

# Reference-Based Framework for Spatio-Temporal Trajectory Compression and Query Processing

Kai Zheng<sup>✉</sup>, *Member, IEEE*, Yan Zhao<sup>✉</sup>, Defu Lian<sup>✉</sup>, Bolong Zheng, Guanfeng Liu<sup>✉</sup>, and Xiaofang Zhou, *Fellow, IEEE*

**Abstract**—The pervasiveness of GPS-enabled devices and wireless communication technologies results in massive trajectory data, incurring expensive cost for storage, transmission, and query processing. To relieve this problem, in this paper we propose a novel framework for compressing trajectory data, REST (Reference-based Spatio-temporal trajectory compression), by which a raw trajectory is represented by concatenation of a series of historical (sub-)trajectories (called reference trajectories) that form the compressed trajectory within a given spatio-temporal deviation threshold. In order to construct a reference trajectory set that can most benefit the subsequent compression, we propose three kinds of techniques to select reference trajectories wisely from a large dataset such that the resulting reference set is more compact yet covering most footprints of trajectories in the area of interest. To address the computational issue caused by the large number of combinations of reference trajectories that may exist for resembling a given trajectory, we propose efficient greedy algorithms that run in the blink of an eye and dynamic programming algorithms that can achieve the optimal compression ratio. Compared to existing work on trajectory compression, our framework has few assumptions about data such as moving within a road network or moving with constant direction and speed, and better compression performance with fairly small spatio-temporal loss. In addition, by indexing the reference trajectories directly with an in-memory R-tree and building connections to the raw trajectories with inverted index, we develop an extremely efficient algorithm that can answer spatio-temporal range queries over trajectories in their compressed form. Extensive experiments on a real taxi trajectory dataset demonstrate the superiority of our framework over existing representative approaches in terms of both compression ratio and efficiency.

**Index Terms**—Reference trajectory, spatio-temporal trajectory, compression

## 1 INTRODUCTION

THE prevalent use of various mobile devices, such as smart-phones, on-board diagnostics, personal navigation devices, and wearable smart devices, has resulted in massive amount of trajectory data. While they contain a wealth of mobility information and offer great opportunities for heightening our understanding about human mobilities,

transmitting and storing raw trajectory data consumes too much network bandwidth and storage capacity [1]. This is calling for effective trajectory compression techniques that can remove redundant and inessential samples from raw trajectory data to reduce the storage cost while preserving the utility of data.

Trajectory data compression approaches can be generally divided into two categories: spatial and spatio-temporal compression. Treating trajectories as polylines, spatial compression methods are also known as line simplification algorithms (e.g., Douglas-Peucker (DP) [2] and Bellman's algorithm [3]), which discard some samples within a given spatial deviation threshold from its original locations. However, trajectories are spatio-temporal records of moving objects, in which temporal information is also critical in many applications such as trajectory monitoring [4] and location tracking [5]. Ignoring temporal information during compression may produce unbounded erroneous results when querying the decompressed data. Therefore, recent studies focus on spatio-temporal compression algorithms [1], [6], [7], [8], [9], which adopt spatio-temporal criteria to bound the compression loss. The common feature of the above approaches is that they just exploit the spatio-temporal characteristics of the single trajectory to be compressed and assume moving objects do not change speed and/or direction frequently while traveling. However, this is quite an optimistic assumption for objects

- K. Zheng is with the School of Computer Science and Engineering, University of Electronic Science and Technology of China, Chengdu 610054, China, and also with Zhejiang Lab, Hangzhou 310000, China. E-mail: zhengkai@uestc.edu.cn.
- Y. Zhao is with the Institute of Artificial Intelligence, School of Computer Science and Technology, Soochow University, Suzhou 215006, China, and also with Zhejiang Lab, Hangzhou 310000, China. E-mail: zhaoyan@suda.edu.cn.
- D. Lian is with the School of Computer Science and Technology and School of Data Science, University of Science and Technology of China, Hefei 230026, China. E-mail: dove.ustc@gmail.com.
- B. Zheng is with the School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan 430074, China. E-mail: bolongzheng@hust.edu.cn.
- G. Liu is with the Department of Computing, Macquarie University, Sydney NSW2109, Australia. E-mail: guanfeng.liu@mq.edu.au.
- X. Zhou is with the University of Queensland, Brisbane Q4072, Australia, and also with Zhejiang Lab, Hangzhou 310000, China. E-mail: zxf@itee.uq.edu.au.

Manuscript received 8 July 2018; revised 15 Apr. 2019; accepted 17 Apr. 2019. Date of publication 2 May 2019; date of current version 6 Oct. 2020. (Corresponding author: Yan Zhao.) Recommended for acceptance by J. M. Phillips. Digital Object Identifier no. 10.1109/TKDE.2019.2914449

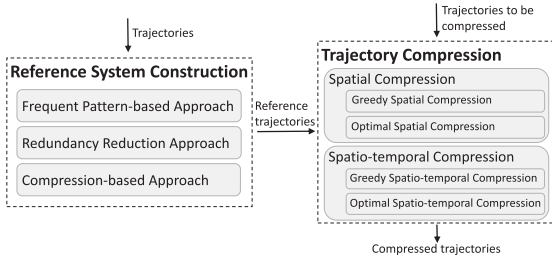


Fig. 1. REST framework overview.

moving with complicated traffic condition, which is why those algorithms cannot achieve high compression ratio on real-world trajectory data. Recently, Song et al [10] propose a data-driven approach, called PRESS, that leverages shortest path and frequent trajectory pattern to compress trajectories in a road network. Nevertheless, there are two key prerequisites for PRESS to work properly. First, object movements must be constrained by a network, which does not apply for a wide range of free-space moving objects such as animals, flying objects, handwriting trajectories and so on. Second, the network should be relatively stable so that the compressed trajectories can be used properly. However, it is not uncommon the topology of road network changes frequently in developing areas (e.g., major cities of China).<sup>1</sup> Since PRESS relies on precomputed all-pair shortest paths, it has to recompute a large number of shortest paths every time the network updates, which is a time consuming process. Moreover, it takes tremendous space to store the all-pair shortest paths for each version of road network.<sup>2</sup>

In this paper, we propose a completely data-driven framework, called REST (Reference-based Spatio-temporal trajectory compression), for trajectory compression. There are a few advantages for REST compared to existing work. First, trajectories can originate from any kind of space (either constrained or non-constrained space). Second, it includes both spatial-only and spatio-temporal compression algorithms. Third, it only takes small amount of memory to store the auxiliary information (i.e., reference set). Fourth, the trajectory data are indexable and usable in their compressed form. This framework is based on some prior studies that find human mobilities have inherent high-level spatial and temporal regularity [11] (i.e., people have high probability to repeat similar travel patterns) and highly skewed travel distribution [12] (i.e., different people often take similar routes when traveling between certain locations). Therefore it is feasible to extract a relatively small collection of trajectories, named reference trajectories, which “covers” most trajectories in the region of interest. Then given a new trajectory in the same area, there is high chance we can use a proper concatenation of reference (sub-)trajectories to resemble, at least partially, the given one. Compared with the raw format that keeps every location sample, this representation saves significant space since only a series of identifiers and offsets of the reference (sub-)trajectories need to be recorded.

As shown in Fig. 1, the REST framework is comprised of two components: reference set construction and reference-based compression. The first component aims to build a reference system, where the challenge is how to trade off high coverage and low redundancy in the reference set such that subsequent compression can be performed more effectively and efficiently. To this end, we present three kinds of approaches including Frequent Pattern-based Approach, Redundancy Reduction Approach and Compression-based Approach, which use different strategies to select a compact yet expressive reference set from a much larger but more redundant training dataset. The second component needs to tackle the computational issue in the great number of reference trajectories we can use to represent a given trajectory. For the sake of efficiency, we propose greedy algorithms that try to represent the longest possible sequence of samples with a single reference trajectory. We also develop optimal algorithms to calculate the minimal storage cost of the compressed trajectory and obtain the corresponding optimal combination of reference trajectories.

Though our preliminary work [13] has already optimally compressed the spatio-temporal information of trajectory data with reference trajectories, it fails to demonstrate the key utility of compressed trajectory data. Trajectory utility mainly depends on the effective and efficient trajectory query processing in trajectory databases, which aims to evaluate the spatio-temporal relationships among spatial data objects. However, it is a challenging task to answer the trajectory queries efficiently due to the inherent difficulties in indexing trajectories as well as the new complexity introduced by the compressed trajectories, which are in the form of sequence of reference trajectories. To tackle this problem, we develop an effective hybrid index structure to support efficient query processing over compressed trajectories without fully decompressing the data and then present the query processing based on the index structure in Section 6.

As a summary, the major value-added extension over our preliminary work [13] in REST framework are four folds.

- 1) We demonstrate REST can support the classical spatio-temporal queries (e.g., range queries) in trajectory databases without fully recovering the compressed trajectories.
- 2) We design an effective hybrid index structure, consisting of both R-tree and inverted index, to support efficient query processing over the compressed trajectories.
- 3) We give more justifications about the custom defined error metrics and discussions on the applicability of our REST framework to other trajectory similarity measures.
- 4) An extensive experimental study is conducted on two real taxi trajectory datasets to validate the effectiveness and efficiency of the query processing on the compressed trajectories.

The remainder of this paper is organized as follows. Section 2 introduces preliminary concepts, error metric and reference-based compression problem. We then present the construction of reference trajectory set in Section 3. The algorithms for spatial-only compression and spatio-temporal compression are presented in Section 4 and Section 5

1. TomTom claims their digital maps have fixes and updates every week. [https://www.tomtom.com/en\\_au/mydrive-connect/](https://www.tomtom.com/en_au/mydrive-connect/)

2. Keeping all-pair shortest path requires  $O(|V|^2)$  space where  $V$  represents the vertex set in a road network.

respectively, followed by the extension of query processing for the compressed trajectories in Section 6. We report the results from empirical study in Section 7 and survey the related work in Section 8. Section 9 concludes this paper.

## 2 PROBLEM STATEMENT

In this section, we will present a set of preliminary concepts, introduce the error metric between raw and compressed trajectories, and finally state our problem and goal.

### 2.1 Preliminary

**Definition 1 (Raw Trajectory).** A raw trajectory of a moving object in 2D euclidean plane, denoted as  $T$ , is a finite sequence of timestamped locations with the form of  $((x_1, y_1, t_1), (x_2, y_2, t_2), \dots, (x_n, y_n, t_n))$ , where  $(x_i, y_i)$  stands for the longitude and latitude information of a sampled location at time stamp  $t_i$ .

**Definition 2 (Sub-Trajectory).** A sub-trajectory, denoted by  $T^{(i,j)}$ , is made of consecutive sample points of  $T$  from the  $i$ th to  $j$ th triplet, i.e.,  $T^{(i,j)} : ((x_i, y_i, t_i), \dots, (x_j, y_j, t_j))$ .

Trajectory compression is a process to reduce storage cost while keeping utility of a trajectory. Normally there are two measures when comparing trajectory compression methods:

- 1) *Compression Ratio* measures how much space has been saved by compressing raw trajectories. It is usually defined as the ratio between space costs of raw trajectories and compressed trajectories, i.e.,  $CR = \frac{\text{space}(T)}{\text{space}(T')}$ .
- 2) *Compression Loss* measures to what extent a compressed trajectory can be reconstructed to its corresponding raw format. It is usually quantified by a distance between raw trajectory and decompressed trajectory based on some predefined distance function.

These two factors are often trade-off: high compression ratio usually leads to greater compression loss and vice versa. Based on compression loss, trajectory compression can be classified into *lossless compression* and *lossy compression*. Due to the extremely fine granularity (hence almost infinite cardinality) of spatio-temporal dimensions in free space, lossless compression algorithms either are not practical or have extremely low compression ratio. Thus in this paper we resort to *bounded lossy compression* algorithm.

**Definition 3 (Bounded Lossy Compression).** Given a deviation threshold  $\epsilon$ , an  $\epsilon$ -Bounded Lossy Compression algorithm transforms a raw trajectory  $T$  into a compressed trajectory  $T'$ , such that the distance between the reconstructed trajectory  $T^*$  and  $T$  does not exceed  $\epsilon$ , i.e.,  $d(T, T^*) \leq \epsilon$ , where  $d$  is some predefined distance function for trajectories.

### 2.2 Error Metric

In this paper, we propose a simple but effective variant of Dynamic Time Warping (DTW) [14] distance, called MaxDTW, to serve the error metric. It works exactly the same way as DTW in looking for the best alignment between two unsynchronized trajectories, with the only exception that it just needs to record the maximum distance among all matched pairs instead of the sum. More formally,

**Definition 4 (MaxDTW).** Given two trajectories  $T_a = (p_1, p_2, \dots, p_n)$  and  $T_b = (q_1, q_2, \dots, q_m)$ , the MaxDTW distance between them is defined as follows:

$$\text{MaxDTW}(T_a, T_b) = \begin{cases} 0, & \text{if } T_a = T_b = \emptyset \\ +\infty, & \text{if } T_a = \emptyset \text{ or } T_b = \emptyset \\ \max\{d(p_n, q_m), Q(p_n, q_m)\}, & \text{otherwise} \end{cases}$$

$$Q(p_n, q_m) = \min \begin{cases} \text{MaxDTW}(T_a^{(1,n-1)}, T_b^{(1,m-1)}) \\ \text{MaxDTW}(T_a^{(1,n-1)}, T_b^{(1,m)}) \\ \text{MaxDTW}(T_a^{(1,n)}, T_b^{(1,m-1)}) \end{cases},$$

where  $d(a, b)$  is a given distance between point  $a$  and  $b$ .

Similar to DTW, we can use a dynamic programming algorithm [14] to compute MaxDTW.

While our framework can be applied to both conventional DTW and the newly defined MaxDTW (will be discussed in Section 4.1), we chose to define a custom metric MaxDTW because this is easier for user to set the error threshold. To see the reason, DTW is calculated based on distance aggregation (sum of distance), so it is dependent on not only the closeness between two trajectories but also their size (i.e., number of sample points). Therefore, when a user wants to compress a set of trajectories, she needs to set a different error threshold for each individual trajectory as they are of different sizes. For instance, we should set a lower error threshold for a trajectory with 10 points and higher error threshold for another trajectory with 100 points if we expect them to be compressed with similar quality. Of course, this problem can be solved by using the average DTW, i.e.,  $\text{DTW}/(\text{size of trajectory})$ , which is essentially the same as our proposed MaxDTW. The rational behind MaxDTW is the furthest pair-wise distance between two trajectories in their best DTW alignment, which is independent of the size. This also gives the user an intuitive view of the closeness they should expect to see between the decompressed trajectory and its corresponding raw trajectory.

### 2.3 Reference-Based Compression

As observed in previous studies [12], there is strong bias when most drivers plan their routes, which means given a new trajectory it is very likely to find from a historical trajectory dataset a few trajectories that resemble, at least partially, the given one. We name these trajectory set as *reference trajectory set*, denoted by  $R$ , which can be generated from a historical trajectory dataset in the region of interest (e.g., where the trajectories to be compressed also reside). While it is difficult to define the optimality of  $R$ , there are two qualitative measures for a good one—high coverage and low redundancy. Here *high coverage* means it has enough power to represent a given trajectory in the same area of interest, which heavily affects the compression ratio. *Low redundancy* means most trajectories in  $R$  are quite unique in terms of their geographical locations since overlapping reference trajectories do not increase the expressive power and make the compression inefficient. In the rest of the paper we use *reference trajectory set* and *reference set* interchangeably when no ambiguity is caused.

**Problem Statement.** Given a reference trajectory set  $R$ , a trajectory  $T$  to be compressed and an error threshold  $\epsilon$ , a reference-based compression algorithm uses a selected



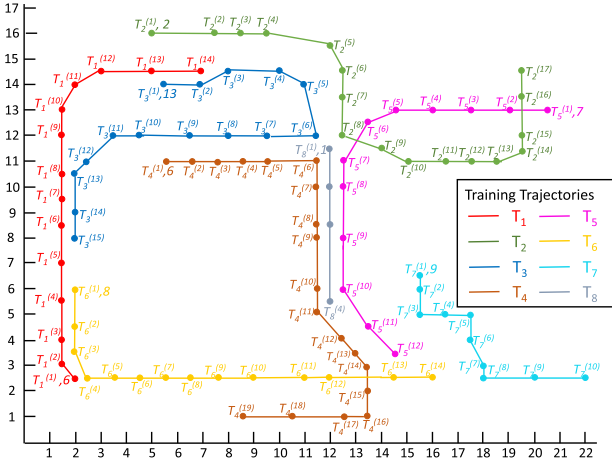


Fig. 2. Running example.

subset of  $R$ , or their sub-trajectories whenever possible, to represent  $T$ , denoted as  $T'$ , and guarantees that the distance between  $T'$  and  $T$  does not exceed  $\epsilon$ .

### 3 REFERENCE SET CONSTRUCTION

In this section, we propose three methods to build a compact and expressive reference set, which is relatively stable and do not require frequent updates to compress new data. Fig. 2 shows a training trajectory set, in which each trajectory sample is labeled by  $T_i^{(j)}$  indicating the  $j$ th sample of trajectory  $T_i$ . Besides, the first sample of each trajectory is companied by a number that indicates the start time stamp. We assume the sampling interval is 1 time unit. For instance,  $T_2^{(2)}$  is the second sample of  $T_2$  whose time stamp is 3.

#### 3.1 Frequent Pattern-Based Approach (FPA)

Previous studies have shown that trajectories of moving objects often follow certain patterns, such as commuter patterns, peak/off peak patterns, weekend patterns, etc. Therefore a natural thought is to leverage these frequent patterns for building a reference set. Given a trajectory dataset, its frequent pattern is a set of sub-trajectories of which the occurrence frequencies exceed a certain support threshold. Inspired by [15], we first introduce *calculating point* and *calculating trajectory* to discretize trajectories into sequential data and then apply Sequential Pattern Mining algorithms [16] to find frequent sub-trajectories. Specifically, given a sample point set  $P = \{p_1, p_2, \dots, p_n\}$ ,  $p_i$  is called a calculating point  $a$  if  $|p_j.x - p_i.x| \leq \epsilon_s$  and  $|p_j.y - p_i.y| \leq \epsilon_s$ . Then all  $p_j$  can be represented by  $a$ . A calculating trajectory is a sequence formed by the calculating points chronologically. The following steps are performed to find all of the frequent sub-trajectories:

- 1) Find all calculating points and calculating trajectories.
- 2) Given the minimum support threshold, the frequent calculating points are obtained by scanning all calculating points.
- 3) Remove non-frequent calculating points, and obtain frequent calculating sub-trajectories.
- 4) Remove the frequent calculating sub-trajectories who is sub-trajectories of another frequent calculating trajectories.
- 5) Return all the frequent sub-trajectories.

### 3.2 Redundancy Reduction Approach

Since only using frequent patterns may result in low coverage, we next present two variant methods to achieve higher coverage and reduce redundancy (to some extent) simultaneously.

#### 3.2.1 Segment Redundancy Reduction (SRR)

Given a training trajectory set, we aim to extract a set of non-redundant sub-trajectories. First we define redundant segment in below,

**Definition 5 (Redundant Segment).** Given a minimum length threshold  $\eta$  and a distance threshold  $\epsilon_s$ , two sub-trajectories (segments) with the same number of samples, i.e.,  $T_a^{(i,i+m)}$  and  $T_b^{(j,j+m)}$ , are said to overlap with each other if  $m \geq \eta$  and their maximum pairwise distance  $d_{max} = \max_{0 \leq k \leq m} d(T_a.p_{i+k}, T_b.p_{j+k}) \leq \epsilon_s$ . If a sub-trajectory  $s$  overlaps with any existing sub-trajectory in a reference set  $R$ ,  $s$  is called a redundant segment.

$\eta$  is to avoid the existence of too many short segments. The basic idea of our approach is to eliminate all the redundant segments of training trajectories and use remaining segments as the reference set.

#### 3.2.2 Trajectory Redundancy Reduction (TRR)

The reference set constructed by SRR algorithm may end up with too many short segments if  $\eta$  is too small, or too many whole trajectories otherwise (since it gets harder to identify long segment overlap). We present an alternative approach to reduce redundancy by treating each trajectory as atomic, i.e., either use the entire one or nothing. The redundant trajectory is defined as follows:

**Definition 6 (Redundant Trajectory).** Given an overlap threshold  $\theta$ , a distance threshold  $\epsilon_s$ , a trajectory  $T$  is called redundant if the overlap portion between  $T$  and  $R$ , denoted as  $L(T, R)$ , exceeds  $\theta$ , where  $L(T, R)$  is calculated as the portion of samples in  $T$  that are sufficiently close to any samples in  $R$ , i.e.,

$$L(T, R) = \frac{|p \in T | \exists q \in R, d(p, q) \leq \epsilon_s|}{|T|} > \theta. \quad (1)$$

The TRR algorithm checks every training trajectory  $T$  and calculates the overlap portion with  $R$ . If the value is below  $\theta$ ,  $T$  is added into  $R$  as a reference trajectory.

### 3.3 Compression-Based Approach (CA)

As the ultimate goal of building a reference set is to achieve high compression ratio, we can compress a training trajectory against the current reference set with the spatial compression algorithm proposed in the following section and record the compression ratio. If the ratio is high enough, that means the training trajectory can be well described by existing reference trajectories, i.e., it is redundant; otherwise, it is non-redundant. The non-redundant training trajectory will be added into the current reference set.

## 4 SPATIAL COMPRESSION

In this section, ignoring the time information, we compress a given trajectory using as few reference trajectories as possible to minimize the space cost.

### 4.1 Matchable Reference Trajectory

We will first introduce *matchable reference trajectory* – a basic concept that will be used throughout our algorithms.

#### Definition 7 (Matchable Reference Trajectory (MRT)).

Given a sub-trajectory  $T^{(i,j)}$  and a spatial deviation threshold  $\epsilon_s$ , its matchable reference trajectory set, denoted as  $M(T^{(i,j)})$ , includes all the reference sub-trajectories with less-than- $\epsilon_s$  MaxDTW distance with  $T^{(i,j)}$ , i.e.,

$$M(T^{(i,j)}) = \left\{ T^{(k,g)} \mid T \in R, 1 \leq k \leq g \leq |T|, \right. \\ \left. \text{MaxDTW}(T^{(i,j)}, T^{(k,g)}) \leq \epsilon_s \right\}. \quad (2)$$

Here each MRT  $T^{(k,g)}$  is recorded as a triplet  $(T.id, k, g)$  (costing 8 bytes). To retrieve the MRTs, we propose an efficient method based on the following observation.

**Lemma 1.** Any sub-trajectory of the MRT of  $T^{(i,j)}$  is also an MRT of sub-trajectory of  $T^{(i,j)}$ .

The proof is directly from the property that MaxDTW distance between longer sequences upper bounds their sub-sequences. Therefore the MRT set of  $T^{(i,j)}$  can be more easily derived by joining the MRT sets of its sub-trajectories. Algorithm 1 is proposed based on this intuition. First, the MRT sets of all the length-2 sub-trajectories is obtained and added to a hash set  $M(T^{(i,j)})$  (lines 1-2). Then for each length- $n$  ( $2 < n \leq |T|$ ) sub-trajectory  $T^{(i,j)}$ , we check if existing MRT of  $T^{(i,j-1)}$  and  $T^{(j-1,j)}$  can also be the MRT of  $T^{(i,j)}$  (lines 6-9), and verify if a longer MRT can be formed for  $T^{(i,j)}$  by joining  $M(T^{(i,j-1)})$  and  $M(T^{(j-1,j)})$  (lines 10-11). The algorithm can be terminated early if we find the MRT sets of all length- $n$  sub-trajectories are empty since all longer sub-trajectories will not have MRT based on Lemma 1. Given  $\epsilon_s = 0.9, \eta = 1$  and the reference set generated by SRR algorithm previously, i.e.,  $R = \{T_1, T_2, T_3^{(3,15)}, T_4, T_5, T_6, T_7\}$ , Fig. 3 illustrates the MRTs for  $T$ , to name a few,  $M(T^{(2,9)}) = \{T_2^{(1,8)}\}$ ,  $M(T^{(10,15)}) = \{T_4^{(6,14)}, T_5^{(7,12)}\}$ .

#### Algorithm 1. Matchable Reference Trajectory Search

---

**Input:**  $T, R, \epsilon_s$   
**Output:**  $M$

- 1 **for each**  $T^{(i,i+1)} \in T$  **do**
- 2    $M(T^{(i,i+1)}) \leftarrow$  MRT set for segment  $T^{(i,i+1)}$ ;
- 3 **for**  $n \leftarrow 3$  **to**  $|T|$  **do**
- 4   **for each** length- $n$  sub-trajectory  $T^{(i,j)} \in T$  **do**
- 5     **for**  $T_a^{(m,n)}, T_b^{(s,t)} \in M(T^{(i,j-1)}), M(T^{(j-1,j)})$  **do**
- 6       **if**  $\text{MaxDTW}(T^{(i,j)}, T_a^{(m,n)}) \leq \epsilon_s$  **then**
- 7         Add  $T_a^{(m,n)}$  into  $M(T^{(i,j)})$ ;
- 8       **if**  $\text{MaxDTW}(T^{(i,j)}, T_b^{(s,t)}) \leq \epsilon_s$  **then**
- 9         Add  $T_b^{(s,t)}$  into  $M(T^{(i,j)})$ ;
- 10      **if**  $a = b$  and  $n = s$  **then**
- 11         Add  $T_a^{(m,t)}$  into  $M(T^{(i,j)})$ ;
- 12   **if no** length- $n$  sub-trajectory has MRT **then**
- 13     Break;
- 14 **return**  $M$ ;

---

*Discussion.* here we briefly discuss the applicability of our framework to other trajectory similarity measure. Since the correctness of our following algorithms rely on Lemma 1, which states the sub-structure optimality in MRT, we only need to examine whether the derived MRT (from Definition 7)

based on a given trajectory similarity measure satisfies Lemma 1 or not. Obviously, our framework can be easily adapted to those distance aggregation based trajectory similarity measure such as euclidean Distance (ED), Edit Distance with Projections (EDwP) [17], Dynamic Time Warping [14] and its variants as they all satisfy Lemma 1. However, the generalization to Longest Common Subsequence (LCSS) [18], Edit distance with Real Penalty (ERP) [19] and Edit Distance on Real sequences (EDR) [20] is not straightforward since they are count-based similarity measure, which means a sub-trajectory of MRT may not meet the similarity threshold any more.

### 4.2 Greedy Spatial Compression

Once the MRT set is obtained, a natural thought is to process the given trajectory in its chronological order and compress the longest possible sub-trajectory with its MRT (selecting an arbitrary MRT if multiple longest MRTs exist), until the last sample point has been reached. This approach is called Greedy Spatial Compression (GSC) algorithm since it seems not a global strategy to combine MRTs in order to achieve the minimal storage cost. However, we will prove later it also yields space optimal compressed trajectory.

Consider the example in Fig. 3. With GSC algorithm, we first compress  $T^{(1,3)}$  with  $T_1^{(12,14)}$  since it is the first longest sub-trajectory with non-empty MRT set. Then the remaining part of  $T$  are represented by  $T_2^{(3,8)}, T_4^{(6,15)}, T^{(17)}, T_7^{(8,10)}$  respectively, resulting in the space cost of 40 bytes, i.e., 25 percent of its original space (160 bytes).

### 4.3 Optimal Spatial Compression

In the sequel, we propose a dynamic programming algorithm, called Optimal Spatial Compression (OSC), that aims to minimize the required storage size for  $T'$ . Specifically, given a trajectory  $T$  and its MRT set  $M(T)$ , we define  $F_T[i]$  as the minimum storage size needed for compressing  $T^{(1,i)}$ , and  $T'^{(1,i)}$  as the corresponding compressed sub-trajectory to achieve this optimum storage.  $F_T[i]$  can be derived by the following recursive formula shown in Equation (3).

$$F_T[i] = \begin{cases} 0 & \text{if } i = 0 \\ \min_{1 \leq j \leq i \wedge M(T^{(j,i)}) \neq \emptyset} \{F_T[j-1] + 8\} & \text{otherwise} \end{cases} \quad (3)$$

where the MRT set of single sample point, i.e.,  $M(T^{(i,i)})$ , is manually set to non empty.

It is trivial  $F_T[0] = 0$  when  $T = \emptyset$ , i.e.,  $i = 0$ . When  $i > 0$ ,  $F_T[i]$  is computed by picking the minimum of: 1) the optimal storage cost of sub-trajectory  $T'^{(1,j-1)}$ , i.e.,  $F_T[j-1]$ , plus the storage (8 bytes) for an MRT of sub-trajectory  $T^{(j,i)}$  if  $M(T^{(j,i)}) \neq \emptyset$  when  $1 \leq j < i$ ; and 2) the optimal storage cost of sub-trajectory  $T'^{(1,i-1)}$ , i.e.,  $F_T[i-1]$ , plus the storage (8 bytes) of original sample point  $p_i \in T$  when  $j = i$ .

With Equation (3), now we can compute the minimum storage size of compressed trajectory, which is presented in Algorithm 2. Note that we introduce another notation  $pre[i]$  for recording the last-to-first sample points having been compressed before achieving  $F_T[i]$  to facilitate the reconstruction of the compressed trajectory  $T'$  with minimum storage size. It first initializes  $T' = null$  and  $F_T[0] = 0$  (lines 1-2). Then the algorithm processes all points of  $T$  in the sampling order





TABLE 1  
Compressed Trajectory from GSTC

$T$	$M_{opt}(T^{(i,j)})$	Time Correction	$C_T^{\epsilon_t}$
$T^{(1,3)}$	$T_1^{(12,14)}$	$T_1^{(12)}.t = 1$	4
$T^{(4,9)}$	$T_2^{(3,8)}$	$T_2^{(4)}.t = 6, T_2^{(7)}.t = 8$	8
$T^{(10,16)}$	$T_4^{(6,15)}$	$T_4^{(6)}.t = 10, T_4^{(9)}.t = 12, T_4^{(11)}.t = 13, T_4^{(13)}.t = 14$	16
$T^{(17)}$	$T^{(17)}$		0
$T^{(18,20)}$	$T_7^{(8,10)}$	$T_7^{(8)}.t = 18$	4

$$F_T[i] = \begin{cases} 0 & \text{if } i = 0 \\ \min \left\{ F_T[i-1] + 12, \min_{1 \leq j < i \wedge M(T^{(j,i)}) \neq \emptyset} \{ F_T[j-1] + C_{T^{(j,i)}}(M_{opt}(T^{(j,i)})) \} + 8 \right\} & \text{otherwise} \end{cases}$$

$F_T[i]$  records the minimal space needed for compressing sub-trajectory  $T^{(1,i)}$ .  $F_T[0] = 0$  defines the termination condition. For  $i \geq 1$ , the algorithm either 1) keeps the original sample  $p_i \in T$ ; or 2) compresses  $T^{(j,i)}$  with  $M_{opt}(T^{(j,i)})$ . In the first case, the final space cost is simply the compressed space of  $T^{(1,i-1)}$  plus 12 bytes (for storing  $p_i$ ). In the second case, the space cost is the minimal sum of compressed space for  $T^{(1,j-1)}$ , time correction cost of  $M_{opt}(T^{(j,i)})$  and space for  $M_{opt}(T^{(j,i)})$  (8 bytes). The process of calculating  $F_T[|T|]$  and construction of  $T'$  is similar to Algorithm 2 and thus omitted here due to space limitation. The time complexity of this algorithm is also  $O(|T|^2)$ .

Taking the example in Fig. 3, we employ OSTC algorithm to calculate the proper MRTs and time correction shown in Table 2, and achieve the minimal storage size of  $T'$  (i.e., 64 bytes), resulting in an approximately 15.79 percent reduction compared to GSTC algorithm.

## 6 QUERYING COMPRESSED TRAJECTORIES

The utility of trajectory compression depends on the space storage reduction as well as the effective and efficient trajectory query processing in trajectory databases. The proposed OSTC algorithm can compress the raw trajectories in such a way that the spatial and temporal information loss can be bounded with a deviation threshold,  $\epsilon$  (i.e.,  $\epsilon_s$  and  $\epsilon_t$ ), where  $\epsilon$  can be set as a very small value. This means the spatial paths and timestamps of a raw trajectory can be captured almost exactly when being compressed. Therefore we can not only directly decompress the compressed trajectories for location based service applications, but also use the compressed trajectories for various queries without fully decompressing the data, in which the error of every point in the partially

 TABLE 2  
Compressed Trajectory from OSTC

$T$	$M_{opt}(T^{(i,j)})$	Time Correction	$C_T^{\epsilon_t}$
$T^{(1)}$	$T^{(1)}$		0
$T^{(2,9)}$	$T_2^{(1,8)}$	$T_2^{(4)}.t = 6, T_2^{(7)}.t = 8$	8
$T^{(10,15)}$	$T_5^{(7,12)}$	$T_5^{(7)}.t = 10$	4
$T^{(16,17)}$	$T_6^{(13,14)}$	$T_6^{(13)}.t = 16$	4
$T^{(18,20)}$	$T_7^{(8,10)}$	$T_7^{(8)}.t = 18$	4

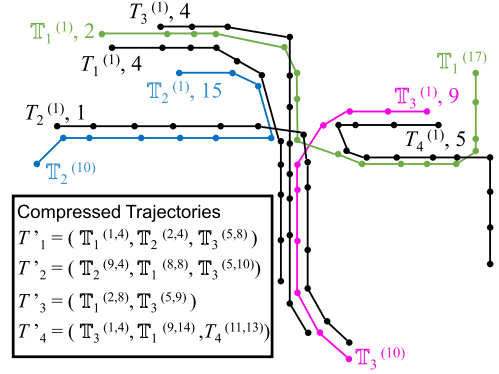


Fig. 4. A set of raw trajectories and reference trajectories.

decompressed trajectories is less than  $\epsilon$ . In this section, we focus on demonstrating that the compressed trajectories can support the range query, a typical trajectory query aiming to evaluate spatio-temporal relationships among spatial data objects. Informally, in Fig. 4, given a raw trajectory set (i.e.,  $\{T_1, T_2, T_3, T_4\}$ ) and reference trajectory set (i.e.,  $R = \{T_1, T_2, T_3\}$ ) generated by SRR algorithm previously, range query searches for raw trajectories that belong to the specified spatio-temporal region.

The query processing is to extract qualitative information from trajectory databases that contain very large numbers of trajectories, whose efficiency depends crucially upon an appropriate index of trajectories. Due to the unique data characteristics of the compressed trajectories (e.g., containing sequence of both reference (sub-)trajectories and original samples of raw trajectories) and the unique query characteristics (e.g., often querying data in an instantaneous/continuous time window), we design a hybrid spatio-temporal index structure by enhancing R-tree [21] with inverted files. The index structure is presented in Section 6.1, followed by the detailed query processing in Section 6.2.

### 6.1 Hybrid Index

In this section, we introduce a hybrid index structure, called IR-tree [22], based on R-tree with each node enhanced with reference to an inverted file for the reference (sub-)trajectories contained in the sub-tree rooted at the node. Compared with the query processing over original trajectories (i.e., process on original data with R-tree), the query processing over compressed trajectories has several non-negligible advantages. First, since the volume of reference trajectories is much less (several orders of magnitude less) than the raw trajectories, there will be less dead space in the IR-tree, which in turn is more efficient for search. Moreover, the R-tree built using reference trajectories is small and compact enough to fully fit into memory to avoid physical I/O operations, so that searching on such an in-memory R-tree is comparatively faster. Our experimental study will also verify these advantages.

In our index structure, a leaf node  $N$  contains a number of entries in the form of  $(T.id, MBR, ifile)$ , in which  $T.id$  is an identifier that points to the reference (sub-)trajectories or the uncompressed part of compressed trajectory,  $MBR$  is a rectangle with two spatial dimensions, and  $ifile$  is a pointer to an inverted file for the reference (sub-)trajectories being indexed. Note that for the uncompressed segments of a compressed trajectory, we regard them as new reference

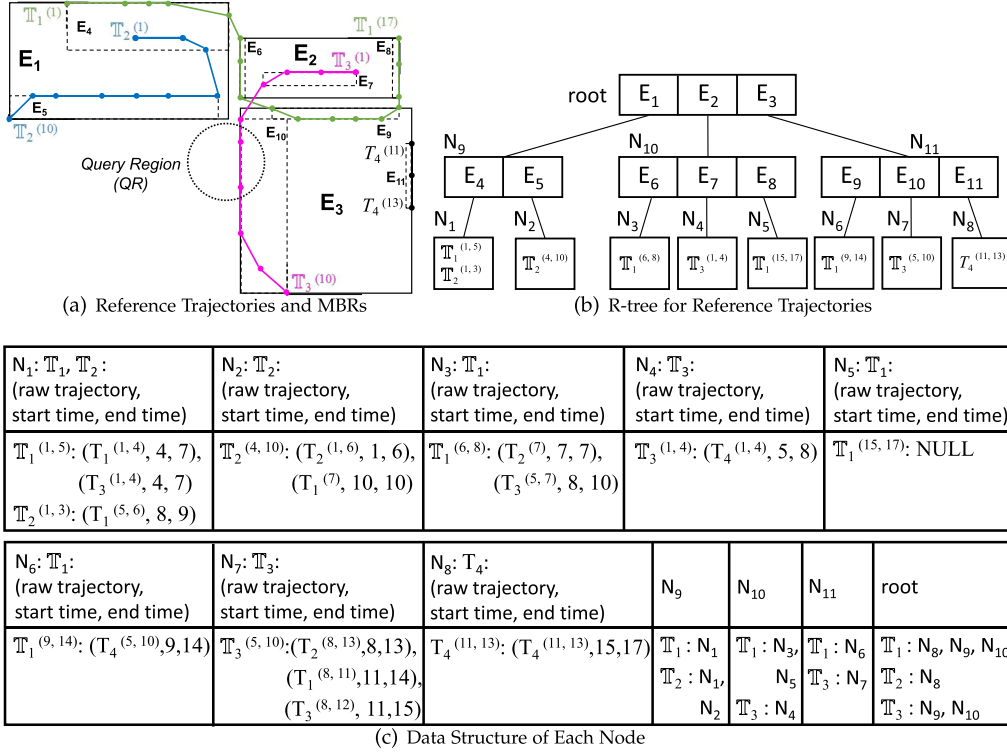


Fig. 5. Illustration of index structure.

trajectories and organize them in the index structure with the same way as reference trajectories. In this way, a single top-down traversal in this unified IR-tree can cover both compressed and uncompressed parts simultaneously. Taking  $T_4$  in Fig. 4 as an example,  $T_4$  can be represented by  $(T_3^{(1,4)}, T_1^{(9,14)}, T_4^{(11,13)})$ , where  $T_4^{(11,13)}$  is the uncompressed part of  $T_4$  that is regarded as a new reference trajectory when building the index. The inverted file consists of a hash table and a set of posting lists, wherein the hash table is used to index all the distinct reference (sub-)trajectories in  $N$  with a lookup time complexity of  $O(1)$ . For each reference (sub-)trajectory  $T^{(i,j)}$  in the hash table, we maintain a list of  $(T^{(m,n)}, T^{(m)}, t, T^{(n)}, t)$  triples for the raw trajectories relevant with  $T^{(i,j)}$ , which are sorted by their start timestamp (i.e.,  $T^{(m)}, t$ ). Note that the timestamps of the compressed part (consisting of a set of reference trajectories) in compressed trajectories may not be the actual timestamps. Therefore, before constructing the inverted lists, we recover the timestamps of the compressed part by scanning from the first point of each compressed sub-trajectory  $T'_{sub}$ . In particular, the timestamp of  $T'^{(i)}_{sub}$  will be replaced with the corrected timestamp when a correction record can be found; otherwise it will be set to  $T'^{(i-1)}_{sub}.t + T^{(i)}.t - T^{(i-1)}.t$ , where  $T$  is the reference trajectory used to compress  $T'_{sub}$ .

On the other hand, the entries in a non-leaf node  $N$  are of form  $(ptr, MBR, ifile)$ , in which  $ptr$  is the pointer to the child nodes of  $N$ ,  $MBR$  is the minimum bounding rectangle covering all the child nodes, and  $ifile$  denotes an inverted file that indexes all the distinct reference (sub-)trajectories within  $N$  in a hash table. The IR-tree can be constructed in the similar way of R-tree (i.e., insertion operation) with an exception that the inverted files have to be updated from the leaf to the root when a new reference (sub-)trajectory is added.'

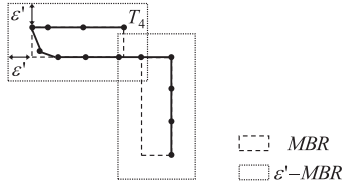
Fig. 5 presents the illustration of the whole hybrid index structure. Given several raw trajectories,  $\{T_1, T_2, T_3, T_4\}$ , which can be totally represented by a set of reference trajectories (see Fig. 4), we first need to obtain the *MBRs* of these reference trajectories, as shown in Fig. 5a. Subsequently, the R-tree is constructed based on the distribution of these reference (sub-)trajectories and *MBRs* in Fig. 5b. The root of the R-tree contains three data entries  $E_1, E_2$  and  $E_3$  referring to the children nodes  $N_8, N_9$  and  $N_{10}$  separately.  $N_8$  represents the minimum bounding rectangle of its children nodes  $N_1$  and  $N_2$ , and the information of the bounding region is contained in  $E_1$ , as with  $N_9$  and  $N_{10}$ . The spatial objects (i.e., the segments of the reference trajectories) are referred by the entries of the leaf nodes in R-tree. In Fig. 5c, the invert lists are given for each reference (sub-)trajectories that are used to compress the given raw trajectories.

## 6.2 Query Processing

After constructing the hybrid index, we next detail the query processing based on the index structure. Since it is expensive to calculate the spatial relationship (e.g., distance, containment, etc) between spatial objects (i.e., a query region and a trajectory in the range query), a query processing algorithm typically adopts a filter-and-refinement approach [23]. In particular, the filter step takes relatively cheap computation cost to find a small set of candidate trajectories that are likely to be the results, which is a super set of the result for the original query. Then these candidate trajectories are further processed using geometric algorithms (e.g., distance calculation) to obtain an actual result at the refinement step.

Subsequently, we discuss how to efficiently process range query based on the IR-tree. The range query aims to retrieve all the raw trajectories within a specified spatio-temporal




 Fig. 6. A query trajectory and its  $\epsilon'$ -MBRs.

window, e.g., a sphere of 10 km radius around a location (see the circular query region in Fig. 5a) in the time interval,  $\{8, 9, 10\}$ . In the filter step, IR-tree is traversed from the root and examines the query region against the *MBR* in each entry visited to check if they are relevant (i.e., contained/overlapped) with each other. The range query processing finally finds all the relevant leaf nodes. Considering the case in Fig. 5a and 5b,  $N_{10}$  is visited since its *MBR* overlaps with the Query Region (QR). Then for the same reason  $N_7$  is visited, and  $T_3^{(5,10)}$  is found from  $N_7$ . Finally, we can easily identify the candidate trajectory  $T_2$  in the specific time interval with the inverted lists in Fig. 5c, and examine whether the candidate actually overlaps with the query region in the refinement step.

For other kinds of queries like Trajectory-based query (T-query), we can first simply convert the queries into range query and then apply the filter-and-refinement approach. For instance, given a query trajectory (e.g.,  $T_4$  in Fig. 4) and a spatial threshold  $\epsilon'$ , a T-query aims to find trajectories that satisfy a given distance function to the query trajectory (e.g., the MaxDTW distance between the result trajectory and the query trajectory is less than  $\epsilon'$ ). Specifically, we can first construct several  $\epsilon'$ -MBRs, each of which extends the original *MBR*'s length and width by  $2\epsilon'$  respectively (as shown in Fig. 6). These  $\epsilon'$ -MBRs can be regarded as the query regions, and then we can process range query based on the IR-tree to identify the result trajectories.

## 7 EXPERIMENT

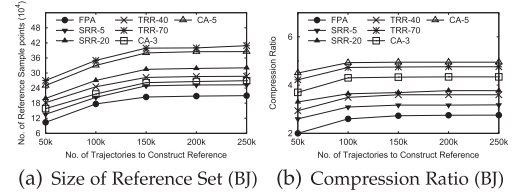
In this section, we conduct extensive experiments to validate the effectiveness of our proposed algorithms. All the algorithms are implemented on an Intel(R) Xeon(R) CPU E5-2630 v2 @ 2.60 GHZ with 256 GB RAM.

### 7.1 Experiment Setup

We use two real trajectory datasets generated by taxis in Beijing and Chengdu respectively over one week. The Beijing

TABLE 3  
Experiment Parameters

Parameters	Values
No. of trajectories to construct reference set $ D_R $	50k, <u>100k</u> , 150k, 200k, 250k
No. of trajectories to be compressed/queried $ D $	<u>200k</u> , 400k, 600k, 800k
Spatial deviation threshold $\epsilon_s$	<u>200m</u> , 400m, 600m, 800m, 1000m
Temporal deviation threshold $\epsilon_t$	30s, <u>60s</u> , 90s, 120s, 150s
Length of trajectory $ T $	50, 100, 150, 200, 250
Size of query window $Q_w$	2km, <u>4km</u> , 6km, 8km, 10km
Compression ratio $CR$	2, 4, 6, 8, 10


 Fig. 7. Performance of reference set construction: Effect of  $|D_R|$ .

dataset, denoted by BJ, contains about 2 million trajectories, and the Chengdu dataset, denoted by CD, contains about 1.4 million trajectories. In both datasets, trajectories in one day are used as training dataset to construct the reference set, and the rest are used to test the performance of different compression algorithms under various parameter settings with the average performance for one day recorded. The ranges and default values (underlined) of all the parameters are summarized in Table 3.

## 7.2 Experiment Results

### 7.2.1 Performance of Reference Set Construction

We first evaluate the performance of reference set construction and its impact to spatial compression on BJ dataset. Two metrics, *Size of Reference Set* (the number of sample points in the reference set) and *Compression Ratio* (CR), are compared in this set of experiments, wherein CR measures how much space has been saved by compressing raw trajectories. Specifically, CR is defined as the ratio between space costs of raw trajectories and compressed trajectories, i.e.,  $CR = \frac{space(T)}{space(T')}$ , where an original sample point contains its longitude, latitude and timestamp (costing 12 bytes), an MRT  $T^{(k,g)}$  used to represent the raw trajectories is recorded as a triplet  $(T.id, k, g)$  (costing 8 bytes), and a time correction costs 4 bytes. We compare these two metrics among the following methods (specified in Section 3) by varying  $|D_R|$ ,  $\epsilon_s$ .

- 1) FPA: FPA with minimum support 100.
- 2) SRR-5: SRR with minimum sub-trajectory length 5.
- 3) SRR-20: SRR with minimum sub-trajectory length 20.
- 4) TRR-40: TRR with overlap threshold 40 percent.
- 5) TRR-70: TRR with overlap threshold 70 percent.
- 6) CA-3: CA with compression ratio threshold 3.
- 7) CA-5: CA with compression ratio threshold 5.

*Effect of  $|D_R|$ .* In this set of experiment, we study the effect of  $|D_R|$ . As shown in Fig. 7a, naturally the sizes of reference sets generated from all approaches increase when more training trajectories are used. However, we also notice that the increase becomes slower when  $|D_R| > 150k$  since with more reference trajectories accumulated there is increasing chance that subsequently added trajectories are redundant. Among those competing methods, FPA generates the smallest reference set while TRR-70 results the largest followed by CA-5, SRR-20, TRR-40, CA-3 and SRR-5. It is found that the size of reference set almost stops growing beyond 400k, which only takes a few megabytes memory space. From Fig. 7b, the compression effectiveness heavily depends on the size of reference set since a larger reference set normally means greater coverage hence better compression ratio. Compression algorithm performs the worst on the reference set generated by FPA, which aims at capturing the major traveling patterns of training trajectories. Even though CA-5

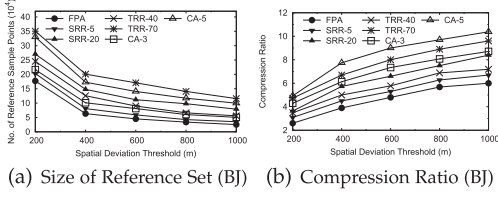


Fig. 8. Performance of Reference Set Construction: Effect of  $\epsilon_s$ .

generates a smaller reference set than TRR-70, it has the best compression ratio, which implies that the reference set generated by CA-5 is of high coverage and low redundancy. This is due to the fact that CA directly optimizes the compression ratio during the construction of reference set while SRR and TRR try to minimize the redundancy.

*Effect of  $\epsilon_s$ .* Next we study the effect of  $\epsilon_s$ . In Fig. 8a, the sizes of reference sets decrease with the increase of  $\epsilon_s$ , since  $\epsilon_s$  affects the granularity of patterns for FPA and the judgment of redundant trajectories for other approaches. Greater  $\epsilon_s$  means more trajectories become redundant, which in turn results in smaller reference set. On the compression ratio aspect, Fig. 8b demonstrates that as  $\epsilon_s$  increases the compression performance improves in spite of smaller reference set since a given trajectory has more chance to match fewer but longer reference trajectories.

### 7.2.2 Performance of Compression Algorithms

In this part we evaluate the effectiveness (*compression ratio*) and efficiency (*running time*) of the proposed compression algorithms, namely GSC, OSC, GSTC and OSTC, based on the same reference set generated by compress-based approach on both trajectory datasets (i.e., BJ and CD datasets). Moreover, we also implement two representative spatio-temporal compression algorithms, namely Normal Douglas-Peucker (NDP) algorithm based on SED [1] (a spatial trajectory simplification algorithm with synchronous euclidean distance) and PRESS [10] (a spatial-temporal trajectory compression algorithm with the constraint of road network), as our competitors. For NDP, we define  $\epsilon_s$  as the maximal allowed SED between a raw trajectory and its compressed trajectory. Since the spatial compression component of PRESS is lossless, we use  $\epsilon_s$  and  $\epsilon_t$  to represent the error metrics, TSND and NSTD, in the temporal compression component, respectively.

*Effect of  $\epsilon_s$ .* As expected, compression ratio of all algorithms gradually increases as  $\epsilon_s$  grows (see Fig. 10a and 10c).

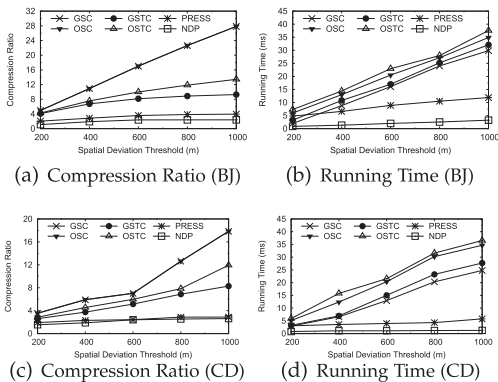


Fig. 9. Performance of compression algorithms: Effect of  $\epsilon_s$ .

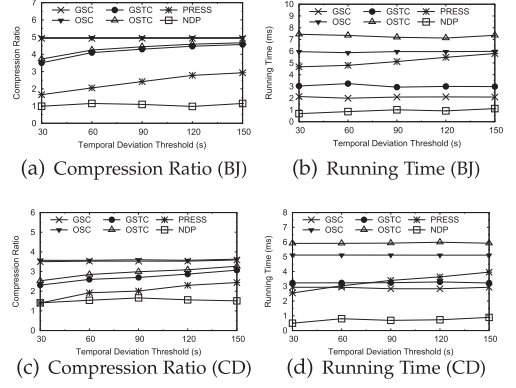
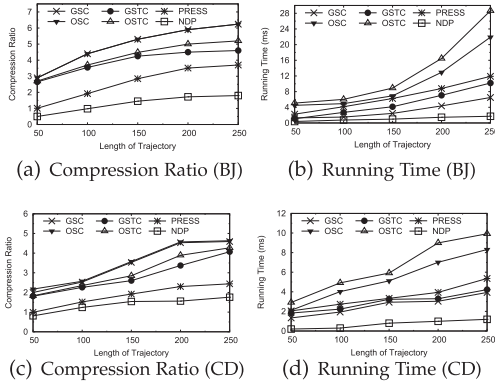


Fig. 10. Performance of compression algorithms: Effect of  $\epsilon_t$ .

Naturally, GSC and OSC achieve the best compression ratio, since they do not consider temporal information. The result also verifies our previous lemma that GSC and OSC have the same power in terms of compression ratio. Moreover, the compression ratio of OSTC and GSTC grows faster than that of PRESS and NDP, showing more benefits as  $\epsilon_s$  increases. OSTC achieves the best compression ratio amongst all spatio-temporal compression methods, confirming the optimality of our proposed algorithm. In terms of running time in Fig. 10b and 10d, NDP is fastest and almost not affected by  $\epsilon_s$ , while OSTC is most time-consuming. The running time of our proposed algorithms increase since a greater  $|\epsilon_s|$  results in more MRT enumerations during the compression. PRESS is also affected by  $|\epsilon_s|$  since it is related to the angular search region during bounded temporal compression. Moreover, GSTC (OSTC) runs slower than GSC (OSC) because of the extra time cost for obtaining the optimal MRT.

*Effect of  $\epsilon_t$ .* Obviously, as depicted by Fig. 10a and 10c, NDP, GSC and OSC are not affected by  $\epsilon_t$  since they do not consider temporal information at all. For GSTC and OSTC, a smaller  $\epsilon_t$  means more time stamps of the MRTs are likely to violate the temporal constraint, leading to more time correction cost, which explains the increasing trend of compression ratio as  $\epsilon_t$  grows. In addition, the compression ratios of OSTC and GSTC are very close and both outperform PRESS and NDP constantly by a large margin. When it comes to efficiency in Fig. 10b and 10d, none of the approaches except PRESS is affected by  $\epsilon_t$ , which is natural for GSC, OSC and NDP. As to GSTC and OSTC,  $\epsilon_t$  affects the number of time stamps to be rectified but the total number of time stamps to be checked remains the same. The efficiency of PRESS slightly decreases with  $\epsilon_t$  as its angular search region increases.

*Effect of  $|T|$ .* To study the effect of the length of trajectory, we select five groups of trajectories from the test dataset, each including 10000 trajectories with about 50, 100, 150, 200, 250 samples respectively and record the average compression ratio and running time for each group. In Fig. 11a and 11c, the compression ratios of all methods increase with  $|T|$ , as longer trajectories tend to have more redundant samples hence there are more room to improve the compression ratio. Regardless of  $|T|$ , proposed methods of REST framework well outperform their competitors constantly. As illustrated in Fig. 11b and 11d, the running time of all methods increase with the length of trajectory, while the growth of computational cost for OSTC and OSC is relatively faster due to the quadratic complexity with respect


 Fig. 11. Performance of compression algorithms: Effect of  $|T|$ .

to  $|T|$  when deriving the optimal storage cost based on dynamic programming.

**Reconstruction Error Analysis.** Furthermore, to evaluate the accuracy of trajectory compression, we adopt two measures, MaxDTW-based and Average Synchronous euclidean Distance [1] (ASED)-based spatial error, which are respectively quantified by the MaxDTW distance and average euclidean distance between a raw trajectory and its corresponding decompressed trajectory. The reason we adopt ASED is that the decompressed trajectory may have different length with its raw correspondence. Due to space limit, we only show the experimental results of BJ dataset. Since PRESS is spatially lossless, we only compare our compression algorithms and NDP by varying  $\epsilon_s$  and  $|T|$ . Fig. 12a shows that all our methods have the similar performances with respect to  $\epsilon_s$ . This is expected as these methods apply the same error metric (i.e., MaxDTW) during spatial compression. Besides, the spatial error of NDP is higher than that of our methods, demonstrating the superiority of our methods. In terms of ASED-based spatial error in Fig. 12b, NDP performs better than our methods because it adopts synchronous euclidean distance when compressing trajectories. From Fig. 12c and 12d we can see that the spatial errors of all approaches have a growing tendency as  $|T|$  is enlarged since a longer trajectory has more chance to deviate from its corresponding decompressed trajectory.

### 7.2.3 Performance of Query Processing

In the final set of experiments, we evaluate the effectiveness (*memory cost of index*) and efficiency (*running time*) of the following methods for processing range queries by varying the number of trajectories (to be queried)  $|D|$  on BJ dataset, the size of query window  $|Q_w|$  and compression ratio  $CR$ :

- 1) Range query processing on raw trajectories with spatial index is denoted as *SI*. This approach of query processing is creating spatial index (i.e., R-tree) on the original trajectories to speed up the query processing.

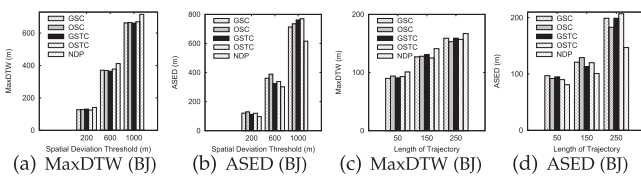
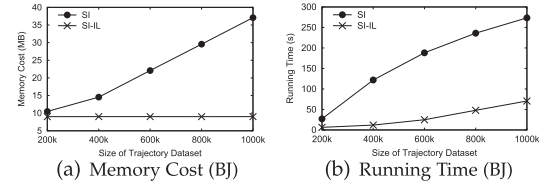


Fig. 12. Spatial error of compression algorithms.


 Fig. 13. Performance of query processing: Effect of  $|D|$ .

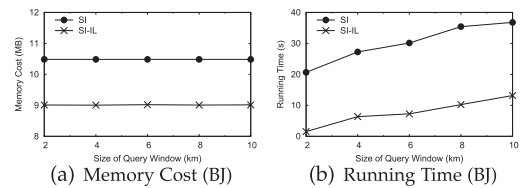
- 2) Range query processing on partially decompressed trajectory data is denoted as *SI-IL*. This method takes IR-tree as the spatial index for compressed trajectories.

In order to compare the performance of the above two approaches for range query processing, we randomly generate  $10^3$  queries on test dataset (i.e., the trajectory set to be queried) and record the total running time. These queries are process in the raw trajectories and compressed trajectories (generated by OSTC) respectively.

**Effect of  $|D|$ .** To study the scalability of the algorithms, we generate 5 datasets containing from 200  $k$  to 1000  $k$  trajectories by random selection from the original trajectory set. Then we run experiments on these five datasets, the results of which are depicted in Fig. 13. As shown in Fig. 13a, the memory cost of SI increases as the size of trajectory dataset is enlarged since its index is mainly established on these trajectories, while SI-IL keeps a constant memory cost when enlarging  $|D|$  due to the fact that its index is mainly established with the reference set whose size is stable. When referring to the running time in Fig. 13b, as expected, though the running time increase as the number of trajectories increases, our proposed algorithm scales well with the size of the dataset generally.

**Effect of  $|Q_w|$ .** Subsequently, we study the effects of the size of query window  $|Q_w|$ , the radius of circular query window, by changing it from 2 km to 10 km. Not surprisingly, as we can see from Fig. 14, SI-IL significantly outperforms SI for all values of  $|Q_w|$  in terms of both memory cost and runtime. For memory cost, SI and SI-IL keep a constant cost as  $|Q_w|$  increases since the number of trajectories used to construct the index is stable, which is demonstrated in Fig. 14a. For efficiency, the performances of both algorithms deteriorate with the increasing size of query window (see Fig. 14b). This is because statistically the relevance of a trajectory will be affected by more points as  $|Q_w|$  increases. In other words, more nodes of the R-tree overlap with the larger query window and need to be examined. Besides, we also notice that the running time of SI-IL is far less than that of SI, which testifies the superiority of SI-IL.

**Effect of  $CR$ .** In the final experiments, we investigate how the compression ratio  $CR$  affects the performance of our proposed SI-IL method. We choose 5 groups of trajectories whose compression ratios are approximately equal to 2, 4, 6, 8, 10 respectively. Generally, a large  $CR$  means that most


 Fig. 14. Performance of query processing: Effect of  $|Q_w|$ .



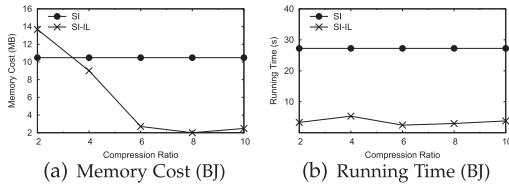


Fig. 15. Performance of query processing: Effect of  $CR$ .

part of the raw trajectory can be compressed by the reference (sub-)trajectories while a small  $CR$  means more original samples of the trajectory are kept for representing it when being compressed. From Fig. 15a we can see that, SI-IL approach performs badly when  $CR$  is small (i.e.,  $CR < 3.5$ ), which consumes more memory than SI, since SI-IL has to construct a complicated index structure for a large number of original samples of the trajectories that are used to represent the trajectories. However, when  $CR \geq 6$ , the memory cost remains relatively small and unchanged. This may be due to the fact that the raw trajectories can be almost compressed by the reference (sub-)trajectories and the index can be only constructed on the reference (sub-)trajectories. Besides, Fig. 15b illustrates SI-IL is consistently more efficient than SI, showing the benefits of our proposed method.

## 8 RELATED WORK

### 8.1 Trajectory Compression

The existing trajectory compression algorithms can be classified into two categories: 1) spatial compression; 2) spatio-temporal compression.

Spatial compression algorithms treat trajectories as polylines. For example, Douglas-Peucker algorithm [2] is a classic line generalization approach that uses a perpendicular distance threshold to reduce the number of points. As DP algorithm is simple and feasible, a variety of applications have been proposed [24] to speed up DP algorithm. Bellman's algorithm [3] fits a finite number of line segments to a curve based on dynamic programming, preserving the most essential spatial features. Due to the pure geometric nature, the above algorithms cannot be applied when the temporal information of trajectories also matters.

Taking the temporal component of trajectories into account, Sliding Window Algorithm [6] and Opening Window Algorithm [1], [25] are designed to keep spatio-temporal information of a trajectory within a sliding window to compress it. Muckell et al. introduce a heuristic method called SQUISH [7], using a priority queue where the priority of each point is defined as an estimate of the error that the removal of that point would introduce, to compress trajectories. Recently, a lossless path compressor in road network is developed in [26], namely Minimum Entropy Labeling, which guarantees a theoretic bound and achieves practically high spatio-temporal compressibility. [27] presents a compressed data structure for moving object trajectories, where the compressed trajectories can be represented as sequences of road edges. Song et al. develop a framework, PRESS [10], that separates a given trajectory in a road network into spatial path and time sequence components. These two components are then compressed by Hybrid Spatial Compression algorithm and error Bounded Temporal Compression algorithm

respectively, achieving spatial lossless and temporal error-bounded compression. However, the pre-computation and storage of all-pair shortest paths and most frequent paths require a stable road network and large memory space to be available, which limits its applied scenarios.

### 8.2 Trajectory Index and Query

As a fundamental trajectory data manipulation, trajectory query is becoming crucial. Building efficient spatio-temporal index structures on trajectories can significantly facilitate query processing in trajectory databases. Koide et al. propose a novel spatio-temporal index structure for Network-Constrained Trajectories (NCTs), namely Suffix-array-based Network-constrained Trajectory index [28]. Then they further design a compressed-index method for NCTs by converting NCTs into a trajectory string, in which Relative Movement Labeling and FM-index are incorporated to accelerate the query processing [27]. However, the above existing work targets network-constrained data and uses this characteristic to improve the efficiency of query processing and the accuracy of the query results. Effective index structures [19], [29], [30], [31] are built to manage trajectories in euclidean space and support high performance trajectory queries, among which R-tree [21] is the most common index adopted to accelerate the query processing. Here, we also process the trajectory query over an R-tree index of all the reference (sub-)trajectories. Notice that we do not consider the STR-tree or TB-tree [29] for trajectory index since they focus more on trajectory preservation and leave other spatial properties like spatial proximity aside, while in the our problem, the R-tree index is only for fast retrieval of nearest trajectories.

For trajectory query, a typical one asks for the information against spatio-temporal relationships between trajectories and other spatial data objects, i.e., points, regions and trajectories [23]. For instance, Chen et al. propose  $k$ -Path Nearest Neighbor ( $k$ -PNN) query to return the  $k$ -NN with respect to shortest path connecting a given point (i.e., a destination); SETI [32] and PA-tree [33] are designed for range queries which seek to find all trajectories that intersect a spatial region; and SECONDO [34] and TrajStore [35] are developed for performing regular  $k$ -NN queries on trajectories using the euclidean distance.

## 9 CONCLUSION AND FUTURE WORK

In this paper we propose a novel data-driven framework, called REST, to compress the spatio-temporal information of trajectory data. In order to achieve high effectiveness and efficiency, we addressed a few challenges by proposing different strategies to construct a compact but expressive reference set, and designing efficient and optimal algorithms to represent a given trajectory with selected matchable reference trajectories. To the best of our knowledge, it is the first data-driven approach to compress trajectories in unconstrained space with both spatial and temporal dimensions considered. In addition, as the compressed trajectories are in the form of sequence of reference trajectories and original samples, we develop an effective index structure to support efficient query processing over compressed trajectories without full decompression. Extensive empirical study based on real trajectories dataset also confirms the superiority of our

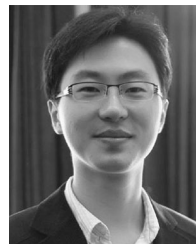
proposed framework over the state-of-the-art approaches in terms of compression ratio, efficiency and space cost.

## ACKNOWLEDGMENTS

The work is supported by the National Natural Science Foundation of China (Grant No. 61532018, 61836007, 61832017, 61772356 and 61872258). Kai Zheng and Yan Zhao contributed equally to this work.

## REFERENCES

- [1] N. Meratnia and A. Rolf, "Spatiotemporal compression techniques for moving point objects," in *Proc. Int. Conf. Extending Database Technol.*, 2004, pp. 765–782.
- [2] D. H. Douglas and T. K. Peucker, "Algorithms for the reduction of the number of points required to represent a digitized line or its caricature," *Int. J. Geographic Inf. Geovisualization*, vol. 10, no. 2, pp. 112–122, 1973.
- [3] R. Bellman, "On the approximation of curves by line segments using dynamic programming," *Archives Internal Med.*, vol. 4, no. 6, 1961, Art. no. 284.
- [4] R. Lange, F. Dürr, and K. Rothermel, "Efficient real-time trajectory tracking," *VLDB J.*, vol. 20, no. 5, pp. 671–694, 2011.
- [5] J. Liu, K. Zhao, P. Sommer, S. Shang, B. Kusy, and R. Jurdak, "Bounded quadrant system: Error-bounded trajectory compression on the go," in *Proc. Int. Conf. Data Eng.*, 2015, pp. 987–998.
- [6] E. Keogh, S. Chu, D. Hart, and M. Pazzani, "An online algorithm for segmenting time series," in *Proc. IEEE Int. Conf. Data Mining*, 2001, pp. 289–296.
- [7] J. Muckell, J.-H. Hwang, V. Patil, C. T. Lawson, F. Ping, and S. Ravi, "Squish: An online approach for gps trajectory compression," in *Proc. 2nd Int. Conf. Comput. Geospatial Res. Appl.*, 2011, Art. no. 13.
- [8] M. Potamias, K. Patroumpas, and T. Sellis, "Sampling trajectory streams with spatiotemporal criteria," in *Proc. 18th Int. Conf. Sci. Statistical Database Manage.*, 2006, pp. 275–284.
- [9] G. Trajcevski, H. Cao, P. Scheuermann, O. Wolfson, and D. Vaccaro, "On-line data reduction and the quality of history in moving objects databases," in *Proc. 5th ACM Int. Workshop Data Eng. Wireless Mobile Access*, 2006, pp. 19–26.
- [10] R. Song, W. Sun, B. Zheng, and Y. Zheng, "Press: A novel framework of trajectory compression in road networks," *VLDB Endowment*, vol. 7, no. 9, pp. 661–672, 2014.
- [11] M. C. González, C. A. Hidalgo, and A.-L. Barabási, "Understanding individual human mobility patterns," *Nature*, vol. 453, no. 7196, pp. 779–782, 2008.
- [12] K. Zheng, Y. Zheng, X. Xie, and X. Zhou, "Reducing uncertainty of low-sampling-rate trajectories," in *Proc. IEEE 28th Int. Conf. Data Eng.*, 2012, pp. 1144–1155.
- [13] Y. Zhao, S. Shang, Y. Wang, B. Zheng, Q. V. H. Nguyen, and K. Zheng, "Rest: A reference-based framework for spatio-temporal trajectory compression," in *Proc. SIGKDD*, 2018, pp. 2797–2806.
- [14] D. J. Berndt and J. Clifford, "Using dynamic time warping to find patterns in time series," in *Proc. 3rd Int. Conf. Knowl. Discovery Data Mining*, 1994, pp. 359–370.
- [15] J. Li, J. Wang, L. Yu, and J. Zhang, "A novel frequent trajectory mining method based on gsp," in *Proc. Int. Conf. Web Inf. Syst. Mining*, 2011, pp. 134–140.
- [16] R. Agrawal and R. Srikant, "Mining sequential patterns," in *Proc. 11th Int. Conf. Data Eng.*, 1995, pp. 3–14.
- [17] S. Ranu, D. P. A. D. Telang, P. Deshpande, and S. Raghavan, "Indexing and matching trajectories under inconsistent sampling rates," in *Proc. IEEE 31st Int. Conf. Data Eng.*, 2015, pp. 999–1010.
- [18] M. Vlachos, D. Gunopoulou, and G. Kollios, "Discovering similar multidimensional trajectories," in *Proc. 18th Int. Conf. Data Eng.*, 2002, pp. 673–684.
- [19] C. Lei and R. Ng, "On the marriage of lp-norms and edit distance," in *Proc. 30th Int. Conf. Very Large Data Bases - Vol. 30*, 2004, pp. 792–803.
- [20] C. Lei, M. T. Ozsu, and V. Oria, "Robust and fast similarity search for moving object trajectories," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2005, pp. 491–502.
- [21] A. Guttman, "R-trees: A dynamic index structure for spatial searching," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 1984, pp. 47–57.
- [22] Z. Li, K. C. Lee, B. Zheng, W. Lee, D. L. Lee, and X. Wang, "Ir-tree: An efficient index for geographic document search," *IEEE Trans. Knowl. Data Eng.*, vol. 23, no. 4, pp. 585–599, Apr. 2011.
- [23] K. Deng, K. Xie, K. Zheng, and X. Zhou, "Trajectory indexing and retrieval," in *Computing with Spatial Trajectories*. New York, NY, USA: Springer, 2011.
- [24] R. B. McMaster, "A statistical analysis of mathematical measures for linear simplification," *Amer. Cartographer*, vol. 13, no. 2, pp. 103–116, 1986.
- [25] J. Muckell, J. H. Hwang, C. T. Lawson, and S. S. Ravi, "Algorithms for compressing gps trajectory data: An empirical evaluation," in *Proc. 18th SIGSPATIAL Int. Conf. Adv. Geographic Inf. Syst.*, 2010, pp. 402–405.
- [26] Y. Han, W. Sun, and B. Zheng, "Compress: A comprehensive framework of trajectory compression in road networks," *ACM Trans. Database Syst.*, vol. 42, no. 2, pp. 1–49, 2017.
- [27] S. Koide, Y. Tadokoro, C. Xiao, and Y. Ishikawa, "Cinct: Compression and retrieval for massive vehicular trajectories via relative movement labeling," in *Proc. IEEE 34th Int. Conf. Data Eng.*, 2018, pp. 1097–1108.
- [28] S. Koide, Y. Tadokoro, and T. Yoshimura, "Snt-index: Spatio-temporal index for vehicular trajectories on a road network based on substring matching," in *Proc. 1st Int. ACM SIGSPATIAL Workshop Smart Cities Urban Anal.*, 2015, pp. 1–8.
- [29] D. Pfoser, C. S. Jensen, and Y. Theodoridis, "Novel approaches to the indexing of moving object trajectories," *Proc. 26th Int. Conf. Very Large Data Bases*, 2000, pp. 395–406.
- [30] B. Zheng, H. Wang, K. Zheng, H. Su, K. Liu, and S. Shang, "Sharkdb: An in-memory column-oriented storage for trajectory analysis," *WWW J.*, vol. 21, no. 2, pp. 1–31, 2017.
- [31] K. Zheng, S. Shang, N. J. Yuan, and Y. Yang, "Towards efficient search for activity trajectories," in *Proc. IEEE 29th Int. Conf. Data Eng.*, 2013, pp. 230–241.
- [32] V. P. Chakka, A. Everspaugh, and J. M. Patel, "Indexing large trajectory data sets with seti," in *Proc. CIDR*, 2003, pp. 169–180.
- [33] J. Ni and C. V. Ravishanker, "Pa-tree: A parametric indexing scheme for spatio-temporal trajectories," in *Proc. Int. Symp. Spatial Temporal Databases*, 2005, pp. 254–272.
- [34] R. H. Gutting, A. Braese, T. Behr, and J. Xu, "Nearest neighbor search on moving object trajectories in second," in *Proc. Int. Symp. Spatial Temporal Databases*, 2009, pp. 427–431.
- [35] P. Cudre-Mauroux, E. Wu, and S. R. Madden, "Trajstore: An adaptive storage system for very large trajectory data sets," in *Proc. IEEE 26th Int. Conf. Data Eng.*, 2010, pp. 109–120.



**Kai Zheng** received the PhD degree in computer science from The University of Queensland, in 2012. He is a professor of computer science with the University of Electronic Science and Technology of China. He has been working in the area of spatial-temporal databases, uncertain databases, social-media analysis, in-memory computing and blockchain technologies. He has published more than 100 papers in prestigious journals and conferences in data management field such as SIGMOD, ICDE, the *VLDB Journal*, ACM Transactions and IEEE Transactions. He is a member of the IEEE.



**Yan Zhao** received the master's degree in geographic information system from the University of Chinese Academy of Sciences, in 2015. She is currently working toward the PhD degree at Soochow University. Her research interests include spatial database and trajectory computing.



He also was a reviewer for several journals such as the the *IEEE Transactions on Knowledge and Data Engineering*, the *ACM Transactions on Information Systems* and the *Knowledge and Information Systems*.

**Defu Lian** received the BE and PhD degrees in computer science from the University of Science and Technology of China (USTC), in 2009 and 2014, respectively. He is currently a research professor with the University of Science and Technology of China. His main research interests include mobile data mining and recommender systems. He has published more than 40 papers in journals and conferences such as KDD, WWW, IJCAI, and the *ACM Transactions on Information Systems*, the *IEEE Transactions on Knowledge and Data Engineering*.



**Bolong Zheng** received the PhD degree from the University of Queensland, in 2017. He is currently working as an associate professor with Huazhong University of Science and Technology. After that, he worked as a postdoctoral research fellow with the University of Queensland and Aalborg University. His research interests include trajectory querying and mining, social-media analysis, spatial-temporal databases.



**Guanfeng Liu** received the PhD degree in computer science from Macquarie University, Australia, in 2013. He is current a lecturer with the Department of Computing, Macquarie University, Australia. His research interests include graph data management, trust computing and social networks. He has published more than 60 papers in the most prestigious journals and conferences such as IJCAI, AAAI, ICDE, CIKM, the *IEEE Transactions on Knowledge and Data Engineering*, the *IEEE Transactions on Services Computing* and ICWS.



**Xiaofang Zhou** received the bachelor's and master's degrees in computer science from Nanjing University, in 1984 and 1987, respectively, and the PhD degree in computer science from the University of Queensland, in 1994. He is a professor of computer science with the University of Queensland. He is the head of the Data and Knowledge Engineering Research Division, School of Information Technology and Electrical Engineering. He is also a specially appointed an adjunct professor with Soochow University, China. His research is focused on finding effective and efficient solutions to managing integrating, and analyzing very large amounts of complex data for business and scientific applications. His research interests include spatial and multimedia databases, high performance query processing, web information systems, data mining, and data quality management. He is a fellow of the IEEE.

▷ **For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/csdl](http://www.computer.org/csdl).**