

Towards Better Detection and Analysis of Massive Spatiotemporal Co-Occurrence Patterns

Yingcai Wu¹, Di Weng¹, Zikun Deng, Jie Bao, Mingliang Xu¹, Zhangye Wang,
Yu Zheng¹, Zhiyu Ding, and Wei Chen¹

Abstract—With the rapid development of sensing technologies, massive spatiotemporal data have been acquired from the urban space with respect to different domains, such as transportation and environment. Numerous co-occurrence patterns (e.g., traffic speed < 10km/h, weather = foggy, and air quality = unhealthy) between the transportation data and other types of data can be obtained with given spatiotemporal constraints (e.g., within 3 kilometers and lasting for 2 hours) from these heterogeneous data sources. Such patterns present valuable implications for many urban applications, such as traffic management, pollution diagnosis, and transportation planning. However, extracting and understanding these patterns is beyond manual capability because of the scale, diversity, and heterogeneity of the data. To address this issue, a novel visual analytics system called *CorVizor* is proposed to identify and interpret these co-occurrence patterns. *CorVizor* comprises two major components. The first component is a co-occurrence mining framework involving three steps, namely, spatiotemporal indexing, co-occurring instance generation, and pattern mining. The second component is a visualization technique called *CorView* that implements a level-of-detail mechanism by integrating tailored visualizations to depict the extracted spatiotemporal co-occurrence patterns. The case studies and expert interviews are conducted to demonstrate the effectiveness of *CorVizor*.

Index Terms—Heterogeneous urban data, spatiotemporal data visualization, co-occurrence pattern analysis.

Manuscript received April 12, 2019; revised October 9, 2019 and January 13, 2020; accepted February 24, 2020. The work was supported in part by the NSFC-Zhejiang Joint Fund for the Integration of Industrialization and Informatization under Grant U1609217, in part by the National Key Research and Development Program of China under Grant 2018YFB1004300 and Grant 2017YFB1002703, in part by the NSFC under Grant 61761136020 and Grant 61822701, in part by the Zhejiang Provincial Natural Science Foundation under Grant LR18F020001, in part by the 100 Talents Program of Zhejiang University, and in part by the Microsoft Research Asia. The Associate Editor for this article was A. Hegyi. (Corresponding authors: Mingliang Xu; Zhangye Wang.)

Yingcai Wu, Di Weng, Zikun Deng, and Wei Chen are with the State Key Lab of CAD&CG, Zhejiang University, Hangzhou 310027, China, and also with the Zhejiang Lab, Hangzhou 311122, China (e-mail: ycwu@zju.edu.cn; dweng@zju.edu.cn; zikun_rain@zju.edu.cn; chenvis@zju.edu.cn).

Jie Bao and Yu Zheng are with the Urban Computing Business Unit, JD Finance, Beijing 101111, China (e-mail: baojie1985@gmail.com; msyuzheng@outlook.com).

Mingliang Xu is with the School of Information Engineering, Zhengzhou University, Zhengzhou 451200, China, and also with the Henan Institute of Advanced Technology, Zhengzhou University, Zhengzhou 451200, China (e-mail: iexumingliang@zzu.edu.cn).

Zhangye Wang is with the State Key Lab of CAD&CG, Zhejiang University, Hangzhou 310027, China (e-mail: zywang@cad.zju.edu.cn).

Zhiyu Ding is with Cloud BU, Huawei Technologies Company, Ltd., Hangzhou 310053, China (e-mail: dingzhiyu@huawei.com).

This article has supplementary downloadable material available at <http://ieeexplore.ieee.org>, provided by the authors.

Digital Object Identifier 10.1109/TITS.2020.2983226

I. INTRODUCTION

THE rapid development of sensing technologies has resulted in a large amount of heterogeneous urban data acquired from different data sources, such as traffic and air quality data. These data inherently comprise numerous interesting *co-occurrence patterns*, i.e., the combinations of the property value ranges that frequently co-occur with each other. These patterns appear frequently within a spatial range and a temporal window and may comprise properties from various data sources. For example, given three data sources, namely, transportation, weather, and air quality, and spatiotemporal constraints (within 3-kilometer range and 2-hour window), a fine-grained co-occurrence pattern like $\{20 < \text{TrafficVolume} < 30, 100\text{m/s} < \text{WindSpeed} < 150\text{m/s}, \text{AirQuality} = \text{healthy}\}$ may be identified. These fine-grained patterns reveal important spatiotemporal insights and anomalies (i.e., counterintuitive co-occurrence patterns) across multiple data sources that support numerous urban decision-making applications, including traffic management and transportation planning.

However, neither have such co-occurrence patterns been systematically studied and detected, nor effectively interpreted and understood. Two challenges arise from the identification and interpretation of these patterns: a) **efficient extraction** and b) **interactive visualization**.

Efficient extraction of the co-occurrence patterns is considerably difficult, particularly from the heterogeneous urban data sources that comprise various properties, such as PM2.5 and PM10 in air quality data and temperature and humidity in meteorological data. Without an efficient approach, exhaustively testing all possible combinations of these properties to find the potential patterns will result in poor computational performance as the number of properties increases. Furthermore, the aforementioned patterns are fine-grained on continuous value domains and may involve both categorical and numerical properties. However, most of the traditional co-occurrence mining techniques, such as Apriori [1], [2], are specifically designed to detect the coarse association rules (i.e., co-occurrence patterns), such as $\{\text{butter}, \text{bread}\} \Rightarrow \{\text{milk}\}$, among categorical attributes only. Some recent techniques [24], [50] that extract patterns from numerical data generally require the continuous domains of property values to be initially discretized, thereby resulting in the severe loss of latent patterns.

The extracted co-occurrence patterns also require a well-designed visualization technique, with which domain experts can examine these frequent co-occurrences, identify spatiotemporal trends, and study anomalies. However, it is difficult to develop an appropriate technique that visualizes these patterns because of three identified design challenges: a) **diversity**: the data properties involved in the visualized co-occurrence patterns may have different types, scales, and semantics; b) **volume**: numerous patterns with overlapping value ranges can be extracted from heterogeneous urban data; c) **organization**: a combination of properties can be shared by many patterns, thereby forming a two-level hierarchical structure where the designed visualizations should enable experts to explore the hierarchy flexibly and efficiently. To the best of our knowledge, none of the existing visualization approaches, including parallel coordinates and parallel sets [25], can be applied directly to address these challenges, which demand a set of considerate visualizations specifically tailored on the basis of the unique characteristics of the extracted patterns.

In this study, we develop a novel data mining model that extracts co-occurrence patterns from massive urban data based on three main modules: a) *spatiotemporal indexing* that builds a unified index structure to accelerate the following mining process, b) *co-occurring instance generation* that identifies co-occurring instances and builds a pruning graph to reduce the search space of patterns, and c) *pattern mining* that aggregates the data by using *value matrices* and extracts distinct patterns via an improved *sweep-line* algorithm. We also propose a new matrix-based visualization technique named *CorView*, which effectively depicts the extracted patterns that comprise properties of different data types, scales, and semantics in an aligned fashion. Particularly, we address the scalability issue by adopting a level-of-detail mechanism that integrates brick-like glyphs, scatterplots, parallel coordinates, and a stacked line chart. Furthermore, we design *CorVizor*, a novel visual analytics system that helps users reliably detect and analyze the co-occurrence patterns in cross-domain urban data. The major contributions of this study are as follows.

- ◊ We characterize the problem of identifying and interpreting the fine-grained spatiotemporal co-occurrence patterns among the cross-domain urban data sources;
- ◊ We develop *CorVizor*, a visual analytics system that integrates a pattern mining framework and a multi-scale visual representation to assist experts in detecting, exploring, and interpreting the co-occurrence patterns effectively;
- ◊ We evaluate the proposed system with the case studies conducted on the real-world data, where the co-occurrence patterns among traffic status, air pollution, and weather conditions are studied.

II. RELATED WORK

This section discusses related studies in the following three parts, namely, co-occurrence pattern mining, spatiotemporal visualization, and co-occurrence visualization.

A. Co-Occurrence Pattern Mining

Co-occurrence pattern mining techniques were proposed to identify frequent co-occurrences among categorical or numerical data properties.

Traditional techniques like Apriori [1], [2] and PrefixSpan [3] attempt to extract the patterns from transactional datasets and thus were limited to handling categorical data. Similar methods have been adapted to transactional urban datasets, where the extraction of co-occurrence patterns from the moving object [10], boolean [35], [59], and event [7], [30], [37] datasets in spatiotemporal contexts were extensively studied. However, these techniques cannot be directly applied to solve our problem since most of the properties like traffic speed and volume in heterogeneous urban data have continuous domains.

Techniques [24], [39], [50] were also proposed to handle numerical data by dividing the continuous domains of properties into a number of bins. However, such discretization may lead to the severe loss of latent patterns. Other studies have attempted to avoid the discretization with topological methods in spatiotemporal contexts. Chirigati *et al.* [14] developed a topology-based method named data polygamy. This method efficiently captures the relationships between extrema in urban datasets. Nonetheless, it can only identify the co-occurrences that comprise the peaks or valleys of data properties.

We define co-occurrence patterns as the flexible value range combinations of both numerical and categorical properties in heterogeneous urban datasets. Based on this definition, our model can extract the fine-grained patterns efficiently without discretization of the domains.

B. Spatiotemporal Visualization

The rapid development of smart cities enables authorities to collect citywide spatiotemporal data via sensors more efficiently than ever, making data-driven solutions possible for urbanization problems like air pollution [16], traffic congestions [56], and transit route planning [51]. To integrate human in the analysis loop, spatiotemporal data visualization has been investigated and applied in many settings, such as billboard location selection [28], public utility analysis [58], home location selection [52], and hotspot prediction [31]. Andrienko *et al.* [4] provided an excellent taxonomy of existing spatiotemporal visualization methods for movement data, which are classified into three categories, namely, direct depiction (e.g. points [20], polylines [5], stacked bands [45], and space-time cubes [6]), summarization (e.g. density maps [41], [53], graphs [46], and flow maps [21]), and pattern extraction [12], [23], [55]. Sun *et al.* [43] also explored the better integration of temporal information in spatial contexts by transforming maps. To handle large-scale spatiotemporal data, many novel methods have been incorporated into visualizations, such as tailored query model [20], topological methods [17], [34], uncertainty analysis [12], and anomaly detection [9]. However, most of the prior studies focus on single-source data only, including trajectory [40], cellphone [55], [60], and weather data [38]. In contrast, our study targets at visualizing multi-source heterogeneous data,

which poses difficult design challenges arisen from the unique characteristics of co-occurrence patterns.

This work establishes a pattern extraction method that aims to detect and visualize an extensive number of fine-grained co-occurrence patterns among heterogeneous datasets. In particular, various types of urban data from multiple domains were analyzed and explored through a mining model and a set of tailored visualization techniques.

C. Co-Occurrence and Correlation Visualization

Visually understanding and analyzing the massive extracted co-occurrence patterns remain a difficult and challenging task. Many co-occurrence and correlation visualization methods targeting categorical data have been proposed based on scalable techniques like 2D plots [27], [29], graphs [18], parallel coordinates [57], and matrices [22], [54]. For numerical datasets, Bothorel *et al.* [8] proposed a visual mining pipeline based on the Apriori algorithm, yet the value ranges must be discretized.

Recently, the visual analytics of spatiotemporal co-occurrence and correlation patterns has attracted wide research interests. Qu *et al.* [38] studied the visualization of correlation between various weather attributes. TelCoVis [55] was designed to illustrate the human co-occurrence patterns with mobile phone data. Furthermore, a few studies have considered the complex co-occurrence and correlation patterns among multiple data sources. Urbane [19] combines datasets from diverse domains for target building selection. VAUD [13] allows users to explore cross-domain correlation based on visual query and reasoning. COPE [26] detects various co-occurrence patterns of spatiotemporal events via a well-designed visual interface. However, most of these visualization techniques neither integrate with an automated mining model nor involve the co-occurrence patterns characterized by combinations of continuous value ranges and categorical sets. Thus, finding and interpreting interesting patterns will become increasingly difficult with the growing size of datasets.

In this paper, we design a novel analytics system that combines several interactive visualizations specifically tailored for the massive fine-grained co-occurrence patterns extracted by the proposed model in spatiotemporal contexts.

III. BACKGROUND AND SYSTEM OVERVIEW

This section presents the background, problem, and over-view of the proposed system.

A. Background

Our study mainly focuses on extracting, visualizing, and evaluating frequent spatiotemporal co-occurrence patterns obtained from heterogeneous urban data. We introduce the following terminologies in the extraction of spatiotemporal co-occurrence patterns. For each annotation, the superscript is used to distinguish different objects, and the subscript is to indicate the association of the current object.

- ◇ **Data source:** A data source $s \in S = \{s^1, s^2, \dots, s^n\}$ comprises a set of spatial locations $\{l^1, l^2, \dots\} \in L$, each of which is associated with a set of time-varying

TABLE I
PROPERTIES OF THREE DATA SOURCES

Data Source	Property	Value Range
Air Quality	PM 2.5, PM 10 (ug/m ³) ¹	[0, 500]
	O ₃ (ug/m ³)	[0, 2300]
	NO ₂ (ug/m ³)	[0, 300]
	CO (ug/m ³)	[0, 70]
	SO ₂ (ug/m ³)	[0, 150]
	AQI Level	6 levels
Weather	Temperature °C	[-20, 40]
	Humidity	[0, 100]
	Wind Speed (m/s)	[0, 300]
	Wind Direction	[1, 24]
	Cloud Conditions	14 conditions
Traffic	Total Cars	[9, 200]
	Low Speed (0 – 20 km/h) %	[6, 80]
	Medium Speed (20 – 50 km/h) %	[17, 74]
	High Speed (above 50 km/h) %	[0, 53]

¹ PM2.5 is the particulate matter measured 2.5 micrometers or less in the diameter. Similarly, PM10 is the particulate matter measured 10 micrometers or less in the diameter.

properties $\{p^1, p^2, \dots\}$ observed at the location. The data sources used in this study are shown in Table I with their properties.

- ◇ **Instance:** An instance φ , associated with a data source s , comprises a spatial location l , a timestamp t , and a value v_p of the property p observed at the time t and location l . For example, Fig. 1(a) presents six instances: $\varphi_{s^a}^1$ and $\varphi_{s^a}^2$ from the property p^1 of the data source s^a , $\varphi_{s^b}^1$ and $\varphi_{s^b}^2$ from the property p^1 of the data source s^b , and $\varphi_{s^\tau}^1$ and $\varphi_{s^\tau}^2$ from the property p^1 of the data source s^τ .
- ◇ **Property value range:** We denote a specific value range of the property p in the data source s as $s_p|\mathcal{C}$. The range \mathcal{C} can either be numeric (e.g. [5, 8]) or ordinal (e.g. {cloudy}), depending on the type of the property p . We say an instance φ satisfies a property value range $s_p|\mathcal{C}$ iff. (1) φ is collected from the property p and (2) the observed value v_p of the property p is within the range of \mathcal{C} , i.e., $v_p \in \mathcal{C}$. For example, in Fig. 1(a), the value range of the property p^1 in the data source s^τ can be written as $s_{p^1}^\tau|[5, 8]$.
- ◇ **Co-occurring instances:** The instances co-occurring w.r.t. a given target instance within the user-specified spatial and temporal thresholds are defined to be co-occurring instances. Fig. 1(a) shows an example of five instances $\varphi_{s^\tau}^2$, $\varphi_{s^a}^1$, $\varphi_{s^a}^2$, $\varphi_{s^b}^1$, and $\varphi_{s^b}^2$ of the data sources s^τ , s^a , and s^b that co-occur w.r.t. the instance $\varphi_{s^\tau}^1$ within the spatial and temporal thresholds d and t .
- ◇ **Co-occurring property value ranges:** By extracting and combining co-occurring instances, we can obtain the co-occurring property value ranges comprising the values of these instances w.r.t. a target property. Fig. 1(b) shows an example of the co-occurring value ranges of the properties in the data sources s^a and s^b w.r.t. the property in the target data source s^τ .
- ◇ **Co-occurrence coverage:** The co-occurrence coverage of a combination of co-occurring property value ranges

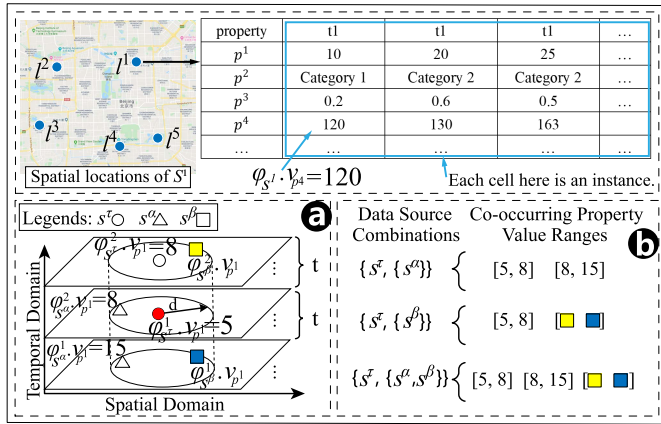


Fig. 1. (a) Co-occurring instances and (b) examples of co-occurring property value ranges.

w.r.t. the target data source s^t is defined as the number of the co-occurring instances in the given property value ranges involving s^t divided by the total number of the instances in s^t (denoted by $|s^t|$).

- ◇ **Co-occurrence pattern:** A co-occurrence pattern is defined as a combination of co-occurring property value ranges, the co-occurrence coverage of which is higher than a given threshold w.r.t. the target data source s^t . A pattern can be written as $\{s_{p^u}^t | C, \{s_{p^x}^a | C, s_{p^y}^b | C, \dots\}\}$, where $s_{p^u}^t | C$, $s_{p^x}^a | C$ and $s_{p^y}^b | C$ are the co-occurring property value ranges w.r.t. the target property $s_{p^u}^t$.

B. Problem Definition

Given a target data source s^t , a group of data sources $S = \{s^1, s^2, \dots, s^n\}$, and a set of mining parameters, including the spatial distance d , temporal window t , and co-occurrence coverage threshold λ_e , the objective is to identify all distinctive co-occurrence patterns efficiently and comprehend them via interactive visualizations.

C. System Overview

We develop CorVizor, a web-based visual analytics system that can assist urban experts in interpreting and analyzing co-occurrence patterns extracted from heterogeneous urban data. CorVizor comprises two components: data mining and interactive visualization. The mining component, implemented in C#, indexes heterogeneous spatiotemporal data, extracts co-occurring instances, and performs pattern mining, transforming raw data into interpretable patterns. The visualization component, implemented in TypeScript, visualizes co-occurrence patterns with four tailored views, thereby enabling users to interactively filter, data compare, and evaluate the patterns across multiple data sources.

IV. MINING FRAMEWORK

Detecting fine-grained co-occurrence patterns from heterogeneous urban data sources is difficult because spatiotemporal urban data generally do not have transactions. Moreover,

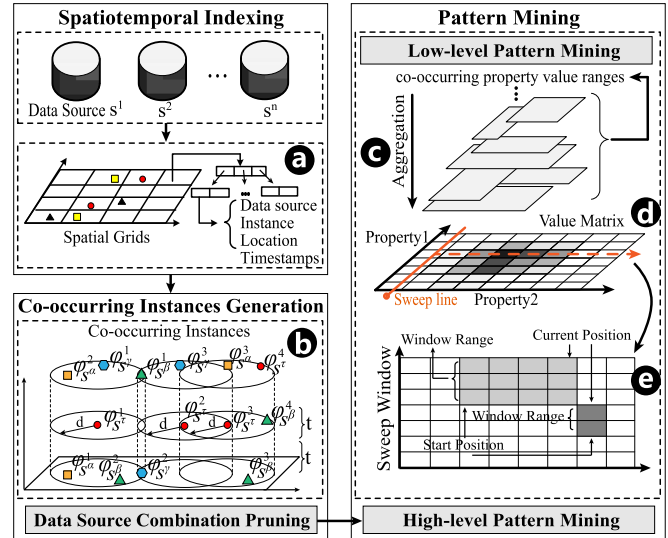


Fig. 2. Overview of the mining framework.

the data sources may comprise diverse value scales, including numeric and categorical scales. Furthermore, characterizing the co-occurrence patterns as the combinations of continuous property value ranges and categories easily leads to a huge combinatorial solution space of the possible co-occurrence patterns. Given these challenges, traditional mining techniques, such as Apriori [2], cannot be applied directly to address our problem. Thus, we propose a mining framework with the following three modules (Fig. 2) to tackle the challenges:

- 1) **indexing spatiotemporal data** with a nested data structure to accelerate the subsequent mining process;
- 2) **generating co-occurring instances** by identifying these instances with the unified indexes and pruning impossible co-occurring instances;
- 3) **mining co-occurrence patterns** with a novel sweep-line algorithm based on the *value matrices* constructed from the generated instances.

A. Indexing Spatiotemporal Data

At this stage, we build a unified index [11], [15], [48] for large-scale heterogeneous spatiotemporal data (Fig. 2(a)) to enable faster data retrieval with different mining parameters specified and accelerate the subsequent mining stages. To construct the index, we divide the map into $n \times m$ spatial grids, each of which has an area of 1km^2 . Each grid maintains the covered instances with a temporal index, where the instances are organized by their timestamps with a B^+ tree. Each leaf node of the B^+ tree records the ID, data source ID, GPS location, and timestamp of an instance.

B. Generating Co-Occurring Instances

With the target data source s^t , spatial threshold d , temporal threshold t , and expected co-occurrence coverage λ_e , the proposed mining framework attempts to extract co-occurring instances with *co-occurrence tables* and *pruning graphs* (Fig. 3) from the spatiotemporal index built at the previous stage.

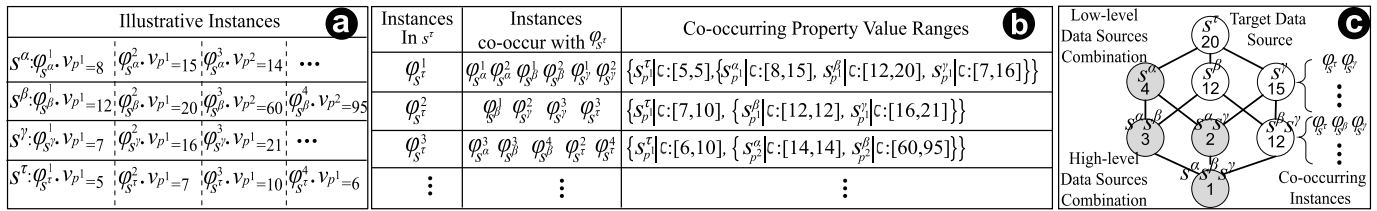


Fig. 3. (a) The instances extracted from data sources (one row for each data source), (b) a co-occurrence table (20 co-occurring instances, each described by a row), and (c) a pruning graph built with 20 co-occurring instances.

First, range queries based on the spatial and temporal thresholds are issued for each instance of target data source s^t to identify co-occurring instances within the same property of the target data source or in other data sources. Based on the given spatial distance and temporal range, the co-occurring instances (Fig. 3(a)) in a spatiotemporal cylinder of each instance in s^t (Fig. 2(b)) are organized into a co-occurrence table (Fig. 3(b)) by the instances in s^t . Additionally, the values of the co-occurring instances associated with the same property are aggregated and represented with a value range.

Next, a pruning graph (Fig. 3(c)) is created to characterize all combinations of the data sources associated with the detected co-occurring instances, such that impossible combinations are pruned at the data source level. The basic idea is that no frequent spatiotemporal patterns of two data sources is present if no sufficient co-occurring instances are found from the two data sources. Each node represents a potential combination of data sources with (1) the IDs of the involved data sources, (2) a list of the co-occurring instances of the involved data sources, and (3) a counter storing the number of the co-occurring instances. For example, the node labeled $\{S^{\alpha}, S^{\beta}, 3\}$ in Fig. 3(c) indicates that three co-occurring instances can be extracted from the data sources S^{α} , S^{β} , and S^{τ} . Links between the nodes depict downward closure relations [1], i.e., an upper-level node contains all the instances of its linked lower-level nodes. Therefore, the insignificant combinations of data sources can be quickly detected and invalidated from the top to the bottom at this stage, by removing the nodes whose number of associated co-occurring instances is lower than the specified threshold ($\lambda_e \cdot |s^t|$, as per the definition of co-occurrence coverage), as illustrated with gray nodes in Fig. 3(c). Corresponding rows in the co-occurrence table are removed thereafter. Hence, the property combinations belonging to invalid data source combinations are eliminated, thereby accelerating the subsequent pattern mining stages.

C. Mining Co-Occurrence Patterns

We propose a two-fold approach to extract frequent patterns from the co-occurrence table and pruning graph. This approach comprises two steps: (1) **Low-Level Pattern Mining** entails identifying low-level patterns that comprise the co-occurring value ranges between the properties of target data source s^t and another data source and (2) **High-Level Pattern Mining** involves identifying high-level patterns that comprise the co-occurring property value ranges across multiple data sources.

1) *Low-Level Pattern Mining*: Low-level pattern mining aims to find significant combinations of the co-occurring value ranges between a property in the target data source s^t and one in another data source by a) aggregating the co-occurring property value ranges discovered at the previous stage with *value matrices* and b) performing a novel sweep-line based algorithm on the value matrices to identify the salient co-occurring property value ranges that satisfy the given threshold λ_a .

a) *Range aggregation*: We first enumerate every possible pair of properties. Then, we aggregate all combinations of co-occurring property value ranges (which we have discovered at the previous step) between two properties in each pair. The aggregation is achieved with a value matrix (Fig. 2(c)). Axes of the value matrix represent the categorical or discretized numerical domains of properties p^x and p^y , while each cell (i, j) in the matrix denotes a property value combination of $v_{p^x} = i$ and $v_{p^y} = j$. We overlay the combinations of co-occurring property value ranges extracted from the co-occurrence table (rectangles in Fig. 2(d)) on the value matrix and maintain the combination and instance counts collected from the covered range combinations for each cell.

b) *Pattern identification*: Given a value matrix, salient patterns appear as the rectangular areas on the matrix in which the instance count of every cell covered by the area is larger than $\lambda_e \cdot |s^t|$, where λ_e indicates the user-desired co-occurrence coverage and $|s^t|$ represents the number of instances of s^t . To detect these areas, we develop a fast algorithm based on the sweep line (Alg. 1). Specifically, we sweep the domain of a property column by column and construct rectangular areas along the way. sw represents the sweep window discovered in the previous column, and col represents the current scanning column. For each column, the algorithm detects vertically continuous *sweep windows* (cf. L3) in which every cell satisfies the constraint and maintains two states, namely, ASW for the active sweep windows detected in the previous column and NSW for the new ones detected in the current column. To replace ASW with NSW, the algorithm considers three cases:

- ◊ **Case 1: Continued.** A sweep window in ASW entirely continues in NSW, thereby remaining active (cf. L5-6).
- ◊ **Case 2: Discontinued.** A sweep window in ASW completely disappears in NSW, thereby being removed from ASW. A new rectangular area will be constructed from the swept area and inserted into the result set RS (cf. L7-9).
- ◊ **Case 3: Partially continued.** A sweep window in ASW only partially continues in NSW. This sweep window will be invalidated and regarded as a discontinued

Algorithm 1 Sweep-Line Pattern Mining Algorithm

Data: Value Matrix VM , desired coverage coverage λ_e .
Result: The result set RS with maximal rectangles in the matrix.

```

1  $ASW \leftarrow \emptyset, NSW \leftarrow \emptyset$ ;
2 for Each column  $col$  in  $VM$  do
3    $NSW \leftarrow$  continuous qualified (satisfying  $\lambda > \lambda_e$ )
   cells in  $col$ ;
4   for sweep window  $sw \in ASW$  do
5     if  $sw$  continues in  $col$  then
6       keep  $sw$  in  $ASW$ ;          /* cont'd */
7     else if sweep window  $sw$  not continue in  $col$  then
8        $RS \leftarrow$  result( $sw$  and  $col$ ); /* discont'd
9       /*
10      remove  $sw$  from  $ASW$ ;
11     else
12        $RS \leftarrow$  result( $sw$  and  $col$ );
13       /* part. cont'd */
14       remove  $sw$  from  $ASW$ ;
15       shrink  $sw$  to the partially overlapped range  $sw'$ ;
16        $ASW \leftarrow sw'$ ;
17 for sweep window  $sw \in NSW$  do
18   if  $sw$  does not have the same window in  $ASW$ 
19   then
20      $ASW \leftarrow sw$ 

```

window (Case 2), and a new shrunk sweep window will be created with the rows that continue from ASW to NSW (cf. L11-14).

In addition, new sweep windows in NSW which are not covered by the above cases will be added to ASW (cf. L15-17). Hence, the low-level patterns between two properties (i.e., the combinations of property value ranges satisfying the given threshold λ_e) are obtained from the result set RS .

2) *High-Level Pattern Mining*: The low-level patterns enable the framework to generate and validate high-level patterns that involve three or more data sources (including the property from the target data source).

a) *Candidate generation*: High-level pattern candidates can be generated by intersecting low-level patterns. For example, pattern candidate $\{s_{p_u}^\tau | (C' \cap C''), \{s_{p_x}^\alpha | C, s_{p_y}^\beta | C\}\}$ can be generated from the intersection of low-level patterns $\{s_{p_u}^\tau | C', \{s_{p_x}^\alpha | C\}\}$ and $\{s_{p_u}^\tau | C'', \{s_{p_y}^\beta | C\}\}$. We only keep candidates whose property value ranges are not empty.

b) *Pattern validation*: A pattern candidate is considered valid only if the number of the instance combinations it covers is larger than $\lambda_e \cdot |s^\tau|$. Valid high-level patterns are inserted into the result set for further interactive analysis.

V. VISUAL DESIGN

This section discusses analysis tasks, design rationales, and the visualization techniques specifically designed for interpreting the extracted patterns.

A. Design Rationales

Although the model can efficiently extract co-occurrence patterns, interpreting massive co-occurrences, detecting unusual anomalies, and obtaining high-level insights remain challenging. Visualization techniques are highly necessary to help expert explore the extracted co-occurrence patterns.

In this study, we have conducted a user-centered design process with three interdisciplinary urban planning experts over the past year. These experts have more than 10 years of experience in developing data-driven solutions for various urban problems, such as location selection, energy planning, and pollution analysis. They approached us to seek an interactive visualization system for interpreting and analyzing the co-occurrence patterns among different heterogeneous data sources, including city-wide meteorological and traffic data, collected in urban environments. Through frequent discussions with the experts, two important analysis tasks, *macro-* and *micro-level analyses*, were identified.

1) *Macro-Level Analysis*: Users select proper mining parameters and wish to see the statistical distribution of all value ranges regarding individual properties. Users also select properties and value ranges that they are interested in for further analysis. A visual summary of co-occurrence patterns should be provided to help users determine a specific property combination and proceed to the micro-level analysis to inspect the co-occurrence patterns in this combination.

2) *Micro-Level Analysis*: A clear overview of the co-occurrence patterns of a given property combination should be provided. Subsequently, users may group interesting patterns for observation and comparison. The spatiotemporal distribution of the instances of the patterns should also be provided for further validation and analysis of the patterns.

Based on these two analysis tasks, the design rationales behind our system are derived and summarized below.

R1 Generating a visual summary of massive patterns

It is challenging to analyze a large number of the patterns individually without a clear visual summary. Such a summary should enable the users to understand how data properties are co-occurring (e.g., finding out whether air pollution is co-occurring with traffic congestion or high traffic volume).

R2 Allowing statistical analysis of properties

Obtaining an intuitive understanding of the overall data range distribution is difficult. Users need a visualization that presents the statistical information of the range distribution of each data property. For instance, the users can determine a pattern does not involve traffic congestion if statistics of the pattern indicates that the value ranges of “Low Speed %” are relatively low.

R3 Enabling interactive visual exploration of patterns

In addition to the summary overview of all patterns, users need a steerable environment to find interesting patterns and conduct further analysis. Hence, the system should enable domain experts to interact with the patterns directly by supporting various interactions like filtering, ranking, and grouping to unfold and inspect the patterns.

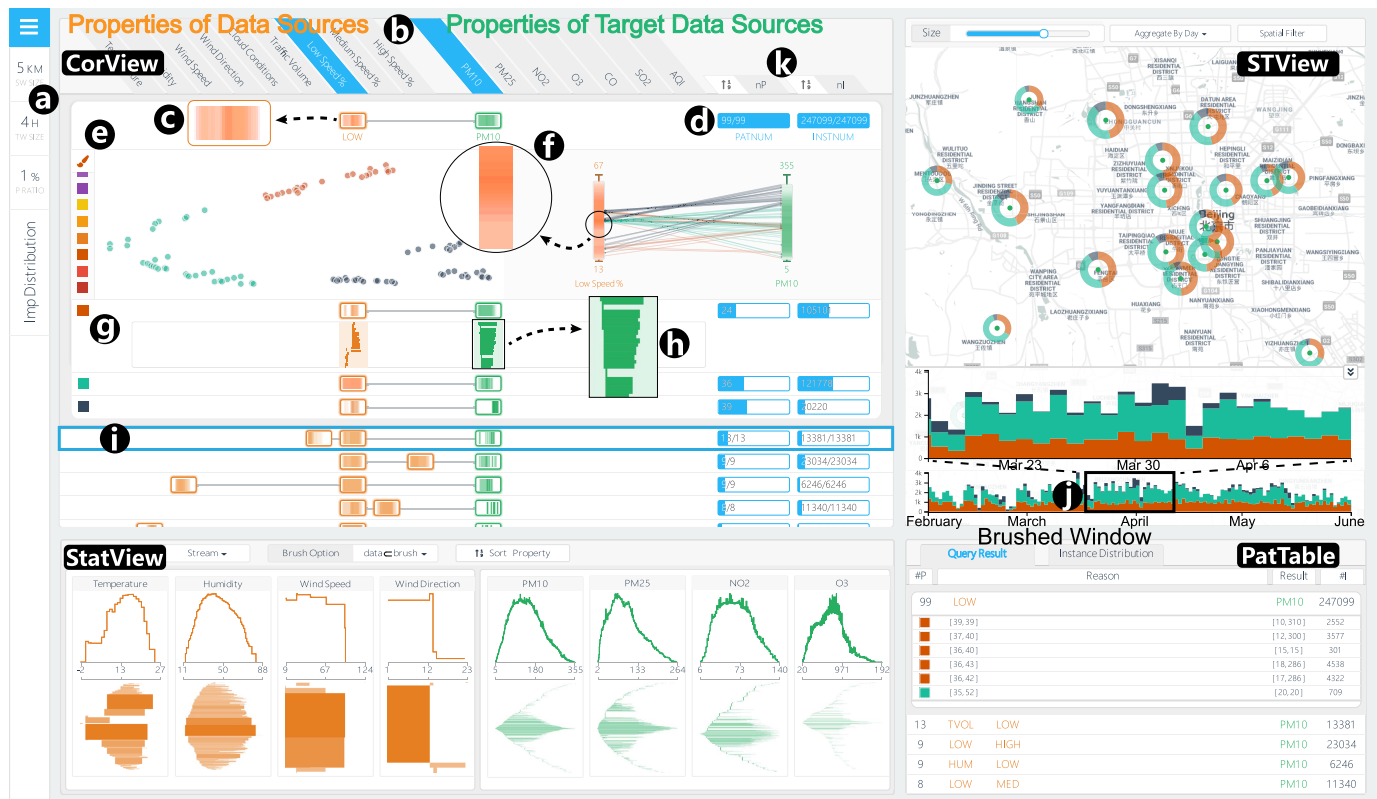


Fig. 4. CorVizor consists of four main views (CorView, STView, StatView, and PatTable) for detecting and understanding co-occurrence patterns.

R4 Visualizing the spatiotemporal information of patterns

A pattern can be associated with many spatiotemporal instances occurring over space and time. The users would like to obtain the spatiotemporal distribution of the pattern's instances in order to answer questions like whether the pattern occurs frequently in the suburbs or whether it happens periodically.

R5 Applying different model parameters

The mining model may not always produce the most desirable results. In order to integrate the domain knowledge into the analysis pipeline and allow users to iteratively improve the result of pattern extraction, user interaction with the model should be supported to select different results of the model.

In the design process, we identified three challenges, namely, diversity, volume, and organization (detailed in Section I). We tackle these challenges by designing CorVizor with four linked views, including CorView, STView, StatView, and PatTable (Fig. 4), based on the aforementioned rationales. CorView is the core component and provides a matrix-style visual summary of the patterns of all property combinations (R1). Multi-level interactive exploration is naturally supported (R3). StatView displays the distributions of value ranges of different data source properties (R2). STView shows the spatiotemporal information of the target instances associated with the patterns (R4). PatTable presents the details of the selected patterns in a table (R3). Choosing different model parameters are also supported (R5) in the Info Panel (Fig. 4(a) and Fig. 10(a)).

B. CorView

This section presents the design of CorView, which visually summarizes the patterns of property combinations (Section V-B1) and interactively unfolds those of a selected property combination (Section V-B2).

1) *Visualization of Property Combinations*: The property combinations can be simply presented with tree visualization (Fig. 3(c)). However, such representation suffers from scalability issues and visual clutter given that a large number of the combinations can be generated. We adopt a scalable matrix-based approach (the **volume** challenge) to visualize the property combinations shared by massive co-occurrence patterns (R1) and provide an unified overview for the diverse data properties among the patterns (the **diversity** challenge). The matrix-based approach is easy to understand and allows users to make efficient visual comparisons of the property combinations in an aligned manner.

Each column in the CorView represents a property, and each row (Fig. 4(i)) represents a group of the co-occurrence patterns with the same property combination (e.g., the patterns that indicate the frequent co-occurrences between “Low Speed %” and “AQI” will be grouped into the same row). In each row, a set of linked rectangular glyphs visualize the summary of the pattern group. Each glyph contains a density map (Fig. 4(c)), which reveals the value range distribution of the property corresponding to the column where the glyph resides. The darker areas indicate that the corresponding value ranges appear more frequently in the patterns. The properties in the target data source and other data sources are shown in green and orange, respectively. In addition, the numbers of

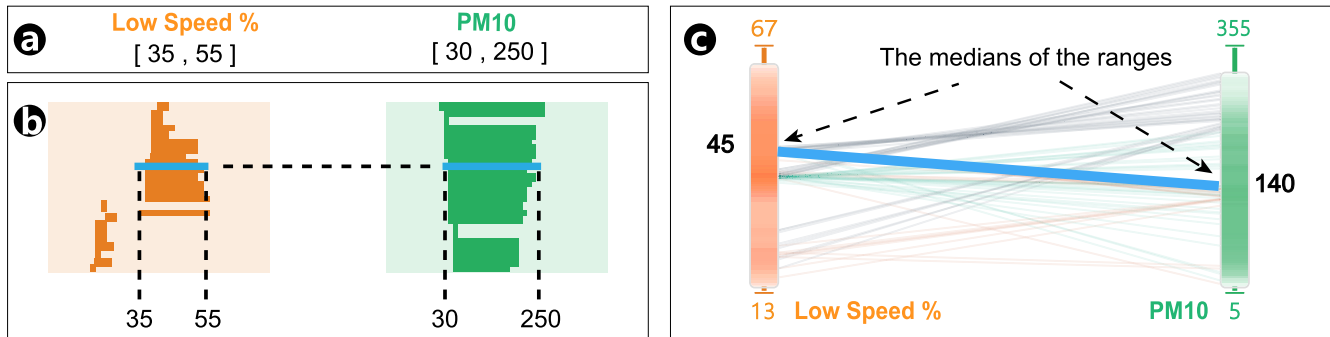


Fig. 5. Example of the visual encodings of the stacked line chart (b) and parallel coordinates (c) for a pattern in (a).

patterns and instances for each group of patterns are visualized with two bars in each row (Fig. 4(d)). By clicking on the column headers (Fig. 4(b, k)), users can filter out the property combinations that do not contain the selected properties or sort the combinations based on the numbers of patterns or instances.

2) *Visualization of Co-Occurrence Patterns*: To maintain scalability of the overview of a property combination's patterns, a collection of highly scalable and well established visualization techniques are chosen to support the coordinated and multifacet analysis of the multidimensional patterns. Users can unfold a property combination in the CorView and analyze the patterns with the selected combination using a similarity-based scatterplot (**RI**), a tailored parallel coordinates plot, and a stacked line chart (Fig. 4(h)) in the expanded view (Fig. 4(e)). The scatterplot provides an overall picture of the similarity among co-occurrence patterns, thereby enabling users to group the patterns and detect anomalies. The parallel coordinates plot depicts the co-occurrence among multiple properties, where the value range distributions of the patterns w.r.t. each property are encoded with a density map on the axis. The stacked line chart further provides a compact visualization of the value ranges of the selected patterns. These scalable visualizations are combined to organize patterns without severe visual clutter and facilitate the effective exploration of both the overview and details (the **volume** and **organization** challenges).

a) *Scatterplot*: Analyzing the relationship between patterns is important for pattern grouping, comparison, and anomaly detection. A scatterplot is used to display the co-occurrence patterns, such that similar patterns are naturally grouped together. This scatterplot provides a concise overview of pattern relationships with less clutter than parallel coordinates. The multidimensional scaling (MDS) is used to create the scatterplot. The distance of patterns i and j is computed with $\sqrt{\sum_k^n d(i_k, j_k)^2}$, where n is the number of properties in the pattern, k denotes property k , and $d(i_k, j_k)$ is the distance between the two ranges with regard to property k of pattern i and j . The Jaccard index and the KL divergence were tested to calculate the distance. However, distance is regarded as a constant value by both measures when two ranges are disjoint regardless of how far the ranges semantically appear. Thus, a new measure is used. In this new measure, four features are extracted from each range: lower bound (lb), upper bound (ub), median (mid), and length (len). All features are nor-

malized into $[0, 1]$ by dividing with the value ranges defined in Table 1. The range distance of property k is measured with the Euclidean distance of the pair of ranges, namely, $d(i_k, j_k) = \sqrt{\Delta lb^2 + \Delta ub^2 + \Delta mid^2 + \Delta len^2}$, where Δ represents the difference between two feature values. For those categorical properties that can be ordered, we assign numeric values for each category starting from 1 by the categorical order and compute the range distance based on these values. For those categorical properties that cannot be ordered, we map the text descriptions of those categories to high-dimensional space with word2vec [32], [33]. The word2vec model can generate a high-dimensional vector for each word considering their semantics in a series of sentences. The Euclidean distance between two vectors indicates the semantic similarity between two corresponding words. As such, the distance between the two categories can be measured. Users can group patterns and highlight anomalies by brushing the corresponding points with various colors.

b) *Parallel coordinates*: Co-occurrence patterns can have a high-level form (Section IV-C2) with more than three properties involved. Thus, parallel coordinates are used as a uniform view to display the multidimensional co-occurrence patterns. Each axis represents a property. The medians of the ranges are used as the end points of the line segments to connect the value ranges in various property axes. Considering that overlaps exist among ranges of the same property, we do not adopt parallel sets as it is more suitable for categorical and disjoint data.

The range distribution is displayed with a density map on its corresponding axis (Fig. 4(f)). A density map is used instead of other methods, such as histogram, because it consumes less space and compactly shows the density information of the property value. In each density map, value ranges are drawn along each coordinate with equal opacity. The ranges are overlaid and their opacity values are combined to encode the density (i.e., dark areas indicate that the corresponding values or categories are covered by many ranges). Fig. 5 (a) shows an example of a pattern and Fig. 5 (c) shows the corresponding visual representation with tailored parallel coordinates.

c) *Stacked line chart*: When a set of patterns are grouped in the scatterplot, a new row summarizing the pattern group is automatically generated and attached under the scatterplot and parallel coordinates. Fig. 4(g) shows one of the three rows

of the grouped patterns. When a user unfolds a row, the row is expanded to show a stacked line chart (Fig. 4(h)). In the chart, the value range for each property is represented by a fine line segment. The segments are stacked to compose a distribution map. Fig. 5(b) shows an example of the chart, where the left and right endpoints of the line segments denote the two endpoints of the range.

3) *Design Alternatives*: In the aforementioned user-centric design process, we attempted to refine the visual design of CorVizor iteratively by proposing and evaluating alternatives. In this section, two design alternatives are discussed to reveal the rationales behind our design choices in terms of the macro- and micro-level analyses of co-occurrence patterns.

a) *Visualization of patterns in many property combinations*: Instead of organizing property combinations with a matrix-based CorView, we attempted to maintain the structure of these combinations with a node-link diagram (Fig. 6(a)). Each node in the diagram represents a property combination. The directed edges in the diagram indicate the composition of subsequent combinations. In each node lies a glyph, which encodes the distribution of property value ranges, and the size of the glyph shows the number of pattern instances. Moreover, we allow users to apply filters to keep the desired combinations by selecting properties on the top. Although such an alternative clearly reveals the inherent structure of property combinations, three major weaknesses prohibit it from being applied in our system, that is, a) the proposed node-link diagram costs excessive screen space; b) the crossing edges introduce serious visual clutters and are thus not scalable; and c) the distributions of property value ranges in different nodes are difficult to compare because they are not aligned.

b) *Visualization of patterns of a property combination*: To help analysts grasp the similarity among massive co-occurrence patterns, we initially projected these patterns into a 2D view via dimensional reduction techniques. Inspired by Liu *et al.* [28], we attempted to depict these patterns with glyphs embedded directly into the view (Fig. 6(b)). On the edge of the glyph lies a circular histogram that encodes the temporal pattern distribution of a single pattern, and the properties involved in the pattern are represented by homocentric donut charts. However, such an approach is not scalable with the number of patterns. The value ranges encoded with radians can also be misleading. Hence, we iterated our design by dissecting the high-dimensional information in these patterns with multiple coordinated views as described previously.

C. StatView

Although CorView shows the property combinations and their co-occurrence patterns, the value range distribution aggregated by properties remains unavailable. This information is essential for high-level exploration. Thus, StatView is used to display the pattern distribution of each property (Fig. 4) (R2). This pattern distribution comprises small multiples that display the distributions of the value ranges for the properties.

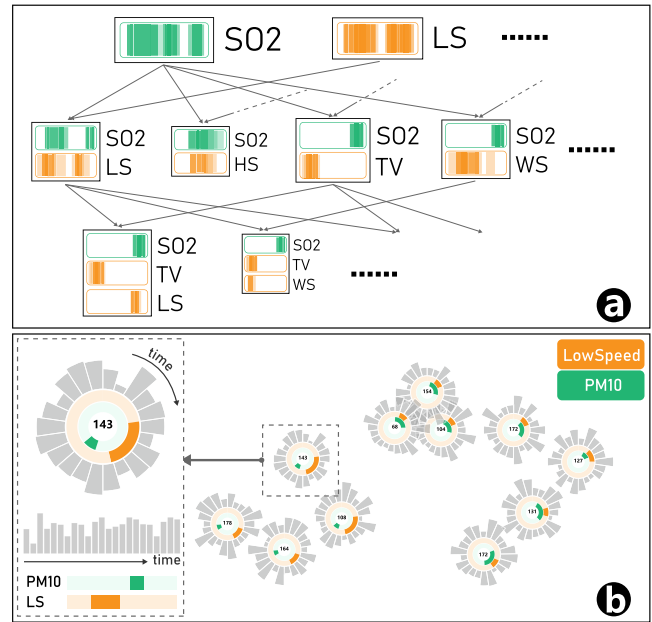


Fig. 6. Design Alternatives for CorView. (a) Node-link visualization for many property combinations; (b) Glyph-based visualization for co-occurrence patterns.

An individual plot of a property displays a value distribution in a line chart (upper part) and presents a range distribution in a stacked rectangle chart (lower part). Both parts have their counterparts in CorView. The upper and lower parts correspond to the density map and stacked line chart, respectively.

StatView uses a line chart to encode the density distribution of the values or categories of each property. Position information is considered perceptually effective in encoding magnitude [36]. Moreover, StatView has more space to display the distribution information. Thus, we utilize position rather than luminance to encode the distribution information. The same value or category ranges are aggregated, and a rectangle is used to represent each unique range in the stacked rectangle chart. The height of a rectangle represents the number of occurrences of the associated range.

Users can brush a span of the property value on the horizontal axis of any line chart to perform filtering. Other views can be subsequently updated. StatView supports three types of filtering interactions: a) removing the patterns that overlap with the selected value ranges; b) showing only the patterns that overlap with the selected value ranges; and c) showing only the patterns that are strictly inside the selected value ranges.

D. STView

STView allows users to gain insights into the spatiotemporal trends of the patterns (R4). The bottom of the view shows a histogram to visualize the temporal distributions of pattern instances with multiple scales. Users can easily select patterns by brushing a temporal window (Fig 4(j)). When the patterns are grouped in CorView, the histogram shows the temporal distributions of pattern instances of different groups by using stacked bars. The top of the view shows the spatial distribution

of the pattern instances on a map. In this study, we use air quality data as the target data source. Each air quality station is represented by a donut chart whose radius encodes the total number of pattern instances in this location. The sectors in different colors indicate the ratios of the pattern instances of each group selected from CorView in each location. Donut charts are used instead of pie charts because the former has a blank center. Users can see through it to observe the details.

When a user hovers his or her mouse on a glyph, a circle around the glyph is displayed to show the coverage of the associated station, namely, the size of the spatial window used in the mining model. The circle covers the spatial area whose co-occurring instances can be viewed as being co-located with the air quality station. Users can select several stations to see the related co-occurrence patterns in other views.

E. PatTable

PatTable is a table-like component that allows users to inspect raw patterns directly on demand. Each row represents a property combination. The combinations can be unfolded to show the corresponding patterns. Detailed information for each individual pattern, such as the number of instances and the co-occurring property value ranges, is depicted in the unfolded view. Moreover, PatTable is coordinated with CorView, where user interactions in one view are reflected in the other view.

F. User Interactions

CorVizor supports various basic and advanced interactions.

- ◊ **Showing overview first and details on demand.** CorView shows a succinct overview of all property combinations. Users can click on a row to explore the corresponding property combination in detail.
- ◊ **Brushing and filtering.** Users are allowed to group patterns or spot anomalies by brushing the patterns with colors in a scatterplot in CorView. Users can filter by spatial area, time range, and property value range in STView and StatView.
- ◊ **Changing the model parameters.** Users can change the parameters including the spatiotemporal threshold and minimum co-occurrence coverage (Fig. 10(a)) and see new results (R5). The histogram (Fig. 10(b)) shows the distribution of the normalized range widths of all co-occurrence patterns.

VI. EXPERIMENTS

This section presents model evaluation, case studies, and expert interview to evaluate the effectiveness and usability of the proposed system. The experimental data contain the 16 data properties from the three data sources listed in Table I in Section III-A. The data were collected from a large city. Data collection was conducted from February 1 to May 31 in 2014. **Weather data** were collected hourly from 20 weather monitoring stations around the city. **Air quality data** were collected hourly from 36 air quality monitoring stations in the city, and **traffic data** were collected from 100,215 segments of the city road network every half hour from a geospatial

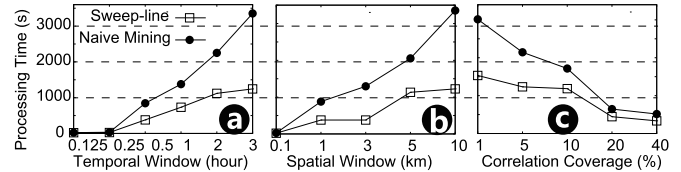


Fig. 7. Comparison of the time performances of the proposed sweep-line algorithm and naive mining method.

mapping platform. To sum up, there are 103 thousand records in the air quality data, 57 thousand records in the weather data, and 577,238 thousand records in the traffic data. All experiments were evaluated on a laptop running Windows 10 with Intel Core i7 3.4GHz CPU, 256GB SSD drive, and 16 GB RAM.

A. Model Evaluation

The proposed sweep-line algorithm is the core component of our pattern mining model. We compared it with a naive approach to demonstrate its performance over the baseline.

1) *Naive Approach*: The naive method to identify distinctive rectangles from a value matrix follows the following steps. First, every cell in the matrix is scanned. Second, if a qualified cell (fulfilling the co-occurrence coverage requirement) is identified, the naive approach considers the cell the left-up corner of certain distinctive rectangles and traverses toward right and down directions to find the rectangles as candidates. Third, each candidate rectangle is tested to see if it is completely covered by other rectangles identified previously. If overlapping cases exist, the candidate rectangle is discarded. Otherwise, the identified rectangle is inserted into the result set. The cost of the approach is prohibitively high. Assuming that an $M \times N$ matrix exists, the approach needs to traverse the entire matrix to identify qualified cells in the outer iteration. For each qualified cell, the approach needs to traverse the remainder of the matrix to identify the distinctive rectangles and test their qualification, which may result in $\mathcal{O}(M^2N^2)$ in the worst case. In contrast, our method has the time complexity of $\mathcal{O}(MN)$ as we only need to scan the matrix once.

The comparison was performed with the varying co-occurrence coverage, spatial window, and temporal window. Fig. 7 shows the results of the time performance comparison.

a) *Temporal window*: Fig. 7(a) shows that the time of both approaches increases with the increase in the temporal window. A large temporal window usually generates a large value range, which increases the probability of finding qualified cells in the value matrix. The sweep-line algorithm outperforms the naive approach, given that a number of redundant cell examinations are avoided during the process.

b) *Spatial window*: Fig. 7(b) provides two observations. First, with a large spatial window, the processing time of both approaches increase because many co-related instances are to be analyzed in a spatiotemporal partition. Second, the processing time of our approach is lower than that of the naive approach, since the naive approach needs to examine more qualified cells and thus incurs more computational time.

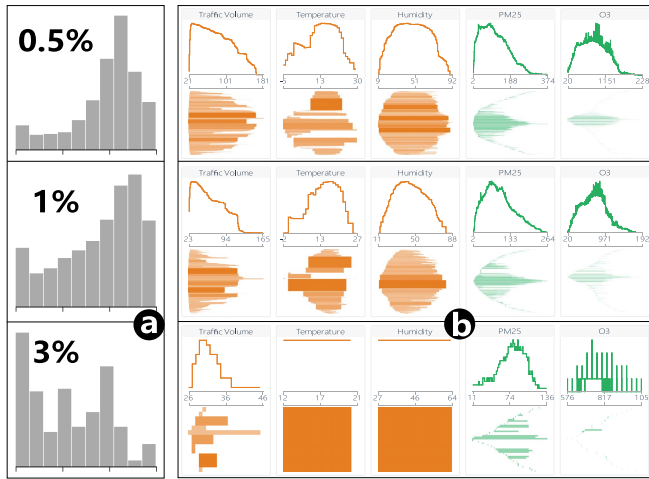


Fig. 8. Selecting proper parameters: (a) histograms showing the distribution of the normalized range widths and (b) statistical information for the co-occurrence coverage thresholds 0.5%, 1% and 3% in StatView.

c) *Co-occurrence coverage*: Fig. 7(c) presents two observations. First, the processing time of both approaches decreases. Second, the sweep-line algorithm performs better than the naive one because a large co-occurrence coverage value results in a small chance of finding a qualified cell (i.e., qualified patterns) in the value matrix.

B. Case Studies

The case studies were conducted with the domain experts to evaluate the effectiveness of our system.

1) *Macro Analysis (Co-Occurrences Relevant to High SO₂)*: This case study demonstrated the effectiveness of CorVizor for the macro-level analysis (detailed in Section V-A).

Selecting proper mining parameters was the first step to explore the co-occurrence patterns. The domain experts suggested 5 km and 4 hours for spatial and temporal window based on their experience. However, the co-occurrence coverage threshold is difficult to choose. The experts attempted the thresholds 0.5%, 1%, and 3%. The distributions of the normalized range widths generated by these thresholds are presented in Fig. 8(a), and the associated statistical information on the extracted patterns is depicted in StatView (Fig. 8(b)). Based on their observations, the experts selected 1% as the threshold because: a) although the histograms generated by the thresholds 0.5% and 1% seemed similar, the patterns represented as stacked rectangles in StatView with the threshold 1% were more organized and meaningful than those with 0.5%; and b) the patterns extracted with the threshold 3% was too coarse to reveal any useful insights. Thus, the threshold 1% (Fig. 10(a)) was selected by the experts for further explorations.

Urban air pollution, which is crucially related to the well-being of city residents, has attracted increasing concerns in recent years. Therefore, the experts attempted to identify the co-occurrence between air quality and other urban data sources with CorVizor. In particular, they were interested in the co-occurrence patterns relevant to high SO₂ because SO₂ was one of the major pollutants produced by human activities in

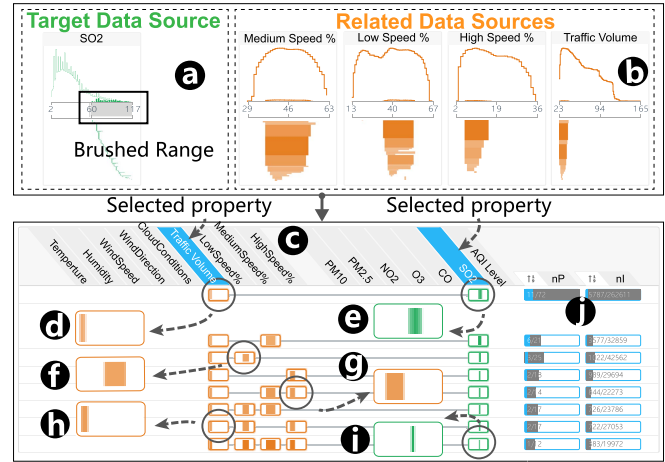


Fig. 9. Macro-level analysis of the co-occurrence patterns that are relevant to high SO₂.

cities. Hence, the experts selected the patterns that comprised high SO₂ in StatView. Fig. 9(b) showed that low traffic volume strongly co-occurred with high SO₂ because the bars in the range distribution view of traffic volume indicated that the corresponding value ranges were relatively low and narrow.

These findings seemed contradicted with the experts' intuition, as they believed that only huge traffic volume would result in severe air pollutant emission. Hence, they selected traffic volume and SO₂ in CorView for further exploration (Fig. 9(c)). Only the property combinations that involved these two selected properties remained in the view. The glyphs in the first row (Fig. 9(d) and 9(e)) validated the aforementioned observation with StatView. By analyzing other rows, the experts discovered that the number of low-speed vehicles (Fig. 9(f)) was considerably larger than that of high-speed vehicles (Fig. 9(g)) while the traffic volume (Fig. 9(h)) was low and SO₂ (Fig. 9(i)) was high. Thus, the experts suggested that the large number of slow vehicles and small traffic volume could be a sign of potential traffic congestions, which resulted in the high SO₂ emission. The co-occurrence patterns among traffic volume, low-speed vehicles, and AQI level were also explored with the identical approach. The result was similar: the air quality appeared to be bad with small traffic volume and the large number of slow vehicles. This insight confirmed that the small traffic volume co-occurred with severe traffic congestions, which were a significant contributing factor, confirmed by the experts, to the deteriorated urban air quality.

2) *Micro Analysis (Co-Occurrences Involving Air Pollution)*: The second case study demonstrates the usefulness of the system in analyzing the co-occurrence patterns associated with a specific property combination.

Road space rationing policies were widely adopted by governments to alleviate serious air pollution. However, the experts doubted the effectiveness of these policies. Hence, they would like to analyze the co-occurrences between traffic and air pollution with our system. The combination of the property "Low Speed (%)" and "PM10" in CorView were selected in this study. The corresponding row was expanded to show its details for an in-depth exploration. Patterns that were

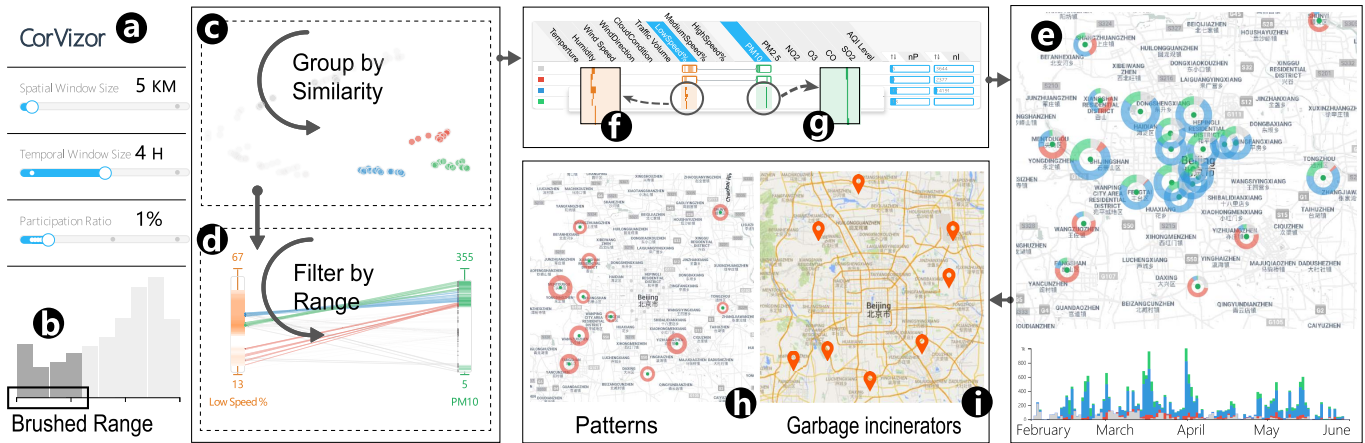


Fig. 10. Micro-level analysis of co-occurrence patterns between traffic congestion and air pollution.

extremely general were filtered out by brushing the histogram of the normalized range width (Fig. 10(b)). The patterns in the scatterplot under the expanded row were grouped in different colors based on the closeness of the patterns in the plot (Fig. 10(c)). The parallel coordinates were colored accordingly (Fig. 10(d)). The red group is considerably different from the blue and green groups in parallel coordinates (Fig. 10(d)). The red group represents “less low speed (%) and high PM10,” whereas the blue and green groups indicate “more low speed (%) and high PM10.”

The experts were particularly interested in the red pattern group. The stacked line chart of the group in Figs. 10(f) and 10(g) also indicates “less low speed (%) (i.e., traffic congestion is unlikely to occur) and high PM10.” STView was used to examine the spatiotemporal distribution of the patterns of the selected groups (Fig. 10(e)). To experts’ surprise, the red group only occurred in the area between Rings 5 and 6 of the expressways, which is the suburban area of the city. One expert indicated that there were several garbage incineration plants in this area. Comparison between the distribution of the red group patterns (Fig. 10(h)) and that of the garbage incinerators (Fig. 10(i)) showed a clear match. They speculated that the number of garbage incinerators could frequently co-occur with the high PM10 in the area, in which traffic congestion did not occur. Further investigation and analysis in the field were required to verify this conjecture and determine its plausible cause.

Furthermore, CorVizor was used to explore the co-occurrence patterns involving air quality index (AQI), which would increase as air quality worsens. The experts were curious about the reasons behind the worst air quality represented by the highest AQI level with the value of 5. Thus, they drew a selection on the AQI property in StatView (Fig. 11(a)) to select the patterns that involved the highest AQI level. From STView, the experts observed that most of these patterns occurred around February and March (Fig. 11(b)). They suggested that air pollution might be caused by coal heating in the winter, which emitted massive pollutants and severely deteriorated the air quality. To confirm this hypothesis, the experts selected the temperature and AQI

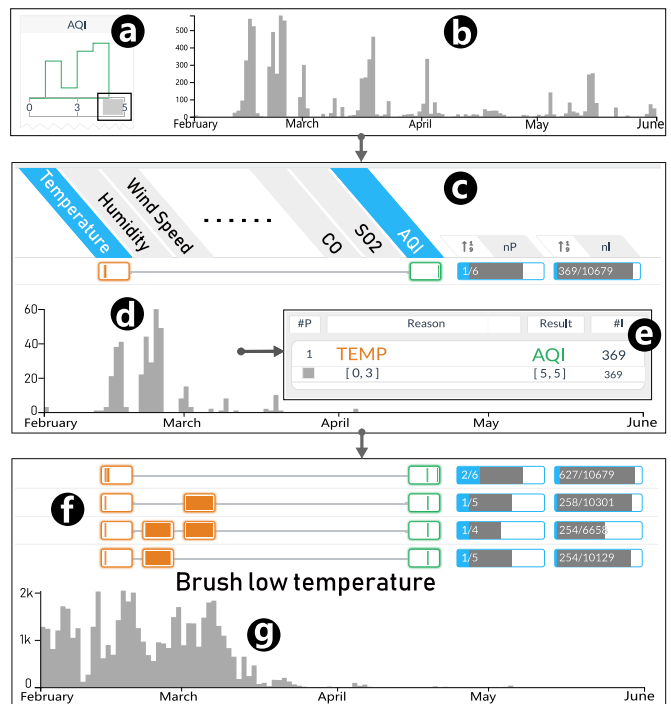


Fig. 11. Micro-level analysis of co-occurrence patterns between temperature and air quality index (AQI).

properties in CorView (Fig. 11(c)) and discovered that the highest AQI level co-occurred with low temperature. Moreover, the temporal distribution of these selected co-occurrence patterns was identical to that of the patterns involving the highest AQI level (Fig. 11(d)). The temperature ranges in PatTable were around 0 °C (Fig. 11(e)), which also provided useful hints for this co-occurrence. Furthermore, the experts attempted to verify the co-occurrence by selecting the patterns with low temperature in StatView. They were satisfied to discover that these patterns all co-occurred with medium and high AQI levels (Fig. 11(f)). These observations supported the experts’ hypothesis and helped them link the deteriorating air quality with coal heating.

In this case study, the co-occurrence patterns involving both numerical and categorical properties were explored and

analyzed in detail. These detailed exploration and analysis demonstrate the effectiveness of CorVizor in handling the micro-level analysis tasks and providing interesting insights into the co-occurrence patterns for further verification and analysis.

C. Interview With Domain Experts

After the case studies, we collected and summarized the feedback from the experts as follows.

1) *Overall System Usability*: CorVizor was well received by the experts. They were pleased to explore and analyze the massive heterogeneous patterns intuitively with the proposed interactive visualizations. “*The visualization system makes the co-occurrence patterns produced by the data mining model much more meaningful,*” an expert said. Both experts acknowledged that the analytical workflow of our system could help them gain considerable insights into the spatiotemporal co-occurrences. Moreover, they believed that our system could be extended to identify interesting co-occurrences in various scenarios, such as business location selection and travel recommendation.

2) *Visual Design and Interactions*: Both the experts were impressed by the visual design and interactions. They praised CorView, which presents the co-occurrences among various data properties explicitly. An expert commented “*the matrix-like layout is familiar to me and the hierarchical visualization method well organizes the exploration process.*” He was also highly satisfied with the intuitive visual summary of the co-occurrence patterns provided in CorVizor. Another expert was deeply impressed by the spatiotemporal view. “*Without this system, it would be impossible to discover interesting cases related to the spatiotemporal distribution of the co-occurrence patterns,*” he said. Both experts appreciated the interactive features of the proposed system. They especially appreciated the usefulness of filtering and brushing. The experts said that these techniques help in anomaly detection and pattern grouping and comparison.

3) *Suggestion*: The usability of our system was confirmed by the experts, who immediately became familiar with the system after a brief training. Nevertheless, they suggested that the design of our system could be simplified further, such as by replacing the scatterplot and parallel coordinate plot with a plain list with numbers and figures, and integrates visual guides to allow average users, such as government officials, to monitor the city dynamics and grasp interesting insights conveniently. We will leave this simplified version of our system as a part of our future work.

VII. DISCUSSION

In this section, we discuss the implications, limitations, and generalization of the proposed system.

A. Implications

CorVizor can identify interesting co-occurrence patterns that may facilitate numerous *transportation applications*, such as traffic management and transportation planning. Important insights revealed by these patterns, including how the

traffic speed and volume in a local area affect the concentrations of air pollutants, provide strong decision-making contexts for urban planners to establish informed road policies and long-term planning strategies in advance. Nevertheless, co-occurrences do not necessarily imply causation. Analysts may not be able to come up with a clear actionable plan with only co-occurrences, and inferring causal relationships remains a challenge. However, the present work still has several important implications with regard to causal inference. First, pattern co-occurrences can reduce the search space of causal inference. Second, the special characteristics of pattern co-occurrences can have significant implications for research on causal inference. Moreover, with CorVizor, data mining researchers can easily obtain an intuitive overview of a large number of co-occurrence patterns while checking the credibility of any specific co-occurrence pattern or group of co-occurrence patterns. As such, researchers can be informed of imperfections of the data mining model, consequently inspiring them to enhance the model’s effectiveness.

B. Limitations

The time performance of the co-occurrence mining framework is not highly optimized. Running the model for our experimental dataset usually requires nearly an hour. Data mining results for possible parameter combinations were computed in advance to support the interactive adjustment of the model results. We plan to optimize the model and adapt it to a high-performance distributed computing platform, such that the interactive adjustment of the model setting is made possible. As for the design part of our system, MDS adopted by the scatterplot is widely used in the visualization literature, but it may be misleading at times [42]. To enhance the scatterplot, the method for visualizing dimensionally-reduced data [42] can be further incorporated into our system.

C. Generalization

CorVizor can be directly applied to various urban analysis applications, such as urban planning, pollution diagnoses, and location selection, to detect and understand the co-occurrence patterns in spatiotemporal datasets that support effective decision-making processes. The case studies we presented were conducted for pollution diagnosis. However, the target data source can be changed to identify other interesting co-occurrences in other domains. For example, traffic congestion [13], [47], [49] can be analyzed efficiently by setting traffic data as the target. In addition, the evolution of business is closely related to many latent co-occurrence patterns extracted from various urban datasets [44], which can also be captured by our framework.

VIII. CONCLUSION AND FUTURE WORK

In this work, we studied the extraction and interpretation of fine-grained spatiotemporal co-occurrence patterns that comprise various properties of different types, scales, and semantics. Based on the proposed data mining framework and interactive multi-scale visualization technique, we developed

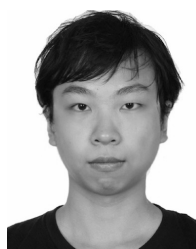
CorVizor, a visual analytics system that assists users in exploring these patterns. This study contributes an important step towards the in-depth understanding of urban dynamics formed by the complex co-occurrence patterns extracted from heterogeneous spatiotemporal data sources, including transportation data.

We will continue on improving our system in several ways as follows. First, we plan to develop the deep learning algorithm [61] for mining co-occurrence patterns and migrate the mining module to a high-performance distributed computing platform. Users can directly interact with the model and see the results instantly in CorVizor. Second, we will deploy CorVizor in the field, such that the streaming datasets collected from diverse sources can be fed into the system in real-time, thereby enabling a proactive analysis workflow of the urban problems.

REFERENCES

- [1] R. Agrawal, T. Imieliński, and A. Swami, "Mining association rules between sets of items in large databases," *ACM SIGMOD Rec.*, vol. 22, no. 2, pp. 207–216, Jun. 1993.
- [2] R. Agrawal and R. Srikant, "Fast algorithms for mining association rules in large databases," in *Proc. VLDB*, 1994, pp. 487–499.
- [3] R. Agrawal and R. Srikant, "PrefixSpan: Mining sequential patterns efficiently by prefix-projected pattern growth," in *Proc. ICDE*, 2001, pp. 215–224.
- [4] G. Andrienko, N. Andrienko, P. Bak, D. Keim, and S. Wrobel, *Visual Analytics of Movement*. Berlin, Germany: Springer-Verlag, 2013.
- [5] G. Andrienko, N. Andrienko, C. Hurter, S. Rinzivillo, and S. Wrobel, "Scalable analysis of movement data for extracting and exploring significant places," *IEEE Trans. Vis. Comput. Graphics*, vol. 19, no. 7, pp. 1078–1094, Jul. 2013.
- [6] N. Andrienko, G. Andrienko, L. Barrett, M. Dostie, and P. Henzi, "Space transformation for understanding group movement," *IEEE Trans. Vis. Comput. Graphics*, vol. 19, no. 12, pp. 2169–2178, Dec. 2013.
- [7] B. Aydin, A. Kucuk, R. A. Angryk, and P. C. Martens, "Measuring the significance of spatiotemporal co-occurrences," *ACM Trans. Spatial Algorithms Syst.*, vol. 3, no. 3, pp. 1–35, Nov. 2017.
- [8] G. Bothorel, M. Serrurier, and C. Hurter, "From visualization to association rules: An automatic approach," in *Proc. Spring Conf. Comput. Graph. (SCCG)*, 2013, pp. 57–64.
- [9] N. Cao, C. Lin, Q. Zhu, Y.-R. Lin, X. Teng, and X. Wen, "Voila: Visual anomaly detection and monitoring with streaming spatiotemporal data," *IEEE Trans. Vis. Comput. Graphics*, vol. 24, no. 1, pp. 23–33, Jan. 2018.
- [10] M. Celik, S. Shekhar, J. P. Rogers, and J. A. Shine, "Mixed-drove spatiotemporal co-occurrence pattern mining," *IEEE Trans. Knowl. Data Eng.*, vol. 20, no. 10, pp. 1322–1335, Oct. 2008.
- [11] V. P. Chakka, A. C. Everspaugh, and J. M. Patel, "Indexing large trajectory data sets with SETI," in *Proc. CIDR*, 2003, pp. 1–12.
- [12] S. Chen *et al.*, "Interactive visual discovering of movement patterns from sparsely sampled geo-tagged social media data," *IEEE Trans. Vis. Comput. Graphics*, vol. 22, no. 1, pp. 270–279, Jan. 2016.
- [13] W. Chen, Z. Huang, F. Wu, M. Zhu, H. Guan, and R. Maciejewski, "VAUD: A visual analysis approach for exploring spatio-temporal urban data," *IEEE Trans. Vis. Comput. Graphics*, vol. 24, no. 9, pp. 2636–2648, Sep. 2018.
- [14] F. Chirigati, H. Doraiswamy, T. Damoulas, and J. Freire, "Data polygamy: The many-many relationships among urban spatio-temporal data sets," in *Proc. ACM SIGMOD*, 2016, pp. 1011–1025.
- [15] P. Cudre-Mauroux, E. Wu, and S. Madden, "TrajStore: An adaptive storage system for very large trajectory data sets," in *Proc. IEEE 26th Int. Conf. Data Eng. (ICDE)*, Mar. 2010, pp. 109–120.
- [16] Z. Deng *et al.*, "AirVis: Visual analytics of air pollution propagation," *IEEE Trans. Vis. Comput. Graphics*, vol. 26, no. 1, pp. 800–810, Aug. 2020.
- [17] H. Doraiswamy, N. Ferreira, T. Damoulas, J. Freire, and C. T. Silva, "Using topological analysis to support event-guided exploration in urban data," *IEEE Trans. Vis. Comput. Graphics*, vol. 20, no. 12, pp. 2634–2643, Dec. 2014.
- [18] G. Ertek and A. Demiriz, "A framework for visualizing association mining results," in *Proc. ICCIS*, 2006, pp. 593–602.
- [19] N. Ferreira *et al.*, "Urbane: A 3D framework to support data driven decision making in urban development," in *Proc. IEEE Conf. Vis. Anal. Sci. Technol. (VAST)*, Oct. 2015, pp. 97–104.
- [20] N. Ferreira, J. Poco, H. T. Vo, J. Freire, and C. T. Silva, "Visual exploration of big spatio-temporal urban data: A study of new york city taxi trips," *IEEE Trans. Vis. Comput. Graphics*, vol. 19, no. 12, pp. 2149–2158, Dec. 2013.
- [21] D. Guo and X. Zhu, "Origin-destination flow data smoothing and mapping," *IEEE Trans. Vis. Comput. Graphics*, vol. 20, no. 12, pp. 2043–2052, Dec. 2014.
- [22] M. Hahsler and S. Chelluboina, "Visualizing association rules in hierarchical groups," in *Proc. 42nd Symp. Interface, Stat., Mach. Learn., Vis. Algorithms*, 2011, pp. 1–11.
- [23] X. Huang, Y. Zhao, C. Ma, J. Yang, X. Ye, and C. Zhang, "TrajGraph: A graph-based visual analytics approach to studying urban network centralities using taxi trajectory data," *IEEE Trans. Vis. Comput. Graphics*, vol. 22, no. 1, pp. 160–169, Jan. 2016.
- [24] Y. Ke, J. Cheng, and W. Ng, "Mining quantitative correlated patterns using an information-theoretic approach," in *Proc. 12th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining (KDD)*, 2006, pp. 227–236.
- [25] R. Kosara, F. Bendix, and H. Hauser, "Parallel sets: Interactive exploration and visual analysis of categorical data," *IEEE Trans. Vis. Comput. Graphics*, vol. 12, no. 4, pp. 558–568, Jul. 2006.
- [26] J. Li, S. Chen, K. Zhang, G. Andrienko, and N. Andrienko, "COPE: Interactive exploration of co-occurrence patterns in spatial time series," *IEEE Trans. Vis. Comput. Graphics*, vol. 25, no. 8, pp. 2554–2567, Aug. 2019.
- [27] X. Lin, A. Mukherji, E. A. Rundensteiner, and M. O. Ward, "SPIRE: Supporting parameter-driven interactive rule mining and exploration," *Proc. VLDB Endowment*, vol. 7, no. 13, pp. 1653–1656, Aug. 2014.
- [28] D. Liu *et al.*, "SmartAdP: Visual analytics of large-scale taxi trajectories for selecting billboard locations," *IEEE Trans. Vis. Comput. Graphics*, vol. 23, no. 1, pp. 1–10, Jan. 2017.
- [29] G. Liu, A. Suchitra, H. Zhang, M. Feng, S.-K. Ng, and L. Wong, "AssocExplorer: An association rule visualization system for exploratory data analysis," in *Proc. 18th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining (KDD)*, 2012, pp. 1536–1539.
- [30] Z. Liu, Y. Huang, and J. R. Trampier, "Spatiotemporal topic association detection on tweets," in *Proc. 24th ACM SIGSPATIAL Int. Conf. Adv. Geograph. Inf. Syst. (GIS)*, 2016, p. 28.
- [31] R. Maciejewski *et al.*, "Forecasting hotspots—A predictive analytics approach," *IEEE Trans. Vis. Comput. Graphics*, vol. 17, no. 4, pp. 440–453, Apr. 2011.
- [32] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," 2013, *arXiv:1301.3781*. [Online]. Available: <http://arxiv.org/abs/1301.3781>
- [33] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Proc. NIPS*, 2013, pp. 3111–3119.
- [34] F. Miranda *et al.*, "Urban pulse: Capturing the rhythm of cities," *IEEE Trans. Vis. Comput. Graphics*, vol. 23, no. 1, pp. 791–800, Jan. 2017.
- [35] P. Mohan, S. Shekhar, J. A. Shine, and J. P. Rogers, "Cascading spatio-temporal pattern discovery," *IEEE Trans. Knowl. Data Eng.*, vol. 24, no. 11, pp. 1977–1992, Nov. 2012.
- [36] T. Munzner, "Marks and channels," in *Visualization Analysis and Design*. Boca Raton, FL, USA: CRC Press, 2014, ch. 5, p. 102.
- [37] K. G. Pillai, R. A. Angryk, and B. Aydin, "A filter-and-refine approach to mine spatiotemporal co-occurrences," in *Proc. 21st ACM SIGSPATIAL Int. Conf. Adv. Geograph. Inf. Syst.*, 2013, pp. 104–113.
- [38] H. Qu, W.-Y. Chan, A. Xu, K.-L. Chung, K.-H. Lau, and P. Guo, "Visual analysis of the air pollution problem in hong kong," *IEEE Trans. Vis. Comput. Graphics*, vol. 13, no. 6, pp. 1408–1415, Dec. 2007.
- [39] R. Rastogi and K. Shim, "Mining optimized association rules with categorical and numeric attributes," *IEEE Trans. Knowl. Data Eng.*, vol. 14, no. 1, pp. 29–50, Jan./Feb. 2002.
- [40] R. Scheepens, C. Hurter, H. Van De Wetering, and J. J. Van Wijk, "Visualization, selection, and analysis of traffic flows," *IEEE Trans. Vis. Comput. Graphics*, vol. 22, no. 1, pp. 379–388, Jan. 2016.
- [41] R. Scheepens, N. Willems, H. van de Wetering, G. Andrienko, N. Andrienko, and J. J. van Wijk, "Composite density maps for multivariate trajectories," *IEEE Trans. Vis. Comput. Graphics*, vol. 17, no. 12, pp. 2518–2527, Dec. 2011.
- [42] J. Stahnke, M. Dörk, B. Müller, and A. Thom, "Probing projections: Interaction techniques for interpreting arrangements and errors of dimensionality reductions," *IEEE Trans. Vis. Comput. Graphics*, vol. 22, no. 1, pp. 629–638, Jan. 2016.

- [43] G. Sun, R. Liang, H. Qu, and Y. Wu, "Embedding spatio-temporal information into maps by route-zooming," *IEEE Trans. Vis. Comput. Graphics*, vol. 23, no. 5, pp. 1506–1519, May 2017.
- [44] G. Sun, R. Liang, F. Wu, and H. Qu, "A Web-based visual analytics system for real estate data," *Sci. China Inf. Sci.*, vol. 56, no. 5, pp. 1–13, May 2013.
- [45] C. Tominski, H. Schumann, G. Andrienko, and N. Andrienko, "Stacking-based visualization of trajectory attribute data," *IEEE Trans. Vis. Comput. Graphics*, vol. 18, no. 12, pp. 2565–2574, Dec. 2012.
- [46] T. von Landesberger, F. Brodtkorb, P. Roskosch, N. Andrienko, G. Andrienko, and A. Kerren, "MobilityGraphs: Visual analysis of mass mobility dynamics via spatio-temporal graphs and clustering," *IEEE Trans. Vis. Comput. Graphics*, vol. 22, no. 1, pp. 11–20, Jan. 2016.
- [47] F. Wang *et al.*, "A visual reasoning approach for data-driven transport assessment on urban roads," in *Proc. IEEE Conf. Vis. Analytics Sci. Technol. (VAST)*, Oct. 2014, pp. 103–112.
- [48] L. Wang, Y. Zheng, X. Xie, and W.-Y. Ma, "A flexible spatio-temporal indexing scheme for large-scale GPS track retrieval," in *Proc. 9th Int. Conf. Mobile Data Manage. (MDM)*, Apr. 2008, pp. 1–8.
- [49] Z. Wang, M. Lu, X. Yuan, J. Zhang, and H. Van De Wetering, "Visual traffic jam analysis based on trajectory data," *IEEE Trans. Vis. Comput. Graphics*, vol. 19, no. 12, pp. 2159–2168, Dec. 2013.
- [50] G. I. Webb, "Discovering associations with numeric variables," in *Proc. 7th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining (KDD)*, 2001, pp. 383–388.
- [51] D. Weng, R. Chen, J. Zhang, J. Bao, Y. Zheng, and Y. Wu, "Pareto-optimal transit route planning with multi-objective monte-carlo tree search," *IEEE Trans. Intell. Transp. Syst.*, early access, Feb. 13, 2020, doi: 10.1109/TITS.2020.2964012.
- [52] D. Weng, H. Zhu, J. Bao, Y. Zheng, and Y. Wu, "HomeFinder revisited: Finding ideal homes with reachability-centric multi-criteria decision making," in *Proc. Conf. Hum. Factors Comput. Syst. (CHI)*, R. L. Mandryk, M. Hancock, M. Perry, and A. L. Cox, Ed., Montreal, QC, Canada, Apr. 2018, p. 247.
- [53] N. Willems, H. van de Wetering, and J. J. van Wijk, "Visualization of vessel movements," *Comput. Graph. Forum*, vol. 28, no. 3, pp. 959–966, Jun. 2009.
- [54] P. Chung Wong, P. Whitney, and J. Thomas, "Visualizing association rules for text mining," in *Proc. IEEE Symp. Inf. Vis. (InfoVis)*, Oct. 1999, pp. 120–128.
- [55] W. Wu *et al.*, "TelCoVis: Visual exploration of co-occurrence in urban human mobility based on telco data," *IEEE Trans. Vis. Comput. Graphics*, vol. 22, no. 1, pp. 935–944, Jan. 2016.
- [56] M. Xu *et al.*, "Traffic simulation and visual verification in smog," *ACM Trans. Intell. Syst. Technol.*, vol. 10, no. 1, pp. 3:1–3:17, Jan. 2019.
- [57] L. Yang, "Visualizing frequent itemsets, association rules, and sequential patterns in parallel coordinates," in *Proc. ICCSA*, 2003, pp. 21–30.
- [58] J. Zhang *et al.*, "Visual analysis of public utility service problems in a metropolis," *IEEE Trans. Vis. Comput. Graphics*, vol. 20, no. 12, pp. 1843–1852, Dec. 2014.
- [59] Z. Zhang and W. Wu, "Composite spatio-temporal co-occurrence pattern mining," in *Proc. WASA Berlin*, Germany: Springer, 2008, pp. 454–465.
- [60] Y. Zheng *et al.*, "TelcoFlow: Visual exploration of collective behaviors based on telco data," in *Proc. IEEE Int. Conf. Big Data (Big Data)*, Dec. 2016, pp. 843–852.
- [61] Z.-H. Zhou, "Abductive learning: Towards bridging machine learning and logical reasoning," *Sci. China Inf. Sci.*, vol. 62, no. 7, pp. 76101:1–76101:3, Jul. 2019.



Di Weng received the B.S. degree in computer science from the Taishan Honored College, Shandong University, in 2016. He is currently pursuing the Ph.D. degree with the State Key Lab of CAD&CG, Zhejiang University. His research interests mainly include the data mining, visualization, and visual analytics of large-scale urban data.

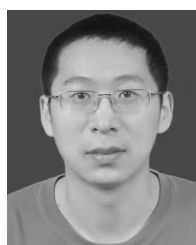


Zikun Deng received the B.S. degree in transportation engineering from Sun Yat-sen University in 2016. He is currently pursuing the Ph.D. degree with the State Key Lab of CAD&CG, Zhejiang University. His research interests mainly include spatiotemporal data mining, visualization, and urban visual analytics.



Jie Bao received the Ph.D. degree from the Department of Computer Science and Engineering, University of Minnesota at Twin Cities, in 2014.

He leads the Data Management Department, JD Intelligent City Business Unit, where he is in charge of the design and development of JD Urban Spatio-Temporal data engine (aka JUST), as well as all the data-centric products in the BU. Before joining JD, he was a Researcher with MSRA, from 2014 to 2017. He has published over 40 research articles in refereed journals and conferences (e.g., SIGKDD, ICDE, AAAI, VLDB, SIGMOD, and SIGSPATIAL). His main research interests include urban computing, spatio-temporal data management/mining, and distributed computing platforms.



Mingliang Xu received the Ph.D. degree in computer science and technology from the State Key Lab of CAD&CG, Zhejiang University, Hangzhou, China. He is currently a Full Professor with the School of Information Engineering, Zhengzhou University. He is also the Director of the Center for Interdisciplinary Information Science Research (CIISR) and the Vice General Secretary of ACM SIGAI China. His current research interests include computer graphics and artificial intelligence. He has authored more than 80 journals and conference papers in these areas, including ACM TOG, ACM TIST, the IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE (TPAMI), the IEEE TRANSACTIONS ON IMAGE PROCESSING (TIP), the IEEE TRANSACTIONS ON CYBERNETICS (TCYB), the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY (TCSVT), the IEEE TRANSACTIONS ON AUTOMATIC CONTROL (TAC), the IEEE TRANSACTIONS ON COMPUTATIONAL INTELLIGENCE AND AI IN GAMES (TCIAIG), ACM SIGGRAPH (Asia), ACM MM, and IJCAI.



Yingcai Wu received the Ph.D. degree in computer science from The Hong Kong University of Science and Technology. He was a Post-Doctoral Researcher with the University of California, Davis, from 2010 to 2012, and a Researcher with Microsoft Research Asia, from 2012 to 2015. He is currently a ZJU100 Young Professor with the State Key Lab of CAD&CG, Zhejiang University. His main research interests are in information visualization and visual analytics, with focuses on urban computing, sports science, immersive visualization, and social media analysis.



Zhangye Wang received the Ph.D. degree from Zhejiang University in 2003. He is currently an Associate Professor with the State Key Lab of CAD & CG, Zhejiang University. His research interests include visual analytics and computer graphics.



Yu Zheng is currently the Vice President of JD.COM, leading the Intelligent Cities Business Unit and JD Intelligent Cities Research. He also serves as the Chief Data Scientist at JD Digits, passionate about using big data and AI technology to tackle urban challenges. Before joining JD.COM, he was a Senior Research Manager with Microsoft Research. He is also a Chair Professor with Shanghai Jiao Tong University and an Adjunct Professor with The Hong Kong University of Science and Technology. He is also a Keynote Speaker of AAAI

2019, KDD 2019 Plenary Keynote Panel, and IJCAI 2019 Industrial Days. His monograph, entitled *Urban Computing*, has been used as the first text book in this field. In 2013, he was named one of the Top Innovators under 35 by MIT Technology Review (TR35) and featured by Time Magazine for his research on urban computing. In 2014, he was named one of the Top 40 Business Elites under 40 in China by Fortune Magazine. In 2017, he is honored as an ACM Distinguished Scientist. He has served as the Chair on over 10 prestigious international conferences, e.g., as the Program Co-Chair of ICDE 2014 (Industrial Track), CIKM 2017 (Industrial Track), and IJCAI 2019 (industrial track). He serves as the Editor-in-Chief of *ACM Transactions on Intelligent Systems and Technology*.



Zhiyu Ding received the Ph.D. degree from the State Key Lab of CAD&CG, Zhejiang University. He is currently the Director of visualization and visual computing, the Manager of WEB Competence Center of EI Department, Cloud BU, Huawei. He mainly focuses on visualization, HCI and graphics, and SAAS applications which utilize these mentioned technologies.



Wei Chen received the Ph.D. degree from Zhejiang University in 2002. He is currently a Professor with the State Key Lab of CAD & CG, Zhejiang University. His research interests include visualization, visual analytics, and biomedical image computing.