
 97 used in the Identification of administrative boundaries task.

98 In summary, the processing flow for urban data using UDL concludes as follows: (1) Convert  
 99 to unified and standardized layer data with *scheme transformation*. (2) Transform to specified  
 100 and aligned layer data by *granularity alignment*. (3) Complete necessary data fusion prior to the  
 101 subsequent tasks through *feature fusion*.

### 102 1.3 Extensibility

103 We adopt a layered structure that is convenient for storing various types of city data. Data of different  
 104 types, domains and times can be quickly inserted into the UDL, making it easy for all researchers to  
 105 extend global urban data layer. Moreover, the unified layer structure allows the expeditious expansion  
 106 of user-defined functions and the data output can be seamlessly accessed to other algorithms with high  
 107 performance, helping subsequent researchers to reproduce it. The extensibility of UDL is promising  
 108 to facilitate research on big data urban computing.

## 109 2 Data Access

110 Our UDL codebase is distributed under the MIT license. Users can easily access the tools through  
 111 interfaces provided in <https://github.com/SJTU-CILAB/udl>.

## 112 3 Detailed Experiment Settings

### 113 3.1 Dataset Description

114 All datasets used in this paper are described here. All these datasets contain information on urban  
115 space. Data in each task have been aligned. The data used in this paper is also provided as ready-to-use  
116 layers as listed in Table 2.

- 117 • *Population*: This dataset from WorldPop records the population counts [6] of each unit. The  
118 original data granularity is  $100\text{m} \times 100\text{m}$ .
- 119 • *Night-time light*: This VIIRS night-time lights data is radiation value measured in  
120  $\text{nanoWatts}/\text{cm}^2/\text{sr}$  from WorldPop [7]. The original data source comes from NOAAs National  
121 Centers for Environmental Information, Visible Infrared Imaging Radiometer Suite (VIIRS). The  
122 original data granularity is  $100\text{m} \times 100\text{m}$ .
- 123 • *PM<sub>2.5</sub>*: This dataset [8] provides an annual global  $0.01^\circ$  surface of concentrations ( $\mu\text{g}/\text{m}^3$ ) of all  
124 composition ground-level fine particulate matter of 2.5 micrometers or smaller (PM<sub>2.5</sub>).
- 125 • *Built-up surface*: The built-up surface ( $\text{m}^2$ ) is the gross surface (including the thickness of the walls)  
126 bounded by the building wall perimeter, which from Global Human Settlement Layer (GHSL).
- 127 • *SMOD*: This dataset from GHSL indicates the Degree of Urbanisation by delineating and classifying  
128 settlement typologies via a logic of cell clusters population size, population and built-up area  
129 densities.
- 130 • *POI*: The POIs (point of interest) of Shanghai and New York City are cleaned from Baidu Maps [9]  
131 and Safe Graph [10], respectively.
- 132 • *Roadnet intersection*: The road network intersection data are collected from OpenStreetMap [11],  
133 representing the number of intersections in each unit.
- 134 • *Districts boundary*: The Shanghai administrative districts boundary is exported from DataV [12]  
135 and New York City from NYC Open Data [13].
- 136 • *El nino*: The El Nino dataset is provided by UCI repository [14] and we use a subset of the data  
137 provided by [15] containing 93,935 samples, measuring oceanographic and surface meteorological  
138 variables.

139 Not all data were updated to the same latest year, and we choose updated latest real data over  
140 predicted data. Since the data used is considered to be relatively static in the city, it does not fluctuate  
141 significantly over a few years.

### 142 3.2 Four Empirical Cases Detailed Settings

#### 143 3.2.1 PM<sub>2.5</sub> concentration prediction.

144 Grid layer type is used in this case. The granularity of the data is  $0.02^\circ \times 0.02^\circ$  per grid in China, and  
145  $0.05^\circ \times 0.05^\circ$  per grid in New York State. The grid aggregation method is an average of the original  
146 values. Missing values in data are replaced by the mean along each column. We treat each grid as a  
147 individual sample in this case. The data is split into training data and test data at a ratio of 9:1.

#### 148 3.2.2 Built-up areas classification.

149 The data used in classification are all organized in grid layer format and well-aligned in granularity  
150  $0.01^\circ \times 0.01^\circ$  in space. As a classification task, we convert the built-up surface which is continuous  
151 data into boolean values by setting a threshold (in Shanghai the threshold is set as 90000, while in  
152 New York State it is 1500) to decide which grid is an urban area so that the proportion of positive and  
153 negative samples is close. The data is split into training data and test data at a ratio of 8:2.

154 **3.2.3 Identification of administrative boundaries.**

155 The given multi-modal data including POIs (point of interest), road intersection density, population  
156 and night-time light consists of both point data and grid data. The feature value of point type data  
157 POI is its site category here, and the other three are grid data of  $0.05^\circ \times 0.05^\circ$  granularity for both  
158 city. Note that, we only consider the coordinates of POIs and sample grid feature data as coordinate  
159 points  $(x, y)$  according to the grid value to perform data fusion. For POI data, we randomly sample  
160  $\frac{1}{100}$  coordinate points from Shanghai and  $\frac{1}{10}$  from New York City due to different POI density. Here,  
161 5749 points remain in Shanghai and 11311 points remain in New York City. For other three grid data,  
162 we randomly generate coordinate points in each grid from their specific distribution determined by  
163 grid value. In each grid  $(i, j)$ , the point number

$$|P(i, j)| = \frac{d_{i,j}}{\max_{\substack{0 \leq s \leq I-1 \\ 0 \leq t \leq J-1}} d_{s,t}} \cdot k \quad (1)$$

164  $d_{i,j}$  indicates the feature value within the latitude and longitude range to which the grid belongs. To  
165 be noted, the highest grid value is selected among New York State.  $k$  is 1000 for Shanghai and 10000  
166 for New York City. After that, we mix all the generated coordinate points as the fused data. The  
167 number of clusters is set according to the administrative districts, where Shanghai is set to 15 and  
168 New York City is set to 5. Chongming district of Shanghai is excluded due to the lack of POI data.

169 **3.2.4 El Nino anomaly detection.**

170 The data before processing is tabular data and consists of the following variables: date, latitude,  
171 longitude and zonal winds etc. twelve attributes. First we convert the date to a standard timestamp  
172 and then build the graph  $\mathcal{G}$  based on the 93,935 sample. Each sample in the data are constructed as a  
173 node instantiated by its own node attributes. The edge sets of  $\mathcal{G}$  contains three types: temporal (on  
174 the same day), spatial (adjacent within the grid distribution) and spatio-temporal (both of the former  
175 two). The constructed graph  $\mathcal{G}$  comprises 93,935 nodes and 56,687,915 edges.

176 Various anomaly detection methods [16] are used to compare the performance of different combina-  
177 tions of node features:

- 178 • LOF [17]: Local outlier factor (LOF) is a method for finding outliers in a multidimensional  
179 dataset. It takes a local instead of a global view on outliers and the degree of isolation depends on  
180 surrounding neighborhood.
- 181 • IF [18]: Isolation Forest (iForest) is a model-based method which builds an ensemble of iTrees for  
182 a given data set and finds anomalies as those instances having short average path lengths on the  
183 iTrees.
- 184 • CoLA [19]: Contrastive self-supervised Learning framework for Anomaly detection on attributed  
185 networks (CoLA) is the first contrastive self-supervised learning-based method for graph anomaly  
186 detection and can efficiently captures the local information of a node and its neighboring substructure  
187 by a novel type of contrastive instance pair. Especially it is friendly to large-scale networked  
188 data.
- 189 • ANOMALOUS [20]: ANOMALOUS innovatively optimize attribute selection and anomaly  
190 detection as a whole. It selects representative instances based on CUR decomposition and uses  
191 residual analysis to measure the normality.
- 192 • GAE [21]: Variational Graph Auto-Encoders (GAE) is a framework for unsupervised learning on  
193 graph-structured data which makes use of latent variables and is capable of learning interpretable  
194 latent representations for undirected graphs.
- 195 • OCGNN [22]: One Class Graph Neural Network (OCGNN) is a novel framework, which combines  
196 the expressive power of Graph Neural Networks (GNNs) with the classical one-class objective.  
197 OCGNN outperforms other baseline methods in terms of anomaly detection accuracy and robustness.  
198 By leveraging the inherent structure and features of the graph, OCGNN effectively captures and  
199 represents the complex relationships and patterns within the data.

- 200 • ONE [23]: A novel unsupervised network embedding algorithm, which is designed to handle  
201 attributed networks with outliers. ONE generates robust network embeddings to minimize the  
202 influence of outlier nodes.

## 203 4 Datasheet

204 This appendix provides a datasheet for the used dataset. The format of this datasheet was introduced  
205 in [24] and consolidates the motivation, creation process, composition, and intended uses of our  
206 dataset as a series of questions and answers.

### 207 4.1 Motivation

208 Q1 **For what purpose was the dataset created?** *Was there a specific task in mind? Was there a*  
209 *specific gap that needed to be filled? Please provide a description.*

210 Urban-related research lacks standard benchmarks with unified data format. The UDL data  
211 aims to fill the gap of absent universal urban data.

212 Q2 **Who created the dataset(e.g.,whichteam,researchgroup) and on behalf of which entity**  
213 **(e.g., company, institution, organization)?**

214 The authors of this work (from Shanghai Jiao Tong University) created the UDL data.

215 Q3 **Who funded the creation of the dataset?** *If there is an associated grant, please provide the*  
216 *name of the grantor and the grant name and number.*

217 This work was sponsored by National Natural Science Foundation of China under Grant No.  
218 62102246, 62272301, and Provincial Key Research and Development Program of Zhejiang  
219 under Grant No. 2021C01034. Part of the work was done when the students were doing  
220 internships at Yunqi Academy of Engineering.

221 Q4 **Any other comments?**

222 No.

### 223 4.2 Composition

224 Q5 **What do the instances that comprise the dataset represent (e.g., documents, photos, people,**  
225 **countries)?** *Are there multiple types of instances (e.g., movies, users, and ratings; people and*  
226 *interactions between them; nodes and edges)? Please provide a description.*

227 UDL data is composed of three components: (1) five base classes of UDL data; (2) a series of  
228 data process APIs to allow data transformation, alignment and fusion; (3) available UDL data  
229 for empirical urban cases.

230 Q6 **How many instances are there in total (of each type, if appropriate)?**

231 As for available UDL data for empirical cases in paper, there are six grid, one point, one polygon  
232 and one graph type UDL data regardless of granularity and region. However, UDL is a universal  
233 data conversion interface. Therefore, it will enable people to create more data instances using  
234 their own data.

235 Q7 **Does the dataset contain all possible instances or is it a sample (not necessarily random)**  
236 **of instances from a larger set?** *If the dataset is a sample, then what is the larger set? Is the*  
237 *sample representative of the larger set (e.g., geographic coverage)? If so, please describe how*  
238 *this representativeness was validated/verified. If it is not representative of the larger set, please*  
239 *describe why not (e.g., to cover a more diverse range of instances, because instances were*  
240 *withheld or unavailable).*

241 UDL sample dataset contain all possible type instance except linestring. Its large set lies in the  
242 expansion of spatial and temporal dimensions. The sample data can be any area and any time of  
243 available raw data.

244 Q8 **What data does each instance consist of? “Raw” data (e.g., unpro cessed text or images)**  
245 **or features?** *In either case, please provide a description.*

- 246 In the each type of UDL data, each value (per grid, node, point, etc.) represents a numerical  
247 value of urban characteristics at geographic location.
- 248 **Q9 Is there a label or target associated with each instance?** *If so, please provide a description.*  
249 No, we do not explicitly define a label or target for the instances.
- 250 **Q10 Is any information missing from individual instances?** *If so, please provide a description,*  
251 *explaining why this information is missing (e.g., because it was unavailable). This does not*  
252 *include intentionally removed information, but might include, e.g., redacted text.*  
253 Some value may be missing if the raw data is missing before transforming to the UDL data.
- 254 **Q11 Are relationships between individual instances made explicit (e.g., users' movie ratings,**  
255 **social network links)?** *If so, please describe how these relationships are made explicit.*  
256 The relationships between individual instances depend on raw data, which can be geographical  
257 adjacency, trajectory sequence, etc.
- 258 **Q12 Are there recommended data splits (e.g., training, development/ validation, testing)?** *If*  
259 *so, please provide a description of these splits, explaining the rationale behind them.*  
260 Not applicable.
- 261 **Q13 Are there any errors, sources of noise, or redundancies in the dataset?** *If so, please provide*  
262 *a description.*  
263 Since the data is transformed from various raw data in different formats and granularities, the  
264 data in UDL may have precision loss due to the geographic alignment. However, this totally  
265 depends on how the users of UDL put the settings.
- 266 **Q14 Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,**  
267 **websites, tweets, other datasets)?** *If it links to or relies on external resources, a) are there*  
268 *guarantees that they will exist, and remain constant, over time; b) are there official archival*  
269 *versions of the complete dataset (i.e., including the external resources as they existed at the*  
270 *time the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated*  
271 *with any of the external resources that might apply to a dataset consumer? Please provide*  
272 *descriptions of all external resources and any restrictions associated with them, as well as*  
273 *links or other access points, as appropriate.*  
274 UDL data relies on available open data source. a) It depends on the source data. However,  
275 since we aim to build urban research benchmark, it will remain constant when it comes to  
276 UDL format and can be reused. b) Yes. c) No.
- 277 **Q15 Does the dataset contain data that might be considered confidential (e.g., data that is**  
278 **protected by legal privilege or by doctor patient confidentiality, data that includes the**  
279 **content of individuals' non-public communications)?** *If so, please provide a description.*  
280 No.
- 281 **Q16 Does the dataset contain data that, if viewed directly, might be offensive, insulting,**  
282 **threatening, or might otherwise cause anxiety?** *If so, please describe why.*  
283 No.
- 284 **Q17 Does the dataset relate to people?** *If not, you may skip remaining questions in this section.*  
285 No.
- 286 **Q18 Does the dataset identify any subpopulations (e.g., by age, gender)?** *If so, please de-*  
287 *scribe how these subpopulations are identified and provide a description of their respective*  
288 *distributions within the dataset.*  
289 No.
- 290 **Q19 Is it possible to identify individuals (i.e., one or more natural persons), either directly or**  
291 **indirectly (i.e., in combination with other data) from the dataset?** *If so, please describe*  
292 *how.*  
293 No.



294 Q20 **Does the dataset contain data that might be considered sensitive in any way (e.g., data**  
295 **that reveals race or ethnic origins, sexual orientations, religious beliefs, political opinions**  
296 **or union member ships, or locations; financial or health data; biometric or genetic data;**  
297 **forms of government identification, such as social security numbers; criminal history)? If**  
298 *so, please provide a description.*

299 No.

300 Q21 **Any other comments?**

301 No.

### 302 4.3 Collection Process

303 Q22 **How was the data associated with each instance acquired? Was the data directly observable**  
304 *(e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly*  
305 *inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or*  
306 *language)? If the data was reported by subjects or indirectly inferred/derived from other data,*  
307 *was the data validated/verified? If so, please describe how.*

308 We gain the data from the open source mentioned in the supplementary, e.g, WorldPop and  
309 GHSL.

310 Q23 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatuses**  
311 **or sensors, manual human curation, software programs, software APIs)? How were these**  
312 *mechanisms or procedures validated?*

313 We collect data mainly by two ways: (1) directly download the data file from the open website  
314 (e.g., tif and shape files); (2) use provided APIs to access raw data such as Baidu Map.

315 Q24 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deter-**  
316 **ministic, probabilistic with specific sampling probabilities)?**

317 UDL data is not a sample from a larger set.

318 Q25 **Who was involved in the data collection process (e.g., students, crowdworkers, contrac-**  
319 **tors) and how were they compensated (e.g., how much were crowdworkers paid)?**

320 Only the authors of this paper have been involved in the data collection process. No extra  
321 payment was made.

322 Q26 **Over what timeframe was the data collected? Does this timeframe match the creation**  
323 **timeframe of the data associated with the instances (e.g., recent crawl of old news**  
324 **articles)? If not, please describe the timeframe in which the data associated with the instances**  
325 *was created.*

326 The road network data was collected during December 2022. Other data was collected during  
327 2023. They are not the latest data, but considering the inherent properties of urban they would  
328 not be changing a lot.

329 Q27 **Were any ethical review processes conducted (e.g., by an institutional review board)? If**  
330 *so, please provide a description of these review processes, including the outcomes, as well as*  
331 *a link or other access point to any supporting documentation.*

332 Not applicable.

### 333 4.4 Preprocessing/cleaning/labeling

334 Q28 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucket-**  
335 **ing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances,**  
336 **processing of missing values)? If so, please provide a description. If not, you may skip the**  
337 *remaining questions in this section.*

338 We aggregate the data at a specified geospatial granularity. The exact form and granularity of  
339 which are mentioned in the paper.

340 Q29 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to**  
341 **support unanticipated future uses)?** *If so, please provide a link or other access point to the*  
342 *“raw” data.*

343 Yes.

344 – <https://www.worldpop.org/>

345 – <https://ghsl.jrc.ec.europa.eu/>

346 – <https://www.openstreetmap.org/>

347 – [https://data.cityofnewyork.us/City-Government/Borough-Boundaries/  
348 tqmj-j8zm](https://data.cityofnewyork.us/City-Government/Borough-Boundaries/tqmj-j8zm)

349 – <https://www.safegraph.com/products/places>

350 – [http://datav.aliyun.com/portal/school/atlas/area\\_selector](http://datav.aliyun.com/portal/school/atlas/area_selector)

351 Q30 **Is the software that was used to preprocess/clean/label the data available?** *If so, please*  
352 *provide a link or other access point.*

353 Yes. The only software that was used to preprocess the data is Python, which is free to all  
354 users.

355 Q31 **Any other comments?**

356 No.

#### 357 4.5 Uses

358 Q32 **Has the dataset been used for any tasks already?** *If so, please provide a description.*

359 Yes. We have presented four illustrative empirical cases, including PM<sub>2.5</sub> concentration  
360 prediction, built-up areas classification, identification of administrative boundaries and El Nino  
361 anomaly detection. See Section 4 in the main paper.

362 Q33 **Is there a repository that links to any or all papers or systems that use the dataset?** *If so,*  
363 *please provide a link or other access point.*

364 N/A

365 Q34 **What (other) tasks could the dataset be used for?**

366 Any urban-related tasks applying the standardized layer can use the data.

367 Q35 **Is there anything about the composition of the dataset or the way it was collected and**  
368 **preprocessed/cleaned/labeled that might impact future uses? For example, is there**  
369 **anything that a dataset consumer might need to know to avoid uses that could result in**  
370 **unfair treatment of individuals or groups (e.g., stereotyping, quality of service issues) or**  
371 **other risks or harms (e.g., legal risks, financial harms)?** *If so, please provide a description.*  
372 *Is there anything a dataset consumer could do to mitigate these risks or harms?*

373 No.

374 Q36 **Are there tasks for which the dataset should not be used?** *If so, please provide a description.*

375 No.

376 Q37 **Any other comments?**

377 No.

#### 378 4.6 Distribution

379 Q38 **Will the dataset be distributed to third parties outside of the entity (e.g., company,**  
380 **institution, organization) on behalf of which the dataset was created?** *If so, please provide*  
381 *a description.*

382 No.

383 Q39 **How will the dataset be distributed (e.g., tarball on websites, API, GitHub)?** *Does the*  
384 *dataset have a digital object identifier (DOI)?*

385 The data will be found in GitHub.

- 386 Q40 **When will the dataset be distributed?**  
387 [The sample data is already available.](#)
- 388 Q41 **Will the dataset be distributed under a copyright or other intellectual property (IP)**  
389 **license, and/or under applicable terms of use (ToU)?** *If so, please describe this license*  
390 *and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant*  
391 *licensing terms or ToU, as well as any fees associated with these restrictions.*  
392 [The UDL dataset will be distributed under the MIT license.](#)
- 393 Q42 **Have any third parties imposed IP-based or other restrictions on the data associated with**  
394 **the instances?** *If so, please describe these restrictions, and provide a link or other access*  
395 *point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated*  
396 *with these restrictions.*  
397 [No.](#)
- 398 Q43 **Do any export controls or other regulatory restrictions apply to the dataset or to indi-**  
399 **vidual instances?** *If so, please describe these restrictions, and provide a link or other access*  
400 *point to, or otherwise reproduce, any supporting documentation.*  
401 [No.](#)
- 402 Q44 **Any other comments?**  
403 [No.](#)
- 404 **4.7 Maintenance**
- 405 Q45 **Who will be supporting/hosting/maintaining the dataset?**  
406 [The authors.](#)
- 407 Q46 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**  
408 [Contact the corresponding author or first author in the author list.](#)
- 409 Q47 **Is there an erratum?** *If so, please provide a link or other access point.*  
410 [N/A.](#)
- 411 Q48 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete**  
412 **instances)?** *If so, please describe how often, by whom, and how updates will be communicated*  
413 *to dataset consumers (e.g., mailing list, GitHub)?*  
414 [Yes, we will keep polishing our processed data and everyone is encouraged to contribute new](#)  
415 [UDL data.](#)
- 416 Q49 **If the dataset relates to people, are there applicable limits on the retention of the data**  
417 **associated with the instances (e.g., were the individuals in question told that their data**  
418 **would be retained for a fixed period of time and then deleted)?** *If so, please describe these*  
419 *limits and explain how they will be enforced.*  
420 [Not applicable.](#)
- 421 Q50 **Will older versions of the dataset continue to be supported/hosted/maintained?** *If so,*  
422 *please describe how. If not, please describe how its obsolescence will be communicated to*  
423 *dataset consumers.*  
424 [Yes. Our UDL data has "year" attribute.](#)
- 425 Q51 **If others want to extend/augment/build on/contribute to the dataset, is there a mech-**  
426 **anism for them to do so?** *If so, please provide a description. Will these contributions*  
427 *be validated/verified? If so, please describe how. If not, why not? Is there a process for*  
428 *communicating/distributing these contributions to dataset consumers? If so, please provide a*  
429 *description.*  
430 [Yes. everyone can release data in UDL format along with their research or contact us to add](#)  
431 [their UDL data in this work repository.](#)
- 432 Q52 **Any other comments?**  
433 [No.](#)

434 **References**

- 435 [1] Aric A. Hagberg, Daniel A. Schult, and Pieter J. Swart. Exploring network structure, dynamics,  
436 and function using networkx. In Gaël Varoquaux, Travis Vaught, and Jarrod Millman, editors,  
437 *Proceedings of the 7th Python in Science Conference*, pages 11 – 15, Pasadena, CA USA, 2008.
- 438 [2] Jia Liu, Tianrui Li, Peng Xie, Shengdong Du, Fei Teng, and Xin Yang. Urban big data fusion  
439 based on deep learning: An overview. *Information Fusion*, 53:123–133, 2020.
- 440 [3] X. Zhang, C. Huang, Y. Xu, L. Xia, P. Dai, L. Bo, J. Zhang, and Y. Zheng. Traffic flow  
441 forecasting with spatial-temporal graph diffusion network. In *AAAI*, 2021.
- 442 [4] Chuishi Meng, Yanhua Li, Yu Zheng, Jieping Ye, Qiang Yang, Philip S Yu, and Ouri Wolfson.  
443 The 12th international workshop on urban computing. In *Proceedings of the 29th ACM SIGKDD*  
444 *Conference on Knowledge Discovery and Data Mining*, pages 5874–5875, 2023.
- 445 [5] Yu Zheng. Methodologies for cross-domain data fusion: An overview. *IEEE transactions on*  
446 *big data*, 1(1):16–34, 2015.
- 447 [6] Global high resolution population denominators project, 2018.
- 448 [7] Worldpop. Geospatial covariate data layers.
- 449 [8] A. van Donkelaar Hammer, M. S. Global annual pm2.5 grids from modis, misr and seawifs  
450 aerosol optical depth (aod). *New York: NASA Socioeconomic Data and Applications Center*  
451 *(SEDAC)*, 2022.
- 452 [9] Baidu map. <https://map.baidu.com/>.
- 453 [10] Global points of interest (poi) data | safegraph places. [https://www.safegraph.com/](https://www.safegraph.com/products/places)  
454 [products/places](https://www.safegraph.com/products/places).
- 455 [11] Openstreetmap. <https://www.openstreetmap.org/>.
- 456 [12] Aliyun. [http://datav.aliyun.com/portal/school/atlas/area\\_selector](http://datav.aliyun.com/portal/school/atlas/area_selector).
- 457 [13] Nyc open data. [https://data.cityofnewyork.us/City-Government/](https://data.cityofnewyork.us/City-Government/Borough-Boundaries/tqmj-j8zm)  
458 [Borough-Boundaries/tqmj-j8zm](https://data.cityofnewyork.us/City-Government/Borough-Boundaries/tqmj-j8zm).
- 459 [14] El Nino. UCI Machine Learning Repository, 1999. DOI: <https://doi.org/10.24432/C5WG62>.
- 460 [15] Guanjie Zheng, Susan L Brantley, Thomas Lauvaux, and Zhenhui Li. Contextual spatial  
461 outlier detection with metric learning. In *Proceedings of the 23rd ACM SIGKDD international*  
462 *conference on knowledge discovery and data mining*, pages 2161–2170, 2017.
- 463 [16] Kay Liu, Yingtong Dou, Yue Zhao, Xueying Ding, Xiyang Hu, Ruitong Zhang, Kaize Ding,  
464 Canyu Chen, Hao Peng, Kai Shu, George H. Chen, Zhihao Jia, and Philip S. Yu. Pygod: A  
465 python library for graph outlier detection. *arXiv preprint arXiv:2204.12095*, 2022.
- 466 [17] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. Lof: identifying  
467 density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference*  
468 *on Management of data*, pages 93–104, 2000.
- 469 [18] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *2008 eighth ieee*  
470 *international conference on data mining*, pages 413–422. IEEE, 2008.
- 471 [19] Yixin Liu, Zhao Li, Shirui Pan, Chen Gong, Chuan Zhou, and George Karypis. Anomaly  
472 detection on attributed networks via contrastive self-supervised learning. *IEEE transactions on*  
473 *neural networks and learning systems*, 33(6):2378–2392, 2021.

- 474 [20] Zhen Peng, Minnan Luo, Jundong Li, Huan Liu, Qinghua Zheng, et al. Anomalous: A joint  
475 modeling approach for anomaly detection on attributed networks. In *IJCAI*, pages 3513–3519,  
476 2018.
- 477 [21] Thomas N Kipf and Max Welling. Variational graph auto-encoders. *arXiv preprint*  
478 *arXiv:1611.07308*, 2016.
- 479 [22] Xuhong Wang, Baihong Jin, Ying Du, Ping Cui, Yingshui Tan, and Yupu Yang. One-class  
480 graph neural networks for anomaly detection in attributed networks. *Neural computing and*  
481 *applications*, 33:12073–12085, 2021.
- 482 [23] Sambaran Bandyopadhyay, N Lokesh, and M Narasimha Murty. Outlier aware network embed-  
483 ding for attributed networks. In *Proceedings of the AAAI conference on artificial intelligence*,  
484 volume 33, pages 12–19, 2019.
- 485 [24] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna  
486 Wallach, Hal Daumé Iii, and Kate Crawford. Datasheets for datasets. *Communications of the*  
487 *ACM*, 64(12):86–92, 2021.