
UrbanDataLayer: A Unified Data Pipeline for Urban Science

Yiheng Wang¹, Tianyu Wang¹, Yuying Zhang¹
Hongji Zhang¹, Haoyu Zheng¹, Guanjie Zheng^{*1}, Linghe Kong^{*.1}
¹ Shanghai Jiao Tong University, Shanghai, China
{yhwang0828, wty500, shjtzyy01, zhanghongji,
langanzheng, gjzheng, linghe.kong}@sjtu.edu.cn

Abstract

1 The rapid progression of urbanization has generated a diverse array of urban
2 data, facilitating significant advancements in urban science and urban computing.
3 Current studies often work on separate problems case by case using diverse data,
4 e.g., air quality prediction, and built-up areas classification. This fragmented
5 approach hinders the urban research field from advancing at the pace observed in
6 Computer Vision and Natural Language Processing, due to two primary reasons.
7 On the one hand, the diverse data processing steps lead to the lack of large-scale
8 benchmarks and therefore decelerate iterative methodology improvement on a
9 single problem. On the other hand, the disparity in multi-modal data formats
10 hinders the combination of the related modal data to stimulate more research
11 findings. To address these challenges, we propose UrbanDataLayer (UDL), a suite
12 of standardized data structures and pipelines for city data engineering, providing a
13 unified data format for researchers. This allows researchers to easily build up large-
14 scale benchmarks and combine multi-modal data, thus expediting the development
15 of multi-modal urban foundation models. To verify the effectiveness of our work,
16 we present four distinct urban problem tasks utilizing the proposed data layer.
17 UrbanDataLayer aims to enhance standardization and operational efficiency within
18 the urban science research community. The examples and source code are available
19 at <https://github.com/SJTU-CILAB/udl>.

20 1 Introduction

21 The accelerated pace of urbanization has enhanced life quality while concurrently inducing issues
22 such as air pollution and traffic congestion. Extensive urban data has been recorded due to the
23 widespread use of advanced sensing technologies [29]. Simultaneously, urban studies have sprung up
24 among various domains of human mobility [15], air quality [6, 20], traffic dynamics [18], climate
25 change [35], spatial planning [48] and poverty [42, 32], etc. However, several *challenges* are posed.
26 *Firstly, numerous urban studies work on separate problems using different datasets case by case or*
27 *performing different processings on the same dataset. This lack of standard benchmarks hinders*
28 *the overall improvement of research.* In urban issues, researchers often self-define the problem and
29 propose methods accordingly. Based on an analysis of 88 papers published in seven AI conferences
30 shown in Fig. 1, three phenomena are observed. (1) Many urban problems are defined within the same
31 domain, yet disparate datasets are used for identical problems. (2) Even if they use the same datasets,

*Corresponding Author.

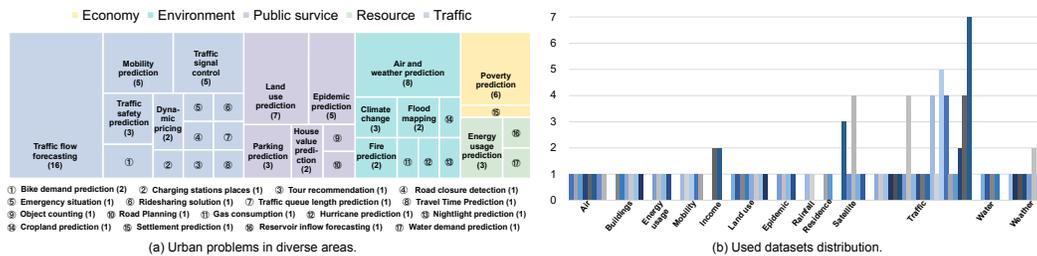


Figure 1: Problems and datasets in published papers. (a) Urban problems studied in five areas: **traffic**, **public service**, **environment**, **economy** and **resource** (from left to right). The numbers below are the count of relevant papers. (b) Dataset types for each category. Each bar represents the number of papers using that dataset. Papers use data of different datasets in similar domains and the distribution of datasets is very decentralized.

32 variations in data processing might lead to inconsistent experimental data. (3) More differences
 33 in the final experimental data may also exist that are not known due to the data not being publicly
 34 available. Even in relatively mature urban tasks such as urban spatial-temporal prediction, only less
 35 than 30% of the papers have made their data public [36]. As shown in Table 1, a significant portion
 36 of experimental data in urban studies remains inaccessible. This phenomenon makes comparisons
 37 between these methods difficult and the results are hard to reproduce due to non-public experimental
 38 data. Furthermore, researchers cannot continuously improve the performance of the methods under
 39 the same standard, which hinders the progress of urban research.

40 *Secondly, urban data exists in multiple modalities, miscellaneous formats, and non-uniform gran-*
 41 *ularity, and involves cumbersome processing; urban research often requires multiple data fusions.*
 42 *Repetitive and intricate data processing is troublesome and prone to errors, making data utilization*
 43 *poor.* Fusing knowledge from different datasets is effective and essential in urban research. Unlike
 44 Computer Vision and Natural Language Processing which have standardized datasets such as Im-
 45 ageNet [7] and WikiText-103 [30], urban datasets frequently adopt distinct storage formats with
 46 diverse granularities, encompassing images, tables, trajectories, points, and beyond. This challenge
 47 impedes researchers especially novices in the domain of efficiently and correctly combining and
 48 leveraging the data, which introduces obstacles in large-scale urban research.

49 Therefore, we propose an effective and efficient urban data management suite named UrbanDataLayer
 50 (UDL), which provides five standard urban data layers and efficient data processing tools with the
 51 following characteristics. (1) **Reproducible benchmark:** People can utilize UDL to easily process
 52 their data, make it a public benchmark, and compare with SOTA methods. (2) **Combinable multi-**
 53 **modal data:** We provide examples of combining urban data with spatio-temporal base data, e.g.,
 54 satellite image and road network data to create the possibility for multi-modal spatio-temporal
 55 foundation model building. (3) **Extensibility:** UDL can be expanded in both spatio-temporal and
 56 feature dimensions and encourages researchers to fill in the gaps of absent universal urban data.

57 2 Related Work

58 In contrast to other domains like Computer Vision, Natural Language Processing or tasks like Graph
 59 Node Classification have common datasets such as ImageNet [7] and CIFAR-10 [17], WikiText-
 60 103 [30], Cora [27], respectively. Regrettably, urban computing research lacks common datasets and
 61 data formats and somewhat inhibits the advancement of this field.

62 It has recently come to our attention that there is a benchmark LibCity [37] for solving urban spatio-
 63 temporal prediction problems. It includes pivotal stages related to traffic prediction into a systematic
 64 pipeline and provides 40 diverse datasets of unified storage format. It merely focuses on scenarios of
 65 urban traffic and does not cover all types of data in urban.

66 Data produced within urban areas typically exhibits an association with either spatial or spatiotemporal
 67 attributes [47]. Datasets originating from diverse domains present different structures, resulting in

Table 1: Data used in published research. The research of the same field works on separate datasets and most of them are not public. Take economy, air, and traffic domain as examples.

Domain	Data	Time span	Spatial coverage	Paper	Type	Used Time	Used Space	Public*
Economy	Digital Globe Worldview Satellite	-	Global	[13]	Polygon	-	South Korea	✗
	Villages images from Google Maps	2011	Global	[32]	Grid	2011	India	✗
	Nightlight from NOAA	2013	Global	[42]	Grid	2013	Africa	✗
	Nightlight from NASA	2012	Global	[28]	Grid	2012	Global	✓
	Expenditure (poverty) from LSMS	2011 - 2012	Uganda	[42]	Grid	2011 - 2012	Uganda	✗
				[2]	Grid	2011 - 2012	Uganda	✗
	Urban LIA (low- income areas) Data	-	Kisumu, Malindi, Nakuru	[19]	Point	-	Kisumu, Malindi, Nakuru	✗
Income statistics from SECC	2011	India	[32]	Grid	2011	India	✗	
Air	KDD CUP of Fresh Air	Jan. 1, 2017 - Apr. 30, 2018	Beijing	[12]	Graph	Jan. 1, 2017 - Apr. 30, 2018	Beijing	✗
	Urban Air data	Aug. 2012 - May. 2015	302 Chinese cities	[49]	Point	Aug. 2012 - May. 2015	Chinese mainland	✗
		Jan. 1, 2015 - Dec. 31, 2018	Chinese mainland	[20]	Point	May. 1, 2014 - Apr. 30, 2015	Beijing	✗
Traffic	NYC-Taxi	Jan. 1, 2015 - Mar. 1, 2015	New York City	[43]	Grid	Jan. 1, 2015 - Mar. 1, 2015	New York City	✓
				[45]	Grid	Jan. 1, 2015 - Mar. 1, 2015	New York City	✗
	Traffic dataset from Caltrans	2015 - 2016	San Francisco	[41]	Graph	2015 - 2016	San Francisco	✗

* Whether the processed data in the paper is public.

68 different representations. When confronting a problem, it is customary to extract knowledge from
69 numerous diverse datasets by data fusion. In particular, the recently proposed time-series large
70 models [40, 11] frequently fuse data from different domains to obtain knowledge.

71 In the last decades, work like Open Geospatial Consortium [1] has been dedicated to establishing
72 standards for geospatial data which is also related to urban data. However, the standards assembled as
73 OGC APIs are designed primarily for geospatial data's release and access, which can be viewed more
74 as a kind of "raw data". Unlike OGC, UDL aims to define an urban data pipeline that can process and
75 fuse data as input directly into the model. In addition, it is not limited to geospatial data and other
76 urban data like time series data are also in this scope.

77 3 UrbanDataLayer: A Data Suite for Urban Research

78 3.1 UDL layer-wise pipeline

79 The UDL (UrbanDataLayer) is a suite of standard data structures and pipelines for city data engineer-
80 ing, which processes city data from various raw data into a unified data format. The datasets used
81 in one research may have different types and formats, and often come from different sources [23].
82 Urban data inputs into the UDL undergo a series of transformations, including conversion from raw
83 data to standardized data layers, re-alignment of granularity, and fusion of disparate datasets, before
84 being utilized and stored. Consequently, we delineate **four stages of data wrapping** and **three data**
85 **processing steps** within the UDL, as depicted in Fig. 2.

86 The four data wrappers represent four stages in the data processing pipeline, transitioning from raw
87 data to fused data that can be directly utilized by models. These stages include the raw data source,
88 standard data layer, granularity-aligned data, and fused data, respectively. In standard data layer which
89 is the main component of UDL, we summarize the urban data into five data structures: grid, graph,
90 point, linestring and polygon. The details of each data layer are provided in the documentation¹.

91 The UrbanDataLayer builds the data layers and user-friendly APIs, simplifying the processing and
92 reuse of city data in urban research. As depicted in Fig. 2, the components of UrbanDataLayer
93 between four data wrappers are scheme transformation, granularity alignment, and feature fusion.

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

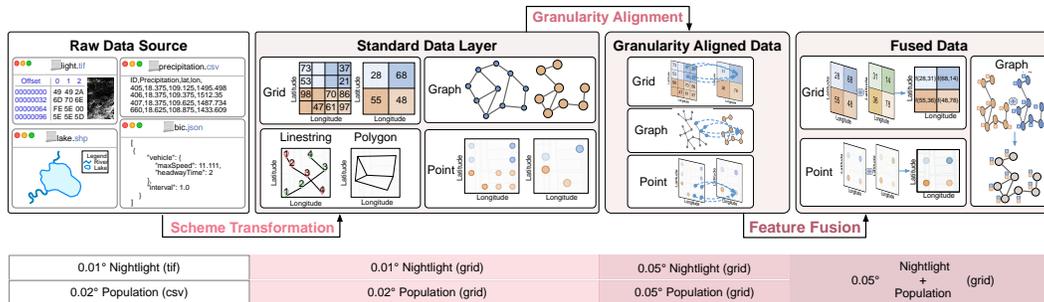


Figure 2: Overview of UrbanDataLayer framework. The words in red are the data processing steps.

94 In contemporary urban computing, datasets from diverse domains increasingly exhibit intercon-
 95 nections influenced by complex underlying relationships [46], underscoring the need for effective data
 96 fusion techniques to capture and leverage these connections. To exemplify the application of UDL
 97 (depicted at the bottom of Fig. 2), let's consider an example. Given nightlight data and population
 98 data in different formats and granularities, we aim to derive fused data for future downstream tasks.
 99 The process unfolds as follows: Firstly, we obtain standard grid data while preserving the original
 100 granularity through Scheme Transformation. Next, we acquire the target granularity data via Granu-
 101 larity Alignment. Subsequently, the fused data can be extracted through Feature Fusion. The entire
 102 process is managed by UDL.

103 3.2 General functionalities

104 For the defined five types of UDL layers, data operations like constructing, modifying and querying
 105 data by coordinates are provided. Besides this, users can easily access common data processing meth-
 106 ods through UDL interfaces. The main types of interfaces are as follows: (1) *Scheme Transformation*:
 107 Facilitates the transfer of raw data to UDL data and between data layers (Fig. 3). (2) *Granularity*
 108 *Alignment*: Converts a standard data layer into different spatial granularities. (3) *Feature Fusion*:
 109 Aggregates cross-domain data. The structure of the UDL interface is shown in Fig. 4.

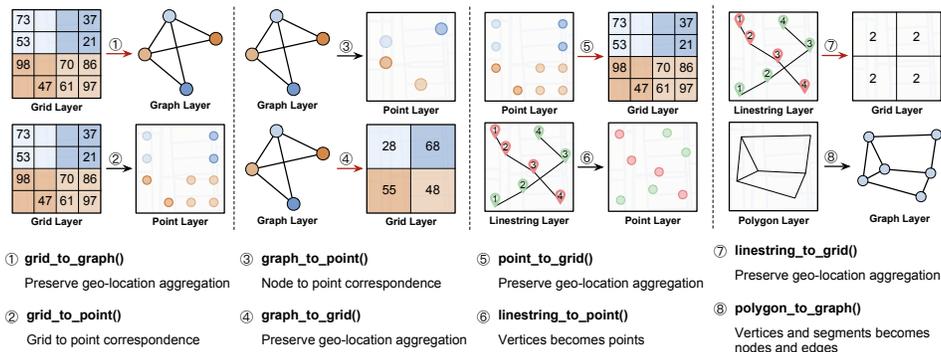


Figure 3: Transformation within layers. Red arrows indicate that there is intra-area aggregation during the transformation process, which may lose some precision.

110 3.3 Productivity

111 Data from diverse domains comprise numerous modalities, each recorded by distinct data types,
 112 distributions, scales, and granularities. For example, satellite images [13] are represented by pixel
 113 intensities, whereas POIs [39, 12] are usually represented by spatial points linked to a static category.
 114 Human mobility data [15] is embodied as trajectories, while road networks are represented as
 115 graph [18] and population data [21] is represented as grid-based data with real-value. The property of
 116 multiple data layers and friendly APIs of UDL well facilitates the combination of features, which

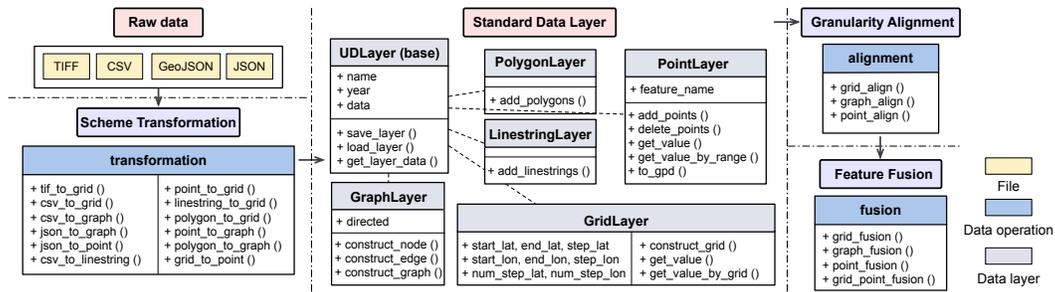


Figure 4: The design and structure of the UDL interface.

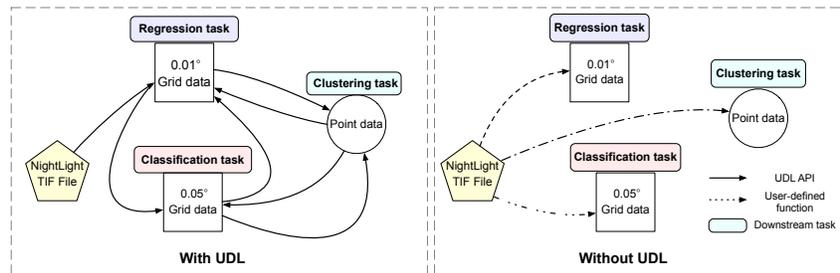


Figure 5: Implementing multiple downstream tasks with the same data through UDL data layers and unified APIs.

117 is evident in two aspects, as depicted in Fig. 5. First, the transition from diverse raw data sources
 118 to standard layer data of specified structures becomes routine and procedural, eliminating the need
 119 for repetitive processing (via user-defined functions) of similar data types. Second, various data
 120 layers can be quickly and easily aligned with or transformed to each other in UDL according to
 121 the geographic coordinate characteristics of urban data, rather than reprocessing from the original
 122 data every time. Both of their outputs can be directly used as inputs to the model or with minor
 123 adjustments.

124 Using the nightlight data from the three experimental cases described in later sections as an example,
 125 we include 0.02° and 0.01° grid data in Shanghai, 0.05° and 0.01° grid data in New York and point
 126 data. Without UDL, processing from the raw data needs to be conducted 6 times. However, with
 127 UDL, only 2 steps are required using the ready-to-use API. Subsequently, only 4 times exist between
 128 the UDL layer, where the users need to complete the conversion from 0 times.

129 Especially with the rise of large language models related to urban computing such as time-series
 130 foundation models, UDL facilitates easy data fusion for these models. To demonstrate this idea with
 131 an example of time-series foundation models UniTS [11] and PatchTST [31], various data types can
 132 be transformed into point data as inputs to both models. And this form is the main data provider for
 133 the current time-series models [50]. We anticipate that it will be a significant tool for urban-related
 134 large language models.

135 4 Empirical Cases

136 In this section, we use four typical downstream tasks to illustrate how UrbanDataLayer (UDL) can
 137 accelerate and enhance urban research. Four cases cover both supervised learning and unsupervised
 138 learning tasks, including $PM_{2.5}$ concentration prediction, built-up areas classification, identification
 139 of administrative boundaries, and El Nino anomaly detection. A more detailed description of data
 140 and implementation of cases are provided in <https://github.com/SJTU-CILAB/udl>.

Table 2: Effectiveness of combining different features in PM_{2.5} prediction problems. The best performance of the combination for each compared method is underlined and the best performance of all is bolded. Overall, in PM_{2.5} prediction, combining more features contributes to better performance.

Region Method Measurement	Shanghai						New York					
	XGBoost			MLP			XGBoost			MLP		
	RMSE	MAE	R^2	RMSE	MAE	R^2	RMSE	MAE	R^2	RMSE	MAE	R^2
Roadnet Intersection Density	3.953	4.861	0.181	3.710	4.779	0.199	0.608	0.778	0.368	0.720	0.898	0.040
Nightlight	4.327	5.127	0.089	4.821	5.859	-0.204	0.743	0.936	0.084	0.721	0.881	0.074
Population	4.374	5.134	0.086	4.373	5.185	0.057	0.740	0.932	0.093	0.676	0.854	0.130
Roadnet + Nightlight	3.672	4.582	0.272	3.404	4.276	0.359	0.591	0.762	0.393	0.648	0.828	0.183
Roadnet + Population	3.669	4.535	0.287	3.464	4.365	0.332	0.591	0.764	0.390	0.677	0.864	0.111
Nightlight + Population	3.974	4.783	0.207	4.044	4.810	0.189	0.713	0.901	0.151	0.619	0.792	0.252
Combining All	3.355	4.235	0.378	3.103	4.075	0.417	0.578	0.753	0.408	0.644	0.817	0.204

- 141 • The major experiments are composed of results of combining various features using classic methods
142 to demonstrate the benefit of unifying diverse data input via UDL. In this sense, innovating advanced
143 methods for each task is not within our scope.
- 144 • Successfully conducting the experiments justifies the fact that: (1) UDL facilitates easy processing
145 of data to build reproducible benchmarks; (2) UDL is applicable across different spatial regions,
146 temporal periods, and feature dimensions, thereby enabling the scaling up of spatial-temporal data.

147 4.1 PM_{2.5} concentration prediction

148 Accurate air quality prediction is of great importance to urban governance and human livelihood [12].
149 In this paper, we study the frequently-discussed PM_{2.5} concentration prediction problem [20, 22, 26].
150 We use XGBoost and MLP models, combining night-time lights, population, and road intersection
151 density as inputs, to conduct experiments in Shanghai, China (120°E ~ 122°E, 30°N ~ 32.4°N)
152 and New York State, United States (80°W ~ 70°W, 40°N ~ 45.5°N). The predicted results
153 are evaluated against the value obtained from the NASA Socioeconomic Data and Applications
154 Center (recognized as ground truth on all grids). Three metrics are considered respectively as RMSE,
155 MAE and R^2 . The data is split into training data and test data at a ratio of 9:1. Table 2 shows the
156 performance of different feature combinations on PM_{2.5} concentration prediction in two regions. It
157 is observed that combining more features performs better on XGBoost and MLP overall. The observed
158 patterns can be attributed to the strong spatial correlation between intersection density, nightlight,
159 population, and the PM_{2.5} (as shown in Fig. 6 and Fig. 7). The figures depict grid aggregation,
160 where each cell value represents the average of the original values within that cell. The granularity
161 of the data is 0.02° × 0.02° per grid in Shanghai, and 0.05° × 0.05° per grid in New York State.
162 An interesting observation is that areas with higher values for the three urban features—intersection
163 density, nightlight, and population—tend to exhibit higher PM_{2.5} concentrations, as seen in New
164 York City and downtown Shanghai. These results indicate that incorporating knowledge from more
165 domain-relevant data sources enhances the accuracy of environmental pollution predictions.

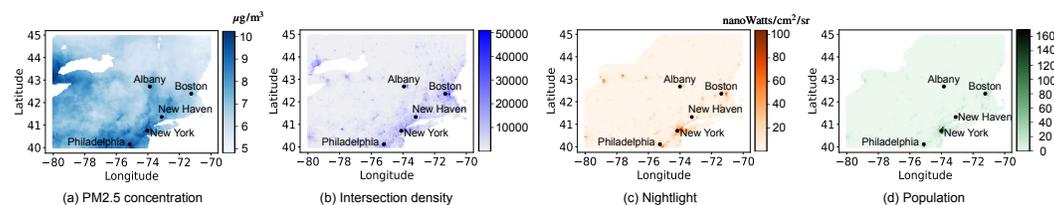


Figure 6: PM_{2.5}, intersection density, night-time light intensity, and population density in areas near New York State. Big cities, e.g., New York and Boston, exhibit high values across all four dimensions. Among these urban features, intersection density emerges as the most significant factor in predicting PM_{2.5} concentrations.

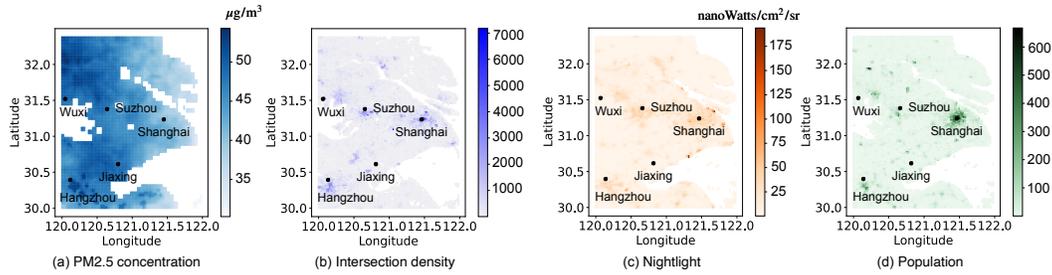


Figure 7: Urban data of PM_{2.5} prediction task in Shanghai. Missing values are imputed by the mean along each column. Values have the same meaning as in Fig. 6.

Table 3: Effectiveness of combining different features in built-up surface classification problems. The best performance of the combination for each compared method is underlined and the best performance of all is bolded.

Region	Shanghai					New York					
	Method	LR	DT	RF	GBDT	Adaboost	LR	DT	RF	GBDT	Adaboost
Accuracy	Nightlight	0.736	0.653	0.653	0.747	0.744	0.783	0.733	0.734	0.808	0.808
	SMOD	0.764	<u>0.767</u>	0.768	0.767	0.760	0.729	0.729	0.729	0.729	0.729
	Population	0.761	0.677	0.677	0.767	0.766	0.861	0.818	0.818	0.869	0.868
	Nightlight+SMOD	0.781	0.706	0.715	0.786	0.782	0.862	0.820	0.821	0.871	0.869
	Nightlight+Population	0.781	0.708	0.710	0.786	0.782	0.862	0.820	0.821	0.871	0.869
	SMOD+Population	0.781	0.707	0.732	0.786	0.782	0.862	<u>0.820</u>	0.823	0.871	0.869
	All	<u>0.782</u>	0.714	<u>0.772</u>	0.794	0.790	<u>0.863</u>	0.819	<u>0.857</u>	0.873	<u>0.871</u>
F1	Nightlight	0.710	0.663	0.664	0.746	0.731	0.814	0.785	0.786	0.853	0.851
	SMOD	<u>0.788</u>	<u>0.786</u>	<u>0.787</u>	0.786	0.737	0.738	0.738	0.738	0.738	0.738
	Population	0.743	0.690	0.690	0.770	0.758	0.884	0.853	0.854	0.896	0.894
	Nightlight + SMOD	0.777	0.718	0.725	0.791	0.792	0.884	0.855	0.856	0.897	0.895
	Nightlight + Population	0.777	0.719	0.721	0.791	0.792	0.884	<u>0.855</u>	0.856	0.897	0.895
	SMOD + Population	0.777	0.719	0.739	0.791	0.792	0.884	0.855	0.858	0.897	0.895
	All	0.776	0.725	0.778	0.800	<u>0.793</u>	<u>0.886</u>	0.854	<u>0.886</u>	0.899	<u>0.896</u>
AUC-ROC	Nightlight	0.742	0.652	0.653	0.749	0.748	0.789	0.717	0.717	0.780	0.784
	SMOD	0.760	<u>0.765</u>	0.765	0.765	0.765	0.765	0.765	0.765	0.765	0.765
	Population	0.765	0.677	0.677	0.768	0.769	0.864	0.807	0.807	0.858	0.860
	Nightlight + SMOD	0.784	0.706	0.715	0.787	0.782	0.865	0.809	0.810	0.859	0.860
	Nightlight + Population	0.784	0.707	0.710	0.787	0.782	0.865	0.809	0.809	0.859	0.860
	SMOD + Population	0.784	0.707	0.733	0.787	0.782	0.865	<u>0.809</u>	0.811	0.859	0.860
	All	<u>0.784</u>	0.714	<u>0.773</u>	0.795	<u>0.790</u>	0.866	0.807	<u>0.845</u>	<u>0.861</u>	<u>0.863</u>

166 4.2 Built-up areas classification

167 Obtaining accurate information about urban built-up areas is crucial for urban planning and man-
 168 agement [34]. In this paper, we investigate the problem of using population, nightlight, and urban
 169 index to classify the urban region functions in the level of $0.01^\circ \times 0.01^\circ$ in space. The experimental
 170 areas of interest are Shanghai and New York State, consistent with the previous section. Five classic
 171 classifiers are chosen for this task: Logistic Regression (LR), Decision Tree (DT), Random Forest
 172 (RF), Gradient Boosting Decision Tree (GBDT) [10] and AdaBoost [9]. To verify the feasibility of
 173 the combination, the accuracy, F1-score (the average harmonic mean of precision and recall), and the
 174 Area Under the Curve (AUC) of the Receiver Operating Characteristic (ROC) are used as the main
 175 metric for the classification tasks.

176 The results are shown in Table 3, from which we have the following observations. (1) Combining
 177 all features achieves the best performance in both regions, which means the nightlight, SMOD (an
 178 indicator showing the degree of urbanization), and population all contribute to the identification of
 179 built-up areas. (2) By further analyzing the SHAP value, we demonstrate the impact of each feature
 180 for individual samples. As observed in Fig. 8 (e), SMOD has more total impact than the other two
 181 features, while for some regions nightlight matters much more. Relation within the data also garners
 182 considerable attention. As depicted in Fig. 8 (a) - (c), SMOD values have a more positive impact on
 183 classification when both SMOD and population values are high in the region. SMOD tends to be
 184 higher when the population is higher, which collectively causes a positive influence. When nightlight
 185 values are the same, the lower the SMOD, the more positive the effect they have on classification.

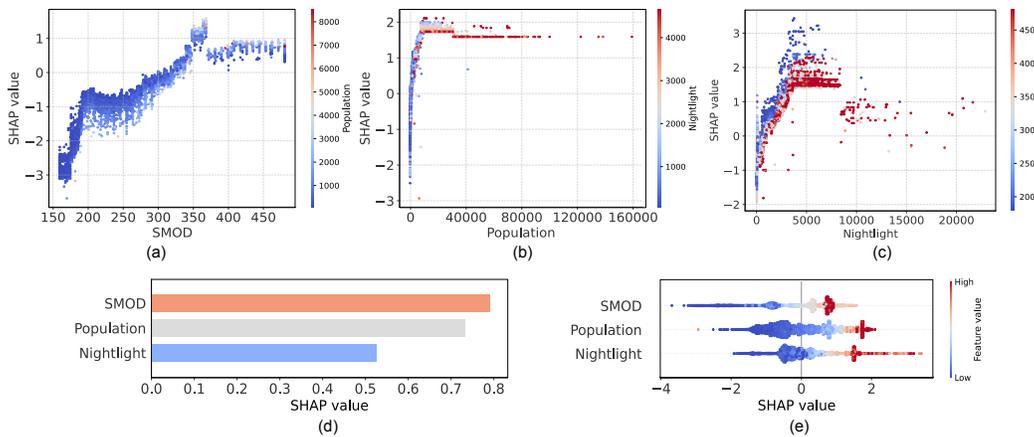


Figure 8: SHAP value analysis of three features for built-up areas classification in Shanghai. (a) - (c) illustrate the interactions between each feature and other features, where each data point represents a sample. In (d), SMOD has the highest mean SHAP value across all given samples, indicating it has the most influence on the results. (e) presents the SHAP values under each feature value, with color representing the level of the feature value.

186 4.3 Identification of administrative boundaries

187 Identifying the boundaries of cities is crucial for urban planning (e.g., infrastructure building) and
 188 urban service arrangement (e.g., delivery). It is believed that using human activity data, e.g., POI,
 189 population, road network data, and nightlight data, can help to identify the city boundary. By utilizing
 190 UDL to unify the aforementioned data to point-wise data, this task can be further formulated as
 191 a clustering problem. Two commonly used clustering methods, K-Nearest Neighbor (KNN) and
 192 Gaussian Mixture Model (GMM) are used, and the clustered boundaries are compared with public
 193 administrative district boundaries. Here, we consider two metrics of this specific task. (1) F1-score:
 194 The F1-score is the harmonic mean of precision and recall. Precision focuses on the number of points
 195 assigned to a district that actually belong to that district while recall is more concerned with how
 196 many points belonging to a district are successfully clustered. (2) IOU: We calculate the Intersection
 197 over Union (IOU) between the obtained clustering boundaries and corresponding administrative
 198 districts.

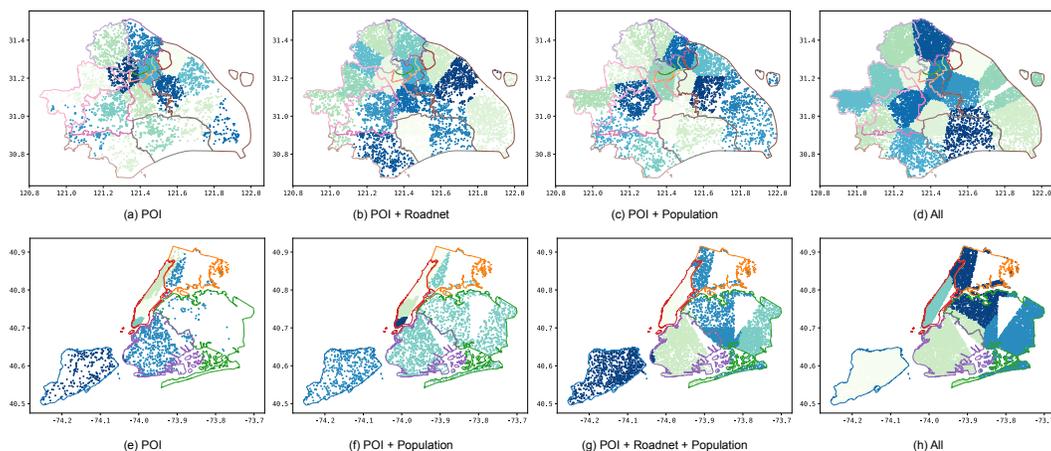


Figure 9: Clustering results in Shanghai using K-means Model ((a) - (d)) and in New York City using Gaussian Mixture Model ((e) - (h)). The x and y coordinates represent latitude and longitude respectively. The points of different colors indicate different clusters predicted, and the polygons of different colors are the ground truth of the administrative divisions.

Table 4: Effectiveness of combining different features in boundary identification problems. The best performance of the combination for each compared method is underlined and the best performance of all is bolded.

City Method Measurement	Shanghai				New York City			
	KNN		GMM		KNN		GMM	
	F1	IOU	F1	IOU	F1	IOU	F1	IOU
POI	0.608	0.338	<u>0.571</u>	0.300	0.557	0.207	0.625	0.303
POI + Roadnet	0.542	0.349	<u>0.497</u>	<u>0.332</u>	0.867	0.330	0.874	0.511
POI + Nightlight	0.499	0.329	0.479	0.294	0.864	0.341	0.771	0.411
POI + Population	0.455	0.231	0.463	0.228	0.467	0.281	0.542	0.306
POI + Roadnet + Nightlight	0.489	0.319	0.473	0.315	0.891	0.370	0.706	0.434
POI + Nightlight + Population	0.509	0.309	0.435	0.279	0.874	0.350	0.702	0.398
POI + Roadnet + Population	0.471	0.286	0.463	0.309	0.914	0.367	<u>0.913</u>	0.577
Combining All	0.475	0.327	0.399	0.280	0.899	<u>0.376</u>	0.909	0.613

Table 5: Effectiveness of combining different features in anomaly detection problems. The best performance of each compared method is bolded and the second best performance is underlined.

Method	EI Nino Dataset						
	LOF	CoLA	ANOMALOUS	GAE	OCGNN	ONE	
AUC-ROC	SP ¹ + ZW ² + MW ³	<u>0.525</u>	<u>0.540</u>	0.469	<u>0.489</u>	0.498	<u>0.469</u>
	SP + Humidity + AT ⁴	0.522	0.450	0.463	0.482	<u>0.496</u>	0.464
	SP + ST ⁵ + AT	<u>0.525</u>	0.542	<u>0.466</u>	0.488	0.493	0.476
	SP + ZW + MW + Humidity + AT	0.540	0.440	<u>0.456</u>	0.499	0.495	0.457
	All	0.538	0.425	0.449	0.478	0.507	0.463

¹ SP: Spatial information contains longitude and latitude.
² ZW: Zonal winds (west < 0, east > 0).
³ MW: Meridional winds (south < 0, north > 0).
⁴ AT: Air temperature.
⁵ ST: Sea surface temperature and subsurface temperatures down to a depth of 500 meters.

199 From Table 4, we observe that using POI information alone achieves the best performance in Shanghai
200 while adding auxiliary data yields better results in New York City. Fig. 9 provides insight into this
201 difference: in Shanghai, POI data effectively differentiates between urban and suburban areas, while
202 the population and road network data distribute more evenly across various districts, which can
203 compromise the distinguishing capability of POI data. Conversely, in New York City, POI data
204 alone is insufficient, and the addition of auxiliary data complements the POI information, leading to
205 improved performance.

206 4.4 EI Nino anomaly detection

207 Detecting urban anomalies (e.g., traffic anomaly, unexpected crowds, environment anomaly, and
208 individual anomaly) holds significant importance in the endeavor to enhance the urban life quality
209 and arrange emergency actions [44]. Here, we use EI Nino dataset as an example to demonstrate
210 how UDL assists in outlier detection tasks. The original dataset is assumed to be without anomalies.
211 Following the approach in [8], we introduce anomalies constituting 2% of the dataset. We then
212 compare the performance of different combinations of node features using various anomaly detection
213 methods [24], including LOF [4], CoLA [25], ANOMALOUS [33], GAE [16], OCGNN [38] and
214 ONE [3]. The evaluation metric utilized is the Area Under the Curve (AUC) of the Receiver Operating
215 Characteristic (ROC).

216 We show AUC values for all methods on all feature combinations in Table. 5. It is observed that the
217 combination of spatial information, zonal winds, and meridional winds achieves relatively better
218 results overall. The best combination of results is sea surface temperature and air temperature using
219 CoLA. Moreover, we shed light on some interesting observations regarding the results to explain why
220 it is more likely to be an outlier. In Fig. 10 (a), considering the properties of air temperature and sea
221 surface temperature, the outlier is similar to its neighbors in one of the attributes while another is
222 much higher or lower. We can observe that the detected anomaly's (upper left) air temperature is
223 around 26°. But its sea surface temperature is higher than 28.5° where its "neighboring" samples
224 with the same air temperature are below 28°. Similar observations in structural aspects can be made

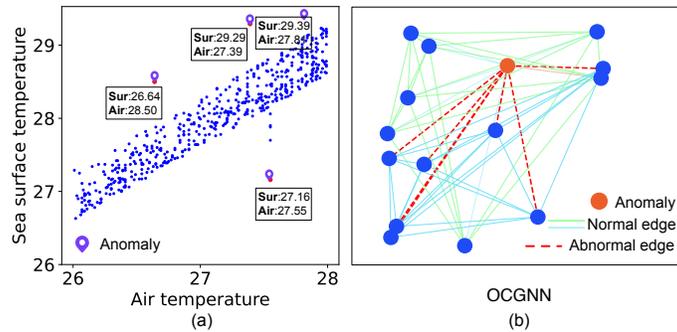


Figure 10: The detected anomaly and surrounding points of El Nino region. (a) Detected anomaly points by CoLA. (b) Detected structural anomaly nodes by OCGNN.

225 in Fig. 10 (b), where an anomaly may be a node whose edges are inaccurately linked. As the edges
 226 are established based on spatio-temporal information, edge relationships exist between nodes that
 227 have the same temporal or spatial information. The node in the graph is recognized as an anomaly
 228 because the spatio-temporal feature is replaced, making the edges in the dotted line unusually present.
 229 Since the features of the anomalies are replaced randomly, the best combination of features may be
 230 stochastic.

231 5 Conclusion and Outlook

232 This paper introduces a unified data pipeline including standard data structures and easy-to-use
 233 processing interfaces on urban research. We define the standard data layers from five common data
 234 organizations used in urban science and provide three components in the pipeline. UDL mitigates
 235 the gap between various urban data and urban computing research by addressing the challenges:
 236 (1) handling dirty and repetitive data processing, (2) establishing a unified standardized format, (3)
 237 integrating alignment and fusion for urban data. This will enable reproducible benchmark construction
 238 and foster the development of the multi-modal databases. The effectiveness and productivity of UDL
 239 have been demonstrated in four instances. We believe it will become a promising data tool to inspire
 240 more researchers to tackle the urban problems our cities face.

241 When data layers are constructed globally, the availability of sufficient data facilitates large-scale
 242 urban research and the development of large models [5, 14]. Despite the high productivity of UDL,
 243 the alignment of urban data is currently limited to geospatial information and future research could
 244 explore more aspects. In the future, we will incorporate tasks across regions and explore solutions to
 245 urban issues on a global scale.

246 Acknowledgments and Disclosure of Funding

247 This work was sponsored by National Natural Science Foundation of China under Grant No.
 248 62102246, 62272301, and Provincial Key Research and Development Program of Zhejiang under
 249 Grant No. 2021C01034. Part of the work was done when the students were doing internships at
 250 Yunqi Academy of Engineering.

References

- [1] Open geospatial consortium. <https://ogcapi.ogc.org/>.
- [2] Kumar Ayush, Burak Uz kent, Kumar Tanmay, Marshall Burke, David Lobell, and Stefano Ermon. Efficient poverty mapping from high resolution remote sensing images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12–20, 2021.
- [3] Sambaran Bandyopadhyay, N Lokesh, and M Narasimha Murty. Outlier aware network embedding for attributed networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 12–19, 2019.
- [4] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. Lof: identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pages 93–104, 2000.
- [5] Shengchao Chen, Guodong Long, Tao Shen, and Jing Jiang. Prompt federated learning for weather forecasting: Toward foundation models on meteorological data. *arXiv preprint arXiv:2301.09152*, 2023.
- [6] Weiyu Cheng, Yanyan Shen, Yanmin Zhu, and Linpeng Huang. A neural attention model for urban air quality inference: Learning the weights of monitoring stations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [8] Kaize Ding, Jundong Li, Rohit Bhanushali, and Huan Liu. Deep anomaly detection on attributed networks. In *Proceedings of the 2019 SIAM International Conference on Data Mining*, pages 594–602. SIAM, 2019.
- [9] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences*, 55(1):119–139, 1997.
- [10] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001.
- [11] Shanghua Gao, Teddy Koker, Owen Queen, Thomas Hartvigsen, Theodoros Tsiligkaridis, and Marinka Zitnik. Units: Building a unified time series model. *arXiv preprint arXiv:2403.00131*, 2024.
- [12] Jindong Han, Hao Liu, Hengshu Zhu, Hui Xiong, and Dejing Dou. Joint air quality and weather prediction based on multi-adversarial spatiotemporal networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 4081–4089, 2021.
- [13] Sungwon Han, Donghyun Ahn, Hyunji Cha, Jeasurk Yang, Sungwon Park, and Meeyoung Cha. Lightweight and robust representation of economic scales from satellite imagery. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 428–436, 2020.
- [14] Langwen Huang and Torsten Hoeffler. Compressing multidimensional weather and climate data into neural networks. *arXiv preprint arXiv:2210.12538*, 2022.
- [15] Renhe Jiang, Xuan Song, Zipei Fan, Tianqi Xia, Quanjun Chen, Satoshi Miyazawa, and Ryosuke Shibasaki. Deepurbanmomentum: An online deep-learning system for short-term urban mobility prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [16] Thomas N Kipf and Max Welling. Variational graph auto-encoders. *arXiv preprint arXiv:1611.07308*, 2016.

- 294 [17] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
295 2009.
- 296 [18] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural
297 network: Data-driven traffic forecasting. *arXiv preprint arXiv:1707.01926*, 2017.
- 298 [19] Zhili Li, Yiqun Xie, Xiaowei Jia, Kara Stuart, Caroline Delaire, and Sergii Skakun. Point-to-
299 region co-learning for poverty mapping at high resolution using satellite imagery. In *Proceedings*
300 *of the AAAI Conference on Artificial Intelligence*, volume 37, pages 14321–14328, 2023.
- 301 [20] Yuxuan Liang, Yutong Xia, Songyu Ke, Yiwei Wang, Qingsong Wen, Junbo Zhang, Yu Zheng,
302 and Roger Zimmermann. Airformer: Predicting nationwide air quality in china with trans-
303 formers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages
304 14329–14337, 2023.
- 305 [21] Ziqian Lin, Jie Feng, Ziyang Lu, Yong Li, and Depeng Jin. Deepstn+: Context-aware spatial-
306 temporal neural network for crowd flow prediction in metropolis. In *Proceedings of the AAAI*
307 *conference on artificial intelligence*, volume 33, pages 1020–1027, 2019.
- 308 [22] Jeremiah Liu, John Paisley, Marianthi-Anna Kioumourtzoglou, and Brent Coull. Accurate
309 uncertainty estimation and decomposition in ensemble learning. *Advances in neural information*
310 *processing systems*, 32, 2019.
- 311 [23] Jia Liu, Tianrui Li, Peng Xie, Shengdong Du, Fei Teng, and Xin Yang. Urban big data fusion
312 based on deep learning: An overview. *Information Fusion*, 53:123–133, 2020.
- 313 [24] Kay Liu, Yingtong Dou, Yue Zhao, Xueying Ding, Xiyang Hu, Ruitong Zhang, Kaize Ding,
314 Canyu Chen, Hao Peng, Kai Shu, George H. Chen, Zhihao Jia, and Philip S. Yu. Pygod: A
315 python library for graph outlier detection. *arXiv preprint arXiv:2204.12095*, 2022.
- 316 [25] Yixin Liu, Zhao Li, Shirui Pan, Chen Gong, Chuan Zhou, and George Karypis. Anomaly
317 detection on attributed networks via contrastive self-supervised learning. *IEEE transactions on*
318 *neural networks and learning systems*, 33(6):2378–2392, 2021.
- 319 [26] Yonghong Luo, Xiangrui Cai, Ying Zhang, Jun Xu, et al. Multivariate time series imputation
320 with generative adversarial networks. *Advances in neural information processing systems*, 31,
321 2018.
- 322 [27] Andrew Kachites McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Automating
323 the construction of internet portals with machine learning. *Information Retrieval*, 3:127–163,
324 2000.
- 325 [28] Chenlin Meng, Enci Liu, Willie Neiswanger, Jiaming Song, Marshall Burke, David Lobell,
326 and Stefano Ermon. Is-count: large-scale object counting from satellite images with covariate-
327 based importance sampling. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
328 volume 36, pages 12034–12042, 2022.
- 329 [29] Chuishi Meng, Yanhua Li, Yu Zheng, Jieping Ye, Qiang Yang, Philip S Yu, and Ouri Wolfson.
330 The 12th international workshop on urban computing. In *Proceedings of the 29th ACM SIGKDD*
331 *Conference on Knowledge Discovery and Data Mining*, pages 5874–5875, 2023.
- 332 [30] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
333 models. *arXiv preprint arXiv:1609.07843*, 2016.
- 334 [31] Yuqi Nie, Nam H Nguyen, Phanwadee Sinthong, and Jayant Kalagnanam. A time series is
335 worth 64 words: Long-term forecasting with transformers. *arXiv preprint arXiv:2211.14730*,
336 2022.

- 337 [32] Shailesh Pandey, Tushar Agarwal, and Narayanan C Krishnan. Multi-task deep learning for
338 predicting poverty from satellite images. In *Proceedings of the AAAI Conference on Artificial*
339 *Intelligence*, volume 32, 2018.
- 340 [33] Zhen Peng, Minnan Luo, Jundong Li, Huan Liu, Qinghua Zheng, et al. Anomalous: A joint
341 modeling approach for anomaly detection on attributed networks. In *IJCAI*, pages 3513–3519,
342 2018.
- 343 [34] Lang Sun, Lina Tang, Guofan Shao, Quanyi Qiu, Ting Lan, and Jinyuan Shao. A machine
344 learning-based classification system for urban built-up areas using multiple classifiers and data
345 sources. *Remote Sensing*, 12(1):91, 2019.
- 346 [35] Thomas Vandal, Evan Kodra, Sangram Ganguly, Andrew Michaelis, Ramakrishna Nemani, and
347 Auroop R Ganguly. DeepSD: Generating high resolution climate change projections through
348 single image super-resolution. In *Proceedings of the 23rd acm sigkdd international conference*
349 *on knowledge discovery and data mining*, pages 1663–1672, 2017.
- 350 [36] Jinyuan Wang, Jiawei Jiang, Wenjun Jiang, Chengkai Han, and Wayne Xin Zhao. Towards ef-
351 ficient and comprehensive urban spatial-temporal prediction: A unified library and performance
352 benchmark. *arXiv preprint arXiv:2304.14343*, 2023.
- 353 [37] Jinyuan Wang, Jiawei Jiang, Wenjun Jiang, Chao Li, and Wayne Xin Zhao. Libcity: An open
354 library for traffic prediction. In *Proceedings of the 29th international conference on advances*
355 *in geographic information systems*, pages 145–148, 2021.
- 356 [38] Xuhong Wang, Baihong Jin, Ying Du, Ping Cui, Yingshui Tan, and Yupu Yang. One-class
357 graph neural networks for anomaly detection in attributed networks. *Neural computing and*
358 *applications*, 33:12073–12085, 2021.
- 359 [39] Zhecheng Wang, Haoyuan Li, and Ram Rajagopal. Urban2vec: Incorporating street view
360 imagery and pois for multi-modal urban neighborhood embedding. In *Proceedings of the AAAI*
361 *Conference on Artificial Intelligence*, volume 34, pages 1013–1020, 2020.
- 362 [40] Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen
363 Sahoo. Unified training of universal time series forecasting transformers. *arXiv preprint*
364 *arXiv:2402.02592*, 2024.
- 365 [41] Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang.
366 Connecting the dots: Multivariate time series forecasting with graph neural networks. In
367 *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery &*
368 *data mining*, pages 753–763, 2020.
- 369 [42] Michael Xie, Neal Jean, Marshall Burke, David Lobell, and Stefano Ermon. Transfer learning
370 from deep features for remote sensing and poverty mapping. In *Thirtieth AAAI Conference on*
371 *Artificial Intelligence*, 2016.
- 372 [43] Huaxiu Yao, Xianfeng Tang, Hua Wei, Guanjie Zheng, and Zhenhui Li. Revisiting spatial-
373 temporal similarity: A deep learning framework for traffic prediction. In *Proceedings of the*
374 *AAAI conference on artificial intelligence*, volume 33, pages 5668–5675, 2019.
- 375 [44] Mingyang Zhang, Tong Li, Yue Yu, Yong Li, Pan Hui, and Yu Zheng. Urban anomaly analytics:
376 Description, detection, and prediction. *IEEE Transactions on Big Data*, 8(3):809–826, 2020.
- 377 [45] X. Zhang, C. Huang, Y. Xu, L. Xia, P. Dai, L. Bo, J. Zhang, and Y. Zheng. Traffic flow
378 forecasting with spatial-temporal graph diffusion network. In *AAAI*, 2021.
- 379 [46] Yu Zheng. Methodologies for cross-domain data fusion: An overview. *IEEE transactions on*
380 *big data*, 1(1):16–34, 2015.

- 381 [47] Yu Zheng, Licia Capra, Ouri Wolfson, and Hai Yang. Urban computing: concepts, method-
382 ologies, and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*,
383 5(3):1–55, 2014.
- 384 [48] Yu Zheng, Yuming Lin, Liang Zhao, Tinghai Wu, Depeng Jin, and Yong Li. Spatial planning
385 of urban communities via deep reinforcement learning. *Nature Computational Science*, pages
386 1–15, 2023.
- 387 [49] Yu Zheng, Xiuwen Yi, Ming Li, Ruiyuan Li, Zhangqing Shan, Eric Chang, and Tianrui Li.
388 Forecasting fine-grained air quality based on big data. In *Proceedings of the 21th ACM SIGKDD*
389 *international conference on knowledge discovery and data mining*, pages 2267–2276, 2015.
- 390 [50] Haoyi Zhou, Shanghang Zhang, Jieqi Peng, Shuai Zhang, Jianxin Li, Hui Xiong, and Wancai
391 Zhang. Informer: Beyond efficient transformer for long sequence time-series forecasting. In
392 *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 11106–11115,
393 2021.

394 **Checklist**

- 395 1. For all authors...
- 396 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
397 contributions and scope? [Yes] See Section 1.
- 398 (b) Did you describe the limitations of your work? [Yes] See Section 5.
- 399 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 400 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
401 them? [Yes]
- 402 2. If you are including theoretical results...
- 403 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 404 (b) Did you include complete proofs of all theoretical results? [N/A]
- 405 3. If you ran experiments (e.g. for benchmarks)...
- 406 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
407 imental results (either in the supplemental material or as a URL)? [Yes] We provide
408 related document and code.
- 409 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
410 were chosen)? [Yes] We include the experiment details in supplemental material.
- 411 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
412 ments multiple times)? [No]
- 413 (d) Did you include the total amount of compute and the type of resources used (e.g., type
414 of GPUs, internal cluster, or cloud provider)? [No]
- 415 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 416 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 4.
- 417 (b) Did you mention the license of the assets? [No]
- 418 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
419 We include detailed anomaly detection models in supplemental material.
- 420 (d) Did you discuss whether and how consent was obtained from people whose data you're
421 using/curating? [N/A]
- 422 (e) Did you discuss whether the data you are using/curating contains personally identifiable
423 information or offensive content? [N/A]
- 424 5. If you used crowdsourcing or conducted research with human subjects...
- 425 (a) Did you include the full text of instructions given to participants and screenshots, if
426 applicable? [N/A]
- 427 (b) Did you describe any potential participant risks, with links to Institutional Review
428 Board (IRB) approvals, if applicable? [N/A]
- 429 (c) Did you include the estimated hourly wage paid to participants and the total amount
430 spent on participant compensation? [N/A]