



Efficient Non-Learning Similar Subtrajectory Search

Jiabao Jin
East China Normal University
Shanghai, China
jiabaojin@stu.ecnu.edu.cn

Peng Cheng
East China Normal University
Shanghai, China
pcheng@sei.ecnu.edu.cn

Lei Chen
HKUST-GZ and HKUST
Guangzhou and Hong Kong, China
leichen@cse.ust.hk

Xuemin Lin
Shanghai Jiaotong University
Shanghai, China
xuemin.lin@gmail.com

Wenjie Zhang
University of New South Wales
Sydney, Australia
wenjie.zhang@unsw.edu.au

ABSTRACT

Similar subtrajectory search is a finer-grained operator that can better capture the similarities between one query trajectory and a portion of a data trajectory than the traditional similar trajectory search, which requires that the two checking trajectories are similar in their entirety. Many real applications (e.g., trajectory clustering and trajectory join) utilize similar subtrajectory search as a basic operator. It is considered that the time complexity is $O(mn^2)$ for exact algorithms to solve the similar subtrajectory search problem under most trajectory distance functions in the existing studies, where m is the length of the query trajectory and n is the length of the data trajectory. In this paper, to the best of our knowledge, we are the first to propose an exact algorithm to solve the similar subtrajectory search problem in $O(mn)$ time for most of widely used trajectory distance functions (e.g., WED, DTW, ERP, EDR and Frechet distance). Through extensive experiments on three real datasets, we demonstrate the efficiency and effectiveness of our proposed algorithms.

PVLDB Reference Format:

Jiabao Jin, Peng Cheng, Lei Chen, Xuemin Lin, and Wenjie Zhang. Efficient Non-Learning Similar Subtrajectory Search. PVLDB, 16(11): 3111 - 3123, 2023.
doi:10.14778/3611479.3611512

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/inabao/trajcSimilar>.

1 INTRODUCTION

The increasing popularity of mobile devices flourishes the generation of trajectory data, which is widely used in many fields (e.g., traffic flow prediction [9, 10], route planning [23]). With the focus of researchers on trajectory data, more and more methods are proposed for analyzing and processing trajectory data.

A significant problem in analyzing trajectory data is to query the most similar trajectory to a given trajectory among the vast amount of trajectories in the database [4, 5, 15, 19, 29, 30]. In real scenarios,

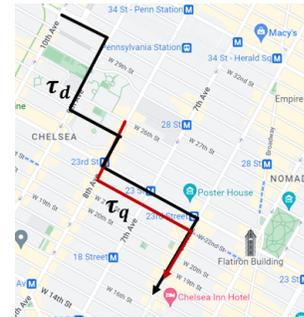


Figure 1: Subtrajectory Search

it is hard to guarantee that the lengths of two trajectories are same or close to each other. Thus, similar subtrajectory search attracts much attention recently as a more practical method [1, 3, 14, 21, 26], which uses a part of a long data trajectory as the basic unit to test its similarity to the short query trajectory. For example, as shown in Figure 1, there are two trajectories: data trajectory τ_d and query trajectory τ_q . They are not similar when the whole trajectories are considered, while τ_q is similar to a portion of τ_d .

Searching similar subtrajectories is usually a basic operator in real applications (e.g., subtrajectory join [21] and subtrajectory clustering [1, 3]) and will be frequently invoked, thus its efficiency is very important. One application scenario of subtrajectory query is to analyze the performance of players by their trajectory data in a sport (e.g., soccer or basketball) [26].

Subtrajectory search is a highly related but different problem from trajectory search [7, 8, 13, 20, 27]. Compared with trajectory search, subtrajectory search has to not only consider the data trajectory itself but also determine whether there are subtrajectories of the data trajectory with a smaller distance from the query trajectory. The state-of-the-art study on similar subtrajectory search utilize reinforcement learning methods to accelerate the detecting speed and achieve the time complexity of $O(mn)$ [27], where m is the length of the query trajectory and n is the length of the data trajectory. However, the reinforcement learning based algorithms are approximation algorithms, which have no theoretical guarantee on the accuracy of the returned results. In this paper, we find that *the similar subtrajectory search problem can be solved exactly with the time complexity of $O(mn)$ for most trajectory distance functions* (e.g., DTW, WED, ERP, EDR and FD. Details will be discussed in Section 5), which had not been discovered to the best of our knowledge.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.
Proceedings of the VLDB Endowment, Vol. 16, No. 11 ISSN 2150-8097.
doi:10.14778/3611479.3611512

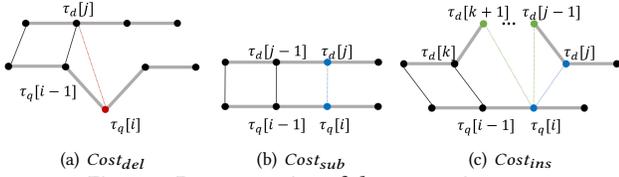


Figure 2: Demonstration of the conversion cost

Challenges. For a data trajectory, the number of its subtrajectories is quadratic to its length. Let n be the length of a data trajectory, there will be $\frac{n(n+1)}{2}$ its subtrajectories. Assuming that the length of the query trajectory is m and the length of the data trajectory is n , the time complexity of directly searching for the most similar subtrajectory is $O(mn^3)$ (through traversal searching $\frac{n(n+1)}{2}$ subtrajectories of the data trajectory, and the time complexity of directly computing the similarity of two trajectories of length x and y by dynamic programming is $O(xy)$). Although a recent work [27] optimizes the time complexity of a single subtrajectory query problem from $O(mn^3)$ to $O(mn^2)$ through dynamic programming techniques, it is still unaffordable for most applications that need to find the optimal subtrajectory in a few seconds. In existing studies, only for *dynamic time wrapping distance* (DTW) and *Frechet distance* (FD), the similar subtrajectory search problem can be exactly solved in $O(nm)$ time complexity with particular algorithms [8, 20]. However, it cannot be extended to other trajectory distance functions.

In this paper, we propose the conversion-matching algorithm (CMA) to find the optimal subtrajectory by computing the minimum cost of converting the query trajectory into the data trajectory. With carefully tailored methods and transformation of the trajectory distance functions, we can incrementally fast track the optimal start position of the optimal subtrajectory in the data trajectory in $O(1)$ time. Given a query trajectory and a data trajectory, we search for the optimal subtrajectory with the time complexity of $O(nm)$. Meanwhile, the algorithm is applicable for the vast majority of distance functions. We use *weighted edit distance* (WED) [13] and *dynamic time warping* (DTW) [30] as examples to analyze the design of the algorithm. We also discuss how to apply our methods to other most popular trajectory distance functions. Experiments show that the performance of our algorithm is better than other existing methods.

To summarize, we make the following contributions:

- We propose CMA with the time complexity of $O(nm)$ to find the most similar subtrajectory for a query trajectory under most order-insensitive trajectory distance functions in Section 4.
- We describe the design idea of the algorithm in detail and simplify the calculation of conversion cost, using WED and DTW as examples in Section 5.
- We conduct experiments on three different real data sets to verify the superiority of our framework with the state-of-the-art similar subtrajectory query methods in Section 6.

2 PROBLEM DEFINITION

2.1 Basic Concepts

There are two types of trajectories for the Similar Subtrajectory Search (SSS) problem: query and data trajectories. We expect to search for the most similar subtrajectory for a given query trajectory

under a specific distance function among a large volume of data trajectories. We first provide the definitions of trajectories and subtrajectories as follows:

Definition 1. (Trajectory) A trajectory τ with the length of n consists of a series of points denoted as $\langle p_1, p_2, p_3, \dots, p_n \rangle$.

We denote the query trajectory as τ_q with the length of m and the data trajectory as τ_d with the length of n . The points of trajectories can be specific physical locations, nodes on a road network or edges on a road network. In particular, we denote a trajectory without any point as τ_\emptyset .

Definition 2. (Subtrajectory) Given a trajectory τ with the length of n , its subtrajectory is a portion of consecutive points, $\tau[i : j] = \langle p_i, p_{i+1}, \dots, p_j \rangle$ ($1 \leq i \leq j \leq n$).

In particular, we denote the i^{th} point in τ as $\tau[i : i]$, abbreviated as $\tau[i]$. If $i > j$, we have $\tau[i : j] = \tau_\emptyset$.

Usually, we have a set of data trajectories. In this paper, we focus on finding the optimal subtrajectory from a data trajectory among many data trajectories to match the query trajectory. We have also implemented two pruning methods, Grid-Based Prune (GBP) and Key Points Filter (KPF), to help filter the irrelevant trajectories quickly. Please refer to our technical report for more details [11].

2.2 Distance Function

The distance function between trajectories represents the cost of converting the points of the query trajectory into the data trajectory plus the cost of inserting prefix subtrajectory and suffix subtrajectory.

Definition 3 (Matching Sequence). For a query trajectory τ_q and a data trajectory τ_d , we define its matching sequence as $\mathcal{A}_{\tau_q, \tau_d} = [a_1, a_2, a_3, \dots, a_m]$. If $\tau_q[i]$'s matching point is $\tau_d[j]$, we let $a_i = j$ to indicate the index of $\tau_d[j]$ in the data trajectory. For any $i \leq j$, we must have $a_i \leq a_j$.

According to the definition, if a trajectory $\tau_d[s : t]$ is a subtrajectory of another $\tau_d[i : j]$, we have $\mathcal{A}_{\tau_q, \tau_d[s:t]} \subseteq \mathcal{A}_{\tau_q, \tau_d[i:j]}$. For example, the matching sequence for Figure 4(a) is [1, 1, 2, 4, 5, 6, 7, 8, 9], and the matching sequence for Figure 4(b) is [1, 1, 2, 2, 3, 3, 5, 6, 9]. Note that, given a data trajectory τ_d and a query trajectory τ_q , there may be many matching sequences. One matching sequence is valid as long as its matching index value is not decreasing (i.e., $a_i \leq a_j, \forall i \leq j$).

Definition 4 (Point Matching-Conversion Cost). For a data trajectory τ_d and a query trajectory τ_q , when $\tau_q[i]$ matches $\tau_d[j]$ (i.e., $a_i = j$), depending on the different matches of $\tau_q[i-1]$ as shown in Figure 2, we define the cost of converting $\tau_q[i]$ into $\tau_d[j]$ in the following three cases:

- $Cost_{del}$. When $a_{i-1} = j$, we need to remove $\tau_q[i]$, and denote the conversion cost as $Cost(\tau_q[i], \tau_d[a_i]) = Cost_{del}(\tau_q[i], \tau_d[j])$.
- $Cost_{sub}$. When $a_{i-1} = j-1$, we replace $\tau_q[i]$ with $\tau_d[j]$, and denote the conversion cost as $Cost(\tau_q[i], \tau_d[a_i]) = Cost_{sub}(\tau_q[i], \tau_d[j])$.
- $Cost_{ins}$. When $a_{i-1} = k$ ($1 \leq k < j-1$), we substitute $\tau_q[i]$ with $\tau_d[j]$ and insert $\tau_d[k+1 : j-1]$. We denote the conversion cost as $Cost(\tau_q[i], \tau_d[a_i]) = Cost_{ins(k)}(\tau_q[i], \tau_d[j])$.

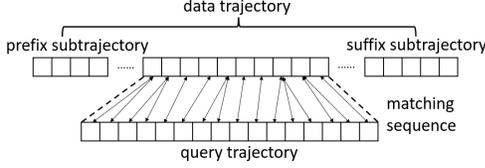


Figure 3: Demonstration of Matching Process

We can calculate the cost of converting the query trajectory into a data trajectory for each matching sequence, which includes the cost of converting each point in the query trajectory into a matching point in the data trajectory and the cost of inserting *prefix trajectory* and *suffix trajectory* of the data trajectory as shown in Figure 3. We denote the cost of inserting a subtrajectory $\tau_d[x : y]$ as $Insert(\tau_d[x : y])$. Given a matching sequence $\mathcal{A}_{\tau_q, \tau_d}$, its *matching-conversion cost* is $\sum_{a_i \in \mathcal{A}_{\tau_q, \tau_d}} Cost(\tau_q[i], \tau_d[a_i]) + Insert(\tau_d[1 : a_1 - 1]) + Insert(\tau_d[a_m + 1 : n])$. We give an example to demonstrate the calculation of the matching-conversion cost in Example 2.

Definition 5 (General Distance Function). We denote the set of all possible matching sequences between the query trajectory τ_q and the data trajectory τ_d as \mathbb{A} . Then, we define the general distance $\Theta(\tau_q, \tau_d)$ between the query trajectory and the data trajectory as follows:

$$\Theta(\tau_q, \tau_d) = \min_{\mathcal{A}_{\tau_q, \tau_d} \in \mathbb{A}} \sum_{a_i \in \mathcal{A}_{\tau_q, \tau_d}} Cost(\tau_q[i], \tau_d[a_i]) + Insert(\tau_d[1 : a_1 - 1]) + Insert(\tau_d[a_m + 1 : n]) \quad (1)$$

There are many trajectory distance functions, such as DTW [30], ERP [4], EDR [5], and WED [13]. In this paper, we use WED and DTW as examples to illustrate our definition. We discuss the generality of the general distance $\Theta(\tau_q, \tau_d)$ in the Appendix A of our technical report [11].

WED. WED is a general distance function that allows the user-defined cost functions and contains several important cost functions (e.g., EDR and ERP). WED defines the distance $wed(\tau_q, \tau_d)$ between τ_q and τ_d as the minimum cost of converting τ_q to τ_d by a finite number of insertion, deletion and substitution. Given two points $\tau_q[i]$ and $\tau_d[j]$, we denote the cost of insertion, deletion and substitution by $ins(\tau_d[j])$, $del(\tau_q[i])$ and $sub(\tau_q[i], \tau_d[j])$. Besides, the cost of deleting the subtrajectory $\tau_q[i : j]$ and inserting the subtrajectory $\tau_d[i : j]$ are denoted as $del(\tau_q[i : j])$ and $ins(\tau_d[i : j])$. We have $del(\tau_q[i : j]) = \sum_{i \leq k \leq j} del(\tau_q[k])$ and $ins(\tau_d[i : j]) = \sum_{i \leq k \leq j} ins(\tau_d[k])$.

Example 1. Given two trajectories τ_q and τ_d as shown in Figure 4, we use WED to calculate the distance between them. The blue points indicate the deleted points in the conversion of τ_q into τ_d , while the green points indicate the inserted points. We set the cost of $ins(\tau_d[j])$, $del(\tau_q[i])$ to 1. In addition, we set the cost of $sub(\tau_q[i], \tau_d[j])$ to 1 if $\tau_q[i] \neq \tau_d[j]$; otherwise, it is set to 0. Figure 4(a) shows an optimal matching sequence that converts τ_q into τ_d by deleting $\tau_q[2]$, inserting $\tau_d[3]$ and substituting $\tau_q[5]$ with $\tau_d[5]$ and $\tau_q[8]$ with $\tau_d[8]$. Since there are no redundant prefix and suffix subtrajectories, we have $Insert(\tau_d[1 : a_1 - 1]) + Insert(\tau_d[a_m + 1 : n]) = 0$. Therefore, the distance between τ_q and τ_d is $4 (= del(\tau_q[2]) + ins(\tau_d[3]) + sub(\tau_q[5], \tau_d[5]) + sub(\tau_q[8], \tau_d[8]))$.

Moreover, we can compute the distance between τ_q and τ_d by a dynamic programming algorithm [12]. We have $wed(\tau_q[i : j], \tau_\emptyset) = del(\tau_q[i : j]) = \sum_{k=i}^j del(\tau_q[k])$ and $wed(\tau_\emptyset, \tau_d[i : j]) = ins(\tau_d[i : j])$

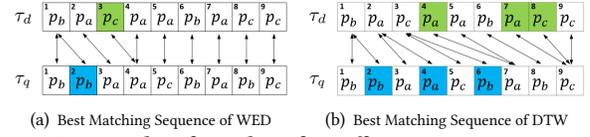


Figure 4: Examples of Matching for Difference Distance Function $j) = \sum_{k=i}^j ins(\tau_d[k])$. The $wed(\tau_q, \tau_d)$ is defined recursively:

$$wed(\tau_q[1 : i], \tau_d[1 : j]) = \min \begin{cases} wed(\tau_q[1 : i - 1], \tau_d[1 : j - 1]) + sub(\tau_q[i], \tau_d[j]) \\ wed(\tau_q[1 : i], \tau_d[1 : j - 1]) + ins(\tau_d[j]) \\ wed(\tau_q[1 : i - 1], \tau_d[1 : j]) + del(\tau_q[i]) \end{cases} \quad (2)$$

DTW. Another well-known distance function is DTW. Unlike WED, there is no deletion and insertion in DTW, instead, multiple points are allowed to be substituted for the same point in another trajectory. However, we try to interpret DTW from a different perspective to make it applicable to the algorithm proposed in this paper. We interpret the original substitution relation as a matching. We consider that only one point $\tau_q[i]$ is substituted for a point $\tau_d[j]$ in another trajectory, while other points that substitute $\tau_d[j]$ are deleted. We can define the insertion in the same way. The cost of deleting a point or inserting a point in the query trajectory is different, depending on which point it matches with, that is, $del(\tau_q[i]) = sub(\tau_q[i], \tau_d[j])$ and $ins(\tau_d[j]) = sub(\tau_q[i], \tau_d[j])$ if $\tau_q[i]$ matches $\tau_d[j]$.

Here, we give an example about the optimal matching when using DTW as distance function.

Example 2. The optimal matching when converting the τ_q into τ_d is shown in Figure 4(b). We set the distance between two points to 1 in case the two points are not equal; otherwise, it is set to 0. When the matching sequence is $[1, 1, 2, 4, 5, 6, 7, 8, 9]$, the conversion cost corresponding to each point is $[0, 0, 1, 1, 1, 0, 0, 1, 0]$. When $\tau_q[4]$ is converted into $\tau_d[4]$, $\tau_d[3]$ needs to be inserted. Therefore, although $\tau_q[4] = \tau_d[4]$, the required cost is still 1. Therefore, the conversion cost of the matching sequence corresponding to Figure 4(a) is 4. When the matching sequence is $[1, 1, 2, 2, 3, 3, 5, 6, 9]$, the conversion cost corresponding to each point is $[0, 0, 0, 0, 0, 1, 0, 0, 1]$. When converting $\tau_q[9]$ into $\tau_d[9]$, the cost of inserting $\tau_d[7]$ is 1. Thus, the conversion cost of this matching sequence corresponding to Figure 4(b) is 2.

Finally, we also give the dynamic process for the calculation of $dtw(\tau_q, \tau_d)$ as follows:

$$dtw(\tau_q[1 : i], \tau_d[1 : j]) = \begin{cases} \sum_{k=1}^j sub(\tau_q[1], \tau_d[k]), & i = 1 \\ \sum_{k=1}^i sub(\tau_q[k], \tau_d[1]), & j = 1 \\ \min \{ dtw(\tau_q[1 : i - 1], \tau_d[1 : j]), \\ dtw(\tau_q[1 : i], \tau_d[1 : j - 1]), \\ dtw(\tau_q[1 : i - 1], \tau_d[1 : j - 1]) \} \\ + sub(\tau_q[i], \tau_d[j]), & \text{else} \end{cases} \quad (3)$$

The algorithm proposed in this paper requires that the distance function to satisfy a specific property: the distance of points between different trajectories is independent of the position of the point in the trajectory. We will explain this in Section 5.3.

2.3 Problem Definition

Definition 6 (Similar Subtrajectory Search Problem, SSS). Given a query trajectory τ_q and a data trajectory τ_d , we expect a closest subtrajectory $\tau_d[i^* : j^*]$ under a specific distance function Θ (e.g., WED or DTW) from the data trajectory for the query trajectory τ_q :

$$(i^*, j^*) = \arg \min_{1 \leq i \leq j \leq n} \Theta(\tau_q, \tau_d[i : j])$$

A more general query is to find the *top-K* similar subtrajectories from massive data trajectories for the query trajectory. Instead,

| Symbol | Description |
|--------------------------------|---|
| τ_d | a data trajectory |
| τ_q | a query trajectory |
| $\tau[i : j]$ | a subtrajectory of τ from i^{th} point to j^{th} point |
| $\tau[i]$ | the i^{th} point in trajectory τ |
| $\mathcal{A}_{\tau_q, \tau_d}$ | a matching sequence between τ_q and τ_d |
| a_i | the matches of $\tau_q[i]$ and $\tau_d[a_i]$ |
| Θ | the distance function |

we can follow such a search process in previous work [26] that maintains the most similar K trajectories and updates it when a more similar subtrajectory appears. Then, we mainly consider querying the most similar subtrajectory from the data trajectory. Details of top-K SSS can be found in Appendix E of our report [11].

Suppose the length of a data trajectory is n , which means that a data trajectory has $\frac{n(n+1)}{2}$ subtrajectories. Assuming that the length of a query trajectory is m and the complexity of computing the distance between the data trajectory and the query trajectory is $O(mn)$ [12]. Therefore, given a query trajectory τ_q and a data trajectory τ_d , the time complexity of searching a subtrajectory of τ_d with the smallest distance from τ_q in τ_d is $O(mn^3)$. Table 1 summarizes the commonly used notations in this paper.

3 REVIEW OF EXISTING SOLUTIONS

We briefly review the existing exact algorithms for the SSS problem. **ExactS.** The vast majority of distance functions [2, 4, 5, 13, 22, 24, 29–31] are defined via recursive processes. Using dynamic programming, we can compute the trajectory distance of a query trajectory and a subtrajectory of the data trajectory in $O(mn)$, where m and n are the lengths of the query trajectory and the data trajectory, respectively. For a query trajectory τ_q and a data trajectory τ_d , let $M_{x,y}$ denote the trajectory distance between $\tau_q[1 : x]$ and $\tau_d[1 : i + y]$ for a given iteration i . ExactS [27] can compute $M_{x,y}$ from $M_{x,y-1}$ using a dynamic programming technique. Thus, line 4 in Algorithm 1 can be solved in $O(mn)$. There are n iterations, thus the overall time complexity of ExactS is $O(mn^2)$. ExactS can be applied to most of the distance functions.

Spring. Spring algorithm [20] is based on the existing dynamic programming computational procedure of DTW and changes the initialization procedure of $dtw(\tau_q[1 : i], \tau_d[1 : j])$ in the Equation 3 when $i = 1$. Spring considers $\tau_d[1 : j - 1]$ to be redundant when $i = 1$; therefore, they modify the equation for $dtw(\tau_q[1 : i], \tau_d[1 : j])$ when $i = 1$ to be as follows:

$$dtw(\tau_q[1 : i], \tau_d[1 : j]) = sub(\tau_q[1], \tau_d[j]) \quad (4)$$

In addition, the authors demonstrate that a modification of the Equation 3 enables it to compute the optimal subtrajectory. However, this trick can only be applied to the DTW function and cannot be extended to other distance functions (e.g., ERP, EDR, and WED). **Greedy Backtracking (GB).** GB [8] investigates finding the optimal subtrajectory in a data trajectory when using FD as the distance function. It constructs a matrix X , where $X_{i,j}$ denotes the Euclidean distance between $\tau_q[i]$ and $\tau_d[j]$. Assuming that $X_{1,1}$ denotes the upper left corner of the matrix, GB finds a path from the top to the right or down until it reaches the bottom. The path's cost is the maximum value in the matrix through which the path passes, and GB finds the optimal subtrajectory by finding the path with

Algorithm 1: *ExactS*(τ_q, τ_d) [27]

Input: a query trajectory τ_q , a data trajectory τ_d
Output: a subtrajectory $\tau_d[i^*, j^*]$

```

1  $i^* \leftarrow 0, j^* \leftarrow 0$ 
2  $score \leftarrow \infty$ 
3 forall  $1 \leq i \leq n$  do
4    $M \leftarrow DP(\tau_q, \tau_d[1 : n])$ 
5    $y^* \leftarrow \arg \min_{1 \leq y \leq n-i+1} M_{m,y}$ 
6   if  $M_{i,y^*} < score$  then
7      $score \leftarrow M_{i,y^*}$ 
8      $i^* \leftarrow i$ 
9      $j^* \leftarrow y^* + i - 1$ 
10 return  $\tau_d[i^*, j^*]$ 

```

the lowest cost. Since FD only considers substitution operations between the trajectory point and trajectory point, it can construct the matrix S . However, the cost of converting $\tau_q[i]$ into $\tau_d[j]$ in other distance functions that consider insertion and deletion operations (e.g., ERP, EDR, and WED) is uncertain; thus, the matrix S cannot be constructed and GB is not suitable.

4 CONVERSION-MATCHING ALGORITHM

This section presents an efficient and exact subtrajectory search algorithm, namely Conversion-Matching Algorithm (CMA). Firstly, we transform the problem of finding the optimal subtrajectory into a problem of finding the optimal matching sequence. Meanwhile, we introduce *the cost of optimal partial matching* $C_{i,j}$ to find the optimal matching sequence. Here, $C_{i,j}$ denotes the minimal cost of converting $\tau_q[1 : i]$ into a subtrajectory of $\tau_d[1 : j]$ when $\tau_q[i]$ matches $\tau_d[j]$ (i.e., $a_i = j$). Note that, converting $\tau_q[1 : i]$ into $\tau_d[1 : j]$ does not mean that $\tau_q[1]$ must match $\tau_d[1]$. Finally, we propose the Conversion-Matching Algorithm (CMA) to calculate $C_{i,j}$ and find the optimal subtrajectory.

4.1 Optimal Matching Sequence

Although previous work [26] has optimized the time complexity of this problem to $O(mn^2)$, it still makes the computational cost increase dramatically when the length of the data trajectory is large. This paper reduces this time complexity to $O(mn)$ by a different dynamic programming algorithm based on a new concept.

Different from existing algorithms, the algorithm is not based on the existing dynamic programming method for calculating the distance. Instead, the basic idea of the algorithm is to calculate the minimum cost of converting the points in the query trajectory to the data trajectory by three operations: insertion, deletion and substitution. Each point in the query trajectory is converted to its matching point in the data trajectory at a specific cost in the conversion process. We can prove that the optimal subtrajectory do not contain redundant prefix trajectories and suffix trajectories by following theorem.

Theorem 4.1. Assume that $\tau_d[i : j]$ is the optimal subtrajectory in τ_d , i.e., $\Theta(\tau_q, \tau_d[i : j]) = \min_{1 \leq s \leq t \leq n} \Theta(\tau_q, \tau_d[s : t])$. Then, we have

$$\Theta(\tau_q, \tau_d[i : j]) = \sum_{a_k \in \mathcal{A}_{\tau_q, \tau_d}^o[i:j]} Cost(\tau_q[k], \tau_d[a_k])$$

where $\mathcal{A}_{\tau_q, \tau_d[i:j]}^o$ is the optimal match sequence of τ_q and $\tau_d[i:j]$.

PROOF. We will prove that $a_1 = i$ and $a_m = j$ in $\mathcal{A}_{\tau_q, \tau_d[i:j]}$ when $\tau_d[i:j]$ is the optimal subtrajectory.

Suppose $a_1 = s$ and $a_m = t$ ($s \geq i, t \leq j$), then $\mathcal{A}_{\tau_q, \tau_d[i:j]}^o$ is also a matching sequence of $\tau_d[s:t]$. Therefore, we have

$$\begin{aligned} \Theta(\tau_q, \tau_d[s:t]) &\leq \sum_{a_k \in \mathcal{A}_{\tau_q, \tau_d[i:j]}^o \setminus \{a_1, a_m\}} \text{Cost}(\tau_q[k], \tau_d[a_k]) \\ &\quad + \text{Cost}(\tau_q[1], \tau_d[s]) + \text{Cost}(\tau_q[m], \tau_d[t]) \\ &= \sum_{a_k \in \mathcal{A}_{\tau_q, \tau_d[i:j]}^o} \text{Cost}(\tau_q[k], \tau_d[a_k]) \\ &\leq \Theta(\tau_q, \tau_d[i:j]) \end{aligned}$$

If $s > i$ or $t < j$, then $\tau_d[i:j]$ is not the optimal subtrajectory, which contradicts what is known. Therefore, we have $a_1 = i$ and $a_m = j$. Further, we can obtain

$$\begin{aligned} \Theta(\tau_q, \tau_d[i:j]) &= \min_{\mathcal{A}_{\tau_q, \tau_d[i:j]} \in \mathbb{A}} \sum_{a_k \in \mathcal{A}_{\tau_q, \tau_d[i:j]}} \text{Cost}(\tau_q[k], \tau_d[a_k]) \\ &\quad + \text{Insert}(\tau_d[i: a_1 - 1]) + \text{Insert}(\tau_d[a_m + 1 : j]) \\ &= \sum_{a_k \in \mathcal{A}_{\tau_q, \tau_d[i:j]}^o} \text{Cost}(\tau_q[k], \tau_d[a_k]) \\ &\quad + \text{Insert}(\tau_d[i: a_1 - 1]) + \text{Insert}(\tau_d[a_m + 1 : j]) \\ &= \sum_{a_k \in \mathcal{A}_{\tau_q, \tau_d[i:j]}^o} \text{Cost}(\tau_q[k], \tau_d[a_k]) \end{aligned}$$

Theorem 4.1 proves that we do not need to consider redundant prefix subtrajectory and suffix subtrajectory in the optimal subtrajectory problem but only need to consider minimizing the conversion cost of all matching points. Then, we prove that the optimal matching sequence of optimal subtrajectory is also optimal among all matching sequences between query and data trajectories. \square

Theorem 4.2. Assume that $\mathcal{A}_{\tau_q, \tau_d[i:j]}^o$ is the optimal matching sequence for the optimal subtrajectory $\tau_d[i:j]$, then it is also the optimal among all matching sequences, i.e. $\mathcal{A}_{\tau_q, \tau_d[i:j]}^o = \arg \min_{\mathcal{A}_{\tau_q, \tau_d} \in \mathbb{A}}$

$$\sum_{a_i \in \mathcal{A}_{\tau_q, \tau_d}} \text{Cost}(\tau_q[i], \tau_d[a_i]).$$

PROOF. We assume that the matching sequence $\mathcal{A}_{\tau_q, \tau_d}^p$ is better than $\mathcal{A}_{\tau_q, \tau_d[i:j]}^o$. If $\mathcal{A}_{\tau_q, \tau_d}^p$ is the matching sequence of subtrajectories $\tau_d[i:j]$ and query trajectories, then it contradicts the condition that $\mathcal{A}_{\tau_q, \tau_d[i:j]}^o$ is the optimal matching sequence for $\tau_d[i:j]$; conversely, if $\mathcal{A}_{\tau_q, \tau_d}^p$ is a matching sequence of the subtrajectory $\tau_d[a_1 : a_m]$, then $\tau_d[i:j]$ is not an optimal subtrajectory, which contradicts what is known. \square

By using the theorems 4.1 and 4.2, we can conclude

$$\min_{1 \leq i \leq j \leq n} \Theta(\tau_q, \tau_d[i:j]) = \min_{\mathcal{A}_{\tau_q, \tau_d} \in \mathbb{A}} \sum_{a_i \in \mathcal{A}_{\tau_q, \tau_d}} \text{Cost}(\tau_q[i], \tau_d[a_i]) \quad (5)$$

According to the Equation 5, we reduce the problem of finding the optimal subtrajectory to finding the optimal matching sequence. We split all match sequences \mathbb{A} of the query trajectory τ_q with the data trajectory τ_d according to the matches at different points. We use $\mathbb{A}[a_i = j]$ to denote the set of all matching sequences in \mathbb{A} that satisfy the condition that $\tau_q[i]$ matches $\tau_d[j]$.

Definition 7 (Optimal Partial Matching-Conversion Cost). We denote by $C_{i,j}$ the minimum value of the cost of converting $\tau_q[1:i]$ into a subtrajectory of $\tau_d[1:j]$ when $\tau_q[i]$ matches $\tau_d[j]$, that is,

$$C_{i,j} = \min_{\mathcal{A}_{\tau_q, \tau_d} \in \mathbb{A}[a_i=j]} \sum_{k=1, a_k \in \mathcal{A}_{\tau_q, \tau_d}}^{k=i} \text{Cost}(\tau_q[k], \tau_d[a_k])$$

Once we have calculated $C_{i,j}$, the distance between the query trajectory and the optimal subtrajectory is the minimum conversion cost when $\tau_q[m]$ matches a point in the data trajectory because

$$\begin{aligned} \min_{1 \leq i \leq j \leq n} \Theta(\tau_q, \tau_d[i:j]) &= \min_{\mathcal{A}_{\tau_q, \tau_d} \in \mathbb{A}} \sum_{a_i \in \mathcal{A}_{\tau_q, \tau_d}} \text{Cost}(\tau_q[i], \tau_d[a_i]) \\ &= \min_{1 \leq j \leq n} \min_{\mathcal{A}_{\tau_q, \tau_d} \in \mathbb{A}[a_m=j]} \sum_{k=1, a_k \in \mathcal{A}_{\tau_q, \tau_d}}^{k=m} \text{Cost}(\tau_q[k], \tau_d[a_k]) \quad (6) \\ &= \min_{1 \leq j \leq n} C_{m,j} \end{aligned}$$

Therefore, we will mainly discuss how to compute $C_{i,j}$ in the subsequent section; meanwhile, we will use DTW and WED as examples to illustrate our algorithm in detail in Section 5.

4.2 Universal Calculation of $C_{i,j}$

We discuss how to calculate $C_{i,j}$ and find the subtrajectory with the shortest distance to the query trajectory from the data trajectory for a given query trajectory τ_q and a data trajectory τ_d .

Calculate $C_{i,j}$. To calculate $C_{i,j}$, we have three cases:

- 1) $i = 1$. When $i = 1$, we substitute $\tau_q[1]$ with $\tau_d[j]$, which is $C_{i,j} = \text{Cost}_{\text{sub}}(\tau_q[1], \tau_d[j])$.
- 2) $j = 1$. There are two possible ways of converting $\tau_q[i]$ when $j = 1$: deleting $\tau_q[i]$, which means that $\tau_q[i-1]$ matches $\tau_d[1]$ so that we have $C_{i,j} = C_{i-1,j} + \text{Cost}_{\text{del}}(\tau_q[i], \tau_d[1])$; the other way is to substitute $\tau_q[i]$ with $\tau_d[1]$, which means that $\tau_q[1:i-1]$ will be deleted, resulting in $C_{i,j} = \text{Cost}_{\text{sub}}(\tau_q[1], \tau_d[1]) + \sum_{k=1}^{i-1} \text{Cost}_{\text{del}}(\tau_q[k], \tau_d[1])$. Therefore, we have $C_{i,j} = \min\{C_{i-1,j} + \text{Cost}_{\text{del}}(\tau_q[i], \tau_d[1]), \text{Cost}_{\text{sub}}(\tau_q[1], \tau_d[1]) + \sum_{k=1}^{i-1} \text{Cost}_{\text{del}}(\tau_q[k], \tau_d[1])\}$.
- 3) $1 < i \leq m, 1 < j \leq n$. Considering that the point $\tau_q[i]$ matches $\tau_d[j]$, there are three different conversion possibilities for $\tau_q[i]$ and $C_{i,j} = \min\{\text{delCost}_{i,j}, \text{subCost}_{i,j}, \text{insCost}_{i,j}\}$:

- (a) $\text{delCost}_{i,j}$: deleting $\tau_q[i]$. When $\tau_q[i]$ is deleted, by the definition of matching, $\tau_q[i-1]$ and $\tau_d[j]$ are matched; thus, we have $\text{delCost}_{i,j} = C_{i-1,j} + \text{Cost}_{\text{del}}(\tau_q[i], \tau_d[j])$.
- (b) $\text{subCost}_{i,j}$: substituting $\tau_q[i]$ with $\tau_d[j]$. In this case, $\tau_q[i-1]$ matches $\tau_d[j-1]$; thus, we have $\text{subCost}_{i,j} = C_{i-1,j-1} + \text{Cost}_{\text{sub}}(\tau_q[i], \tau_d[j])$.
- (c) $\text{insCost}_{i,j}$: substituting $\tau_q[i]$ and inserting $\tau_d[k+1:j-1]$. In this situation, $\tau_q[i-1]$ may match $\tau_d[k]$ ($1 \leq k < j-1$). We insert $\tau_d[k+1:j-1]$ and have $C_{i,j} = C_{i-1,k} + \text{Cost}_{\text{ins}(k)}(\tau_q[i], \tau_d[j])$. Considering all possible values of k , $\text{insCost}_{i,j} = \min_{1 \leq k < j-1} C_{i-1,k} + \text{Cost}_{\text{ins}(k)}(\tau_q[i], \tau_d[j])$.

In case 3.(b), our substitution of $\tau_q[i]$ for $\tau_d[j]$ can be seen as inserting an empty trajectory along with the substitution. Therefore, we will discuss 3.(b) and 3.(c) together in the subsequent sections.

To record the start position the optimal subtrajectory, we use $s_{i,j}$ to denote the index of $\tau_q[1]$'s matched point in τ_d , when $\tau_q[i]$ matches $\tau_d[j]$, i.e., the start position of the subtrajectory. Based on the computation process of $C_{i,j}$, we are able to determine which point $\tau_q[i-1]$ matches when $\tau_q[i]$ matches $\tau_d[j]$. Suppose $\tau_q[i-1]$

Algorithm 2: $CMA(\tau_q, \tau_d)$

Input: a query trajectory τ_q , a data trajectory τ_d **Output:** a subtrajectory $\tau_d[i^*, j^*]$

```
1 forall  $1 \leq i \leq m$  do
2   forall  $1 \leq j \leq n$  do
3     if  $i = 1$  then
4        $C_{i,j} \leftarrow Cost_{sub}(\tau_q[i], \tau_d[j])$ 
5        $s_{i,j} \leftarrow j$ 
6     else if  $j = 1$  then
7        $C_{i,j} \leftarrow \min\{C_{i-1,j} + Cost_{del}(\tau_q[i], \tau_d[j]),$ 
8          $Cost_{sub}(\tau_q[i], \tau_d[j]) + \sum_{k=1}^{i-1} Cost_{del}(\tau_q[k], \tau_d[j])\}$ 
9        $s_{i,j} \leftarrow 1$ 
10    else
11       $C_{i,j} \leftarrow \min\{delCost_{i,j}, subCost_{i,j}, insCost_{i,j}\}$ 
12      update  $s_{i,j}$  according to the matches of  $\tau_q[i-1]$ 
13  $j^* \leftarrow \arg \min_{1 \leq j \leq n} C_{m,j}$ 
14  $i^* \leftarrow s_{m,j^*}$ 
15 return  $\tau_d[i^*, j^*]$ 
```

matches $\tau_d[k]$ ($1 \leq k \leq j$), then we have $s_{i,j} = s_{i-1,k}$. Finally, we propose CMA to solve the SSS problem as shown in Algorithm 2.

Complexity. Since $Cost_{sub}(\tau_q[i], \tau_d[j])$ and $Cost_{del}(\tau_q[i], \tau_d[j])$ involve only the substitution and deletion of one trajectory point, their time complexity is $O(1)$; therefore, when $i = 1$, the time complexity of $C_{i,j}$ is $O(1)$. We can calculate $\sum_{k=1}^{i-1} Cost_{del}(\tau_q[k], \tau_d[j])$ when $j = 1$ in advance for any i by preprocessing, and thus we can compute $C_{i,j}$ within the time complexity of $O(1)$. In other cases, we need to calculate $\min\{delCost_{i,j}, subCost_{i,j}, insCost_{i,j}\}$. Therefore, the time complexity of CMA is $O(mn)$. We will discuss how to compute $C_{i,j}$ in $O(1)$ time complexity for a specific distance function (e.g., DTW and WED) in Section 5.

Discussion. Spring and GB can also achieve $O(mn)$ time complexity for SSS problem under DTW and FD distance function, respectively. CMA is different from them. CMA and Spring are all DP methods. The main difference between CMA and Spring is that their recursive formulas are different: Spring's recursive formula is well-designed for DTW function, and just can support DTW; CMA's recursive formula is a more general one, which can support abstract insertion, substitution and deletion operations thus can be applied under most commonly used trajectory distance functions (e.g., DTW, WED and FD). Secondly, Spring will output all the subtrajectories with distances less than a given threshold to the query trajectory, thus some additional computations are involved in the process of Spring, which do not exist in CMA. Greedy Backtracking is in general a breadth-first search method with memorizing techniques.

5 FAST CALCULATING $C_{i,j}$ ON SPECIFIC Θ

In this section, we discuss how to calculate the conversion cost and $C_{i,j}$ for each point of the query trajectory with WED and DTW. Meanwhile, we will explain how $insCost_{i,j}$ can be computed in $O(1)$ time for the two distance functions WED and DTW.

5.1 Minimum Cost $C_{i,j}$ of WED

By introducing the concept of matching, we can convert the distance between trajectories into the cost required to convert points in τ_q

into points in τ_d . Let's discuss the cost of converting each point $\tau_q[i]$ to its matched point $\tau_d[j]$ in τ_d .

Conversion Cost. There are three cases:

(a) $\tau_q[i-1]$ matches $\tau_d[j]$. We delete $\tau_q[i-1]$ so that $Cost_{del}(\tau_q[i], \tau_d[j]) = del(\tau_q[i])$.

(b) $\tau_q[i-1]$ matches $\tau_d[j-1]$. We substitute $\tau_q[i]$ with $\tau_d[j]$, i.e., $Cost_{sub}(\tau_q[i], \tau_d[j]) = sub(\tau_q[i], \tau_d[j])$.

(c) $\tau_q[i-1]$ matches $\tau_d[k]$, where $1 \leq k < j-1$. The cost of converting $\tau_q[i]$ to $\tau_d[j]$ is the summation of $sub(\tau_q[i], \tau_d[j])$ and the cost of inserting the trajectory $\tau_d[k+1 : j-1]$. Therefore, we have $Cost_{ins(k)}(\tau_q[i], \tau_d[j]) = ins(\tau_d[k+1 : j-1]) + sub(\tau_q[i], \tau_d[j])$.

Example 3. Consider the example in Figure 4, where τ_q is converted into τ_d . Since $\tau_q[1]$ has no predecessor node, $\tau_q[1]$ is only substituted for $\tau_d[1]$ with the cost of $sub(\tau_q[1], \tau_d[1])$. $\tau_q[2]$ matches $\tau_d[1]$, but since $\tau_q[1]$ matches $\tau_d[1]$, $\tau_q[2]$ has to be deleted, with the cost of $del(\tau_q[2])$. $\tau_q[4]$ matches $\tau_d[4]$ and $\tau_q[3]$ matches $\tau_d[2]$, thus $\tau_q[4]$ is converted to $\tau_d[4]$ with the cost of $sub(\tau_q[4], \tau_d[4]) + ins(\tau_d[3])$.

Calculate $C_{i,j}$. After obtaining the conversion cost of the points using WED as a distance function, we can calculate $C_{i,j}$. We will discuss the relational equation for $C_{i,j}$ in three cases:

1) $i = 1$. $C_{i,j} = sub(\tau_q[1], \tau_d[j])$.

2) $j = 1$. $C_{i,j} = \min\{C_{i-1,j} + del(\tau_q[i]), sub(\tau_q[i], \tau_d[1]) + del(\tau_q[1 : i-1])\}$.

3) $1 < i \leq m, 1 < j \leq n$. Considering that the point $\tau_q[i]$ matches $\tau_d[j]$, $\tau_q[i]$ may need to be deleted or substituted. Thus, the update of $C_{i,j}$ depends mainly on whether $\tau_q[i]$ is deleted or replaced:

(a) $delCost_{i,j}$. When $\tau_q[i]$ is deleted, by the definition of matching, $\tau_q[i-1]$ and $\tau_d[j]$ are matched and we have $delCost_{i,j} = C_{i-1,j} + del(\tau_q[i])$.

(b) $subCost_{i,j}$ and $insCost_{i,j}$. $\tau_q[i-1]$ may match $\tau_d[k]$ ($1 \leq k < j-1$) while substituting $\tau_q[i]$. In this case, we insert $\tau_d[k+1 : j-1]$ and have $C_{i,j} = C_{i-1,k} + ins(\tau_d[k+1 : j-1]) + sub(\tau_q[i], \tau_d[j])$. Considering all possible values of k , $insCost_{i,j} = \min_{1 \leq k < j-1} C_{i-1,k} + ins(\tau_d[k+1 : j-1]) + sub(\tau_q[i], \tau_d[j])$. Another situation is that $\tau_q[i-1]$ matches $\tau_d[j-1]$ and we have $subCost_{i,j} = C_{i-1,j-1} + sub(\tau_q[i], \tau_d[j])$. Combining these two situation, we have $C_{i,j} = \min\{insCost_{i,j}, subCost_{i,j}\} = \min_{1 \leq k < j} C_{i-1,k} + ins(\tau_d[k+1 : j-1]) + sub(\tau_q[i], \tau_d[j])$. Then, the calculation of $C_{i,j}$ can be simplified by follows:

$$\begin{aligned} C_{i,j} &= \min_{1 \leq k < j} C_{i-1,k} + ins(\tau_d[k+1 : j-1]) \\ &\quad + sub(\tau_q[i], \tau_d[j]) \\ &= \min\{ \min_{1 \leq k < j-1} C_{i-1,k} + ins(\tau_d[k+1 : j-1]) \\ &\quad + sub(\tau_q[i], \tau_d[j]), C_{i-1,j-1} + sub(\tau_q[i], \tau_d[j]) \} \\ &= \min\{C_{i,j-1} + ins(\tau_d[j-1]) - sub(\tau_q[i], \tau_d[j-1]) \\ &\quad + sub(\tau_q[i], \tau_d[j]), C_{i-1,j-1} + sub(\tau_q[i], \tau_d[j]) \} \end{aligned}$$

By the above analysis, we can obtain the expression for the calculation of $C_{i,j}$ while using WED as distance function

$$C_{i,j} = \begin{cases} sub(\tau_q[i], \tau_d[j]), & i = 1 \\ \min\{C_{i-1,j} + del(\tau_q[i]), sub(\tau_q[i], \tau_d[1]) \\ + del(\tau_q[1 : i-1])\}, & j = 1, i \neq 1 \\ \min\{C_{i-1,j} + del(\tau_q[i]), \\ C_{i,j-1} + ins(\tau_d[j-1]) - \\ sub(\tau_q[i], \tau_d[j-1]) + sub(\tau_q[i], \tau_d[j]), \\ C_{i-1,j-1} + sub(\tau_q[i], \tau_d[j])\}, & otherwise \end{cases} \quad (9)$$

Finally, we illustrate algorithm with an example as follows:

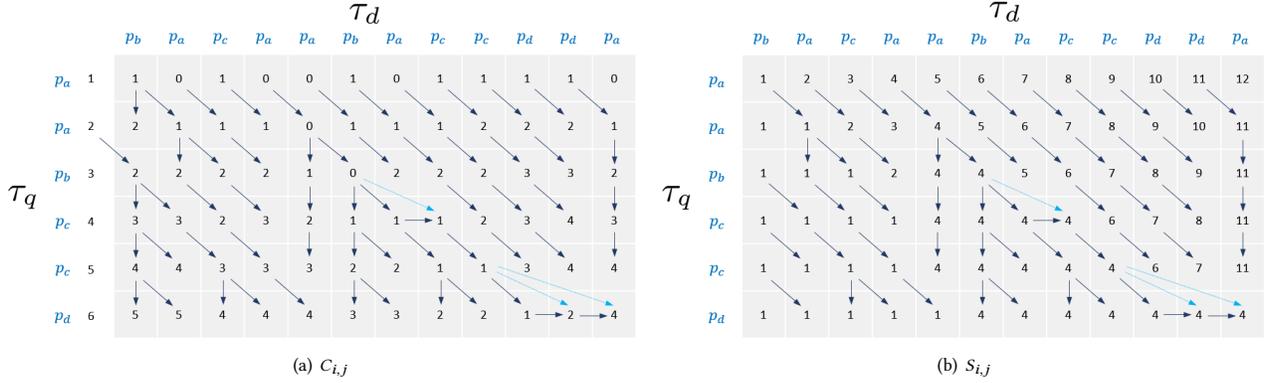


Figure 5: Demonstration of Calculating $C_{i,j}$ and $S_{i,j}$ when using WED as distance function

Example 4. Given two trajectories as shown in Figure 5, we need to find the subtrajectory from τ_d that is closest to τ_q . The insertion, deletion, and substitution costs are the same as the settings in Example 4(a). At the beginning, we will initialize $C_{1,j}$ (i.e., $C_{1,j} = \text{sub}(\tau_q[1], \tau_d[j])$). Then initialize $C_{i,1}$ based on whether $\tau_q[i]$ matches $\tau_d[1]$. Figure 5(a) shows that $\tau_d[1] = b$, thus only $\tau_q[3]$ is substituted with it. For $\tau_q[4]$, when it matches $\tau_d[8]$, we need to determine which point is optimal for $\tau_q[3]$ to match with. From Figure 5(a), we can see that the cost of $\tau_q[3]$ when matching with $\tau_d[6]$ is 0, being the minimum, which means we need to insert $\tau_d[7]$. Therefore, considering $\tau_q[4] = \tau_d[8]$, we can compute the result of $C_{4,8}$ from $C_{3,6}$, i.e., $C_{4,8} = C_{3,6} + \text{ins}(\tau_d[7]) + \text{sub}(\tau_q[4], \tau_d[8]) = 0 + 1 + 0 = 1$. In the actual implementation of the algorithm 2, we will compute $C_{4,8}$ by $C_{4,7}$, i.e. $C_{4,8} = C_{4,7} + \text{ins}(\tau_d[7]) - \text{sub}(\tau_q[4], \tau_d[7]) + \text{sub}(\tau_q[4], \tau_d[8]) = 1 + 1 - 1 + 0 = 1$.

On the other hand, the algorithm updates $S_{i,j}$ as it executes. For example, when $\tau_q[4]$ matches $\tau_d[8]$, $\tau_q[3]$ is matched with $\tau_d[6]$ and we have $S_{4,8} = S_{3,6}$ as shown in Figure 5(b).

5.2 Minimum Cost $C_{i,j}$ of DTW

Unlike WED, the cost to delete a point or insert a point in DTW is different. We analyze the cost of converting each point $\tau_q[i]$ in the query trajectory to its matched point $\tau_d[j]$ in the data trajectory.

Conversion Cost. There are three cases:

(a) $\tau_q[i-1]$ matches $\tau_d[j]$. The cost of deleting $\tau_q[i]$ is equal to the cost of substituting $\tau_q[i]$ with $\tau_d[j]$, thus $\text{Cost}_{del}(\tau_q[i], \tau_d[j]) = \text{sub}(\tau_q[i], \tau_d[j])$.

(b) $\tau_q[i-1]$ matches $\tau_d[j-1]$. We substitute $\tau_q[i]$ with $\tau_d[j]$, i.e., $\text{Cost}_{sub}(\tau_q[i], \tau_d[j]) = \text{sub}(\tau_q[i], \tau_d[j])$.

(c) $\tau_q[i-1]$ matches $\tau_d[k]$, where $1 \leq k < j-1$. The cost of converting $\tau_q[i]$ to $\tau_d[j]$ is the summation of $\text{sub}(\tau_q[i], \tau_d[j])$ and the cost of inserting the trajectory $\tau_d[k+1 : j-1]$. The cost for inserting subtrajectories $\tau_d[k+1 : j-1]$ depends on the points matched by $\tau_d[k+1 : j-1]$ at τ_q . Suppose $\tau_q[i-1]$ matches $\tau_d[k]$ and $\tau_q[i]$ will be matched into $\tau_d[j]$. Thus, the cost to insert $\tau_d[k+1 : j-1]$ is $\min_{k \leq t \leq j-1} \sum_{p=k+1}^t \text{sub}(\tau_q[i-1], \tau_d[p]) + \sum_{p=t+1}^{j-1} \text{sub}(\tau_q[i], \tau_d[p])$. Thus, we have $\text{Cost}_{ins(k)}(\tau_q[i], \tau_d[j]) = \min_{k \leq t \leq j-1} \sum_{p=k+1}^t \text{sub}(\tau_q[i-1], \tau_d[p]) + \sum_{p=t+1}^{j-1} \text{sub}(\tau_q[i], \tau_d[p]) + \text{sub}(\tau_q[i], \tau_d[j])$.

Example 5. Let's take Figure 4(b) as an example, the cost of converting $\tau_q[1]$ to $\tau_d[1]$ when $i = 1$ and $j = 1$ is $\text{sub}(\tau_q[1], \tau_d[1]) = \text{sub}(b, b)$. When $i = 2$ and $j = 1$, $\tau_q[2]$ can only be converted to

$\tau_d[1]$, thus the cost of the conversion is $\text{sub}(\tau_q[2], \tau_d[1]) = \text{sub}(b, b)$. By the time $\tau_q[4]$ matches $\tau_q[2]$, we need to delete $\tau_q[4]$ requiring a cost of $\text{del}(\tau_q[4]) = \text{sub}(\tau_q[4], \tau_q[2])$ because $\tau_q[3]$ matches $\tau_q[2]$. For $i=9$ and $j=9$, since $\tau_q[8]$ matches $\tau_d[7]$, the cost of converting $\tau_q[9]$ to $\tau_d[9]$ consists of not only the cost of the substitution $\text{sub}(\tau_q[9], \tau_d[9])$, but also the cost of inserting $\tau_d[8]$, that is, $\min_{7 \leq t \leq 8} \sum_{p=8}^t \text{sub}(\tau_q[8], \tau_d[p]) + \sum_{p=t+1}^8 \text{sub}(\tau_q[9], \tau_d[p])$. It is equal to $\min\{\text{sub}(\tau_q[8], \tau_d[8]), \text{sub}(\tau_q[9], \tau_d[8])\}$. It can be understood in another way that when $\tau_q[8]$ matches $\tau_d[7]$ and $\tau_q[9]$ matches $\tau_d[9]$, inserting $\tau_d[8]$ is equivalent to replacing $\tau_d[8]$ with $\tau_q[8]$ or $\tau_q[9]$.

Calculate $C_{i,j}$. After analyzing the conversion cost, similarly, we discuss the computation of $C_{i,j}$ in three cases:

(a) $i = 1$. When $i = 1$, $\tau_q[1]$ can only be substituted with $\tau_d[j]$ as the same as WED and we have $C_{i,j} = \text{sub}(\tau_q[1], \tau_d[j])$.

(b) $j = 1$. Considering that the cost of deleting $\tau_q[i]$ and substituting $\tau_q[i]$ is the same when $j = 1$, we have

$$\begin{aligned} C_{i,j} &= \min\{\text{Cost}_{sub}(\tau_q[i], \tau_d[j]) + \sum_{k=1}^{i-1} \text{Cost}_{del}(\tau_q[k], \tau_d[j]), \\ &\quad C_{i-1,j} + \text{Cost}_{del}(\tau_q[i], \tau_d[j])\} \\ &= \min\{\sum_{k=1}^i \text{sub}(\tau_q[k], \tau_d[j]), C_{i-1,j} + \text{sub}(\tau_q[i], \tau_d[1])\} \\ &= C_{i-1,j} + \text{sub}(\tau_q[i], \tau_d[1]) \end{aligned}$$

(c) $1 < i \leq m, 1 < j \leq n$. If we delete $\tau_q[i]$, we have $\text{delCost}_{i,j} = C_{i-1,j} + \text{sub}(\tau_q[i], \tau_q[j])$. Another conversion is substitution. $\tau_q[i-1]$ may be matched with any $\tau_d[k]$ ($1 \leq k < j$), and $C_{i,j}$ denotes the smallest of all possible values. Thus, we have

$$\begin{aligned} C_{i,j} &= \min\{\text{insCost}_{i,j}, \text{subCost}_{i,j}\} \\ &= \min\{\min_{1 \leq k < j-1} C_{i-1,k} + \text{Cost}_{ins(k)}(\tau_q[i], \tau_d[j]), \\ &\quad C_{i-1,j-1} + \text{Cost}_{sub}(\tau_q[i], \tau_d[j])\} \\ &= \min_{1 \leq k < j} C_{i-1,k} + \min_{k \leq t < j-1} \sum_{p=k+1}^t \text{sub}(\tau_q[i-1], \tau_d[p]) \\ &\quad + \sum_{p=t+1}^{j-1} \text{sub}(\tau_q[i], \tau_d[p]) + \text{sub}(\tau_q[i], \tau_d[j]) \\ &= \min_{1 \leq k < j} \min_{k \leq t < j-1} C_{i-1,k} + \sum_{p=k+1}^t \text{sub}(\tau_q[i-1], \tau_d[p]) \\ &\quad + \sum_{p=t+1}^{j-1} \text{sub}(\tau_q[i], \tau_d[p]) + \text{sub}(\tau_q[i], \tau_d[j]) \end{aligned}$$

The time complexity of computing $C_{i,j}$ ($1 < i < m, 1 < j < n$) directly from the above expression is very high, and therefore we are required to simplify the computation of $C_{i,j}$ by Theorem 5.1.

Theorem 5.1. *When $i \geq 2, j \geq 2$, we have $C_{i,j} = \min_{1 \leq k < j} C_{i-1,k} + \sum_{t=k+1}^j sub(\tau_q[i], \tau_d[t])$.*

PROOF. We use mathematical induction to prove this theorem. To simplify the proof, we denote $\sum_{t=k+1}^j sub(\tau_q[i], \tau_d[t])$ as $sub(i, k+1 : j)$. For $\forall j \geq 2$ when $i = 2$, we have

$$\begin{aligned} C_{2,j} &= \min_{1 \leq k < j} \min_{k \leq t < j} C_{1,k} + sub(1, k+1 : t) + sub(2, t+1 : j) \\ &= \min_{1 \leq t < j} \min_{1 \leq k \leq t} sub(1, k : t) + sub(2, t+1 : j) \\ &= \min_{1 \leq t < j} sub(\tau_q[1], \tau_d[t]) + sub(2, t+1 : j) \\ &= \min_{1 \leq t < j} C_{1,t} + \sum_{k=t+1}^j sub(\tau_q[2], \tau_d[k]) \\ &= \min_{1 \leq k < j} C_{1,k} + \sum_{t=k+1}^j sub(\tau_q[2], \tau_d[t]) \end{aligned}$$

Suppose $i = h - 1$, and we have $C_{h-1,j} = \min_{1 \leq k < j} C_{h-2,k} + \sum_{t=k+1}^j sub(\tau_q[h-1], \tau_d[t]) = \min_{1 \leq k < j} C_{h-2,k} + sub(h-1, k+1 : j)$. Next, we have to prove that the theorem also holds when $i = h$.

$$\begin{aligned} C_{h,j} &= \min_{1 \leq k < j} \min_{k \leq t < j} C_{h-1,k} + sub(h-1, k+1 : t) + sub(h, t+1 : j) \\ &= \min_{1 \leq t < j} \min_{1 \leq k \leq t} C_{h-1,k} + sub(h-1, k+1 : t) + sub(h, t+1 : j) \\ &= \min_{1 \leq t < j} \min_{1 \leq k \leq t} \min_{1 \leq l < k} C_{h-2,l} + sub(h-1, l+1 : k) \\ &\quad + sub(h-1, k+1 : t) + sub(h, t+1 : j) \\ &= \min_{1 \leq t < j} \min_{1 \leq l < t} C_{h-2,l} + sub(h-1, l+1 : t) + sub(h, t+1 : j) \\ &= \min_{1 \leq t < j} C_{h-1,t} + sub(h, t+1 : j) \\ &= \min_{1 \leq k < j} C_{h-1,k} + \sum_{t=k+1}^j sub(\tau_q[h], \tau_d[t]) \end{aligned}$$

The above analysis shows that the theorem holds when $i = 2$ and the theorem holds when $i = h - 1$ can infer that the theorem holds when $i = h$. Therefore, the theorem holds. \square

After obtaining the expression for $C_{i,j}$ from the theorem 5.1, we can further simplify it.

$$\begin{aligned} C_{i,j} &= \min_{1 \leq k < j} C_{i-1,k} + \sum_{t=k+1}^j sub(\tau_q[i], \tau_d[t]) \\ &= \min\left\{ \min_{1 \leq k < j-1} C_{i-1,k} + \sum_{t=k+1}^{j-1} sub(\tau_q[i], \tau_d[t]), \right. \\ &\quad \left. C_{i-1,j-1} + sub(\tau_q[i], \tau_d[j]) \right\} \\ &= \min\{C_{i,j-1}, C_{i-1,j-1}\} + sub(\tau_q[i], \tau_d[j]) \end{aligned}$$

Finally, integrating all the previous analysis results, we can get the computational expression of $C_{i,j}$. With the Equation 8, we can quickly adapt the Algorithm 2 to get the optimal subtrajectory using DTW as the distance function.

$$C_{i,j} = \begin{cases} sub(\tau_q[i], \tau_d[j]), & i = 1 \\ C_{i-1,j} + sub(\tau_q[i], \tau_d[1]), & j = 1, i \neq 1 \\ \min\{C_{i-1,j}, C_{i,j-1}, C_{i-1,j-1}\} \\ \quad + sub(\tau_q[i], \tau_d[j]), & otherwise \end{cases} \quad (8)$$

5.3 Other Similarity Functions

In addition to DTW and WED, our method is also valid for other order-insensitive distance functions. EDR and ERP are specific cases of WED functions. Therefore, we only need to define *sub*, *ins*, and *del* in Equation 7. We denote the euclidean distance between two

points $\tau_q[i]$ and $\tau_d[j]$ as $d(\tau_q[i], \tau_d[j])$. We can convert WED to ERP and EDR by defining *sub*, *ins* and *del*: (i) ERP. We can convert WED into ERP by making $sub(\tau_q[i], \tau_d[j]) = d(\tau_q[i], \tau_d[j])$, $del(\tau_q[i]) = d(\tau_q[i], q_c)$, $ins(\tau_d[j]) = d(\tau_d[j], q_c)$, where q_c is a fixed point on the map (e.g., the center of the region). (ii) EDR. $ins(\tau_d[j])$ and $del(\tau_q[i])$ in EDR are both 1, while $sub(\tau_q[i], \tau_d[j])$ takes a value of 0 if and only if $d(\tau_d[j], q_c) < \epsilon$ holds; otherwise, $sub(\tau_q[i], \tau_d[j]) = 1$.

FD is similar to DTW. In the same way, we can obtain the expressions for $C_{i,j}$ when FD is the distance function.

$$C_{i,j} = \begin{cases} sub(\tau_q[i], \tau_d[j]), & i = 1 \\ \max\{C_{i-1,j}, sub(\tau_q[i], \tau_d[1])\}, & j = 1, i \neq 1 \\ \max\{C_{i-1,j}, C_{i,j-1}, C_{i-1,j-1}\}, & otherwise \\ sub(\tau_q[i], \tau_d[j]), & \end{cases} \quad (9)$$

When the order-insensitive distance functions are used, the calculation of the conversation cost does not depend on the position of the current point in the trajectory.

Unfortunately, our method cannot be applied to the subtrajectory search problem when an order-sensitive trajectory distance function (such as LCSS) is used. This is because we do not consider the position from which the subtrajectory starts when computing $C_{i,j}$. When $\tau_q[i]$ matches $\tau_d[j]$, the cost of converting $\tau_q[i]$ to $\tau_d[j]$ is only related to the matching relationship between $\tau_q[i-1]$ and $\tau_d[k]$ ($1 \leq k \leq j$). However, when LCSS is used as the distance function, the cost of converting $\tau_q[i]$ to $\tau_d[j]$ is also related to the matching relation of $\tau_q[1]$, i.e., the starting position of the subtrajectory. The starting position of the subtrajectories has a great influence on judging the distance between the points in two trajectories when LCSS is used as the distance function. Therefore, our algorithm is not suitable for a class of distance functions that considers the position of points in the trajectory, such as LCSS.

6 EXPERIMENTAL STUDY

6.1 Experimental Settings

Data Sets. We conduct experiments on three real data sets: (i) Porto [18] is a dataset describing a whole year (i.e., from July 1st, 2013 to June 30th, 2014) of the trajectories for all the 442 taxis running in the city of Porto (i.e., size: $23.44km \times 24.7km$, longitude: $-8.75^\circ \sim -8.47^\circ$, latitude: $41.02^\circ \sim 41.25^\circ$). There are 1,710,670 trajectories with 15-seconds point intervals, whose average length is 67. (ii) Xi'an Taxi Trip Dataset. DiDi Chuxing GAIA Open Dataset [17] provides a dataset of taxi trips in Xi'an area (i.e., size: $33.43km \times 23.5km$, longitude: $108.78^\circ \sim 109.05^\circ$, latitude: $34.14^\circ \sim 34.38^\circ$). We use the taxi trip records on October 1st. There are 149,742 trajectories with 3-seconds point intervals, whose average length is 401. (iii) T-Drive Data [32, 33]. T-Drive Data provides taxi trips in Beijing area (i.e., size: $49.80km \times 42.11km$, longitude: $116.15^\circ \sim 116.60^\circ$, latitude: $39.75^\circ \sim 40.10^\circ$). There are 10,357 trajectories with 300-seconds point intervals, whose average length is 1705.

In this experiment, we generate Q query trajectories from all trajectories and take the average of the results. We select Q trajectories in uniform random as query trajectory, while the other trajectories are used as data trajectories. We set Q to 100 by default. **Searching Algorithms.** We mainly compare our CMA algorithm with the following existing methods:

1) ExactS. When it computes the distance between the query trajectory and some subtrajectories of the data trajectory, it records

Table 2: Effectiveness of Algorithms.

| Dataset | Algorithm | DTW | | | EDR | | | ERP | | | FD | | |
|---------|-----------|----------|----------|-----------|----------|----------|-----------|----------|----------|-----------|----------|----------|-----------|
| | | AR | MR | RR |
| Porto | POS | 3.0335 | 351.01 | 13.09% | 1.4327 | 321.91 | 15.35% | 1.4982 | 58.99 | 1.83% | 2.9362 | 210.86 | 5.02% |
| | PSS | 1.9760 | 128.91 | 6.80% | 1.3527 | 237.01 | 9.11% | 2.5320 | 139.10 | 5.71% | 1.3780 | 154.50 | 2.60% |
| | RLS | 1.7391 | 97.56 | 5.13% | 1.3435 | 190.54 | 7.32% | 2.2324 | 114.62 | 4.70% | 1.3848 | 134.26 | 2.25% |
| | RLS-Skip | 2.0330 | 142.68 | 7.01% | 1.3545 | 234.13 | 9.28% | 2.4470 | 134.79 | 5.48% | 1.6435 | 173.68 | 3.08% |
| | CMA | 1 | 1 | 0% |
| | ExactS | 1 | 1 | 0% | 1 | 1 | 0% | 1 | 1 | 0% | 1 | 1 | 0% |
| | Spring | 1 | 1 | 0% | - | - | - | - | - | - | - | - | - |
| | GB | - | - | - | - | - | - | - | - | - | - | 1 | 1 |
| Xi'an | POS | 35.5635 | 10505.00 | 18.12% | 1.5162 | 286.84 | 1.14% | 1.4532 | 34.51 | 0.15% | 20.5025 | 3771.50 | 5.31% |
| | PSS | 4.3745 | 676.99 | 2.99% | 1.4608 | 378.33 | 1.34% | 1.7038 | 41.49 | 0.17% | 1.3838 | 25.90 | 0.03% |
| | RLS | 3.6131 | 511.53 | 2.26% | 1.4340 | 304.03 | 1.08% | 1.5648 | 34.32 | 0.14% | 1.3898 | 22.68 | 0.02% |
| | RLS-Skip | 7.3181 | 1567.09 | 4.25% | 1.4596 | 352.07 | 1.26% | 1.6914 | 41.94 | 0.17% | 3.5313 | 434.34 | 0.60% |
| | CMA | 1 | 1 | 0% |
| | ExactS | 1 | 1 | 0% | 1 | 1 | 0% | 1 | 1 | 0% | 1 | 1 | 0% |
| | Spring | 1 | 1 | 0% | - | - | - | - | - | - | - | - | - |
| | GB | - | - | - | - | - | - | - | - | - | - | 1 | 1 |

these intermediate results. Then, ExactS can utilize a dynamic programming technique to optimize the time complexity of searching the optimal subtrajectory from a data trajectory to $O(mn^2)$.

2) PSS and POS. PSS traverses each point of a data trajectory to find the appropriate splitting position. The current optimal subtrajectory is updated by comparing the distance between the subtrajectory before the splitting point and the subtrajectory after the splitting point and the query trajectory. Then, the next suitable splitting point is found starting from the current splitting point. PSS can find an approximate solution with the time complexity of $O(mn)$. As a variant of PSS, POS does not consider the subtrajectory after the splitting point. Thus, the efficiency of POS is substantially improved compared with PSS, but the result quality of PSS is better than that of POS.

3) RLS and RLS-Skip. RLS is an algorithm based on reinforcement learning to determine whether to split the current point, and RLS takes a different action based on the state of the current point. RLS-Skip adds a new action to RLS by skipping the next point to traverse the entire trajectory faster. As a result, RLS-Skip can get a solution in less time, while RLS can find a better solution.

4) Spring and Greedy Backtracking (GB). Both algorithms are of time complexity $O(mn)$. However, Spring and GB can only be applied to specific distance functions, DTW and FD, respectively.

Considering that there are a large number of data trajectories in the database, to improve the efficiency of searching the optimal subtrajectories, we use the pruning methods in the subsequent experiments to filter out the data trajectories that are different from the query trajectories. This paper proposes two modules to filter data trajectories that are not similar to the query trajectory: Filter with Key Points (FKP) and Grid-Based Pruning (GBP). They are compared with the SOTA pruning method, OSF [13]. The details of the pruning methods and their experimental results can be found in the Appendixes B and C of our technical report [11].

Metrics. We will compare our CMA algorithm with the existing algorithms regarding efficiency and effectiveness. For a given query trajectory, we evaluate the efficiency of an algorithm in terms of the time to find the most similar subtrajectory from all data trajectories. We use four evaluation metrics identical to those used in previous

work to evaluate the solutions found by different algorithms in this experiment: (1) Distance. It refers to the raw distance between a query trajectory and the optimal subtrajectory of the data trajectory found by the search algorithm. (2) Approximate Ratio (AR). Given a distance function, AR represents the ratio of the distance between the query trajectory and the subtrajectory found by an approximate algorithm to the distance between the query trajectory and the optimal solution. (3) Mean Rank (MR). It denotes the rank of the distance between the optimal subtrajectory found by the algorithm and the query trajectory among all subtrajectories of the original data trajectory. In particular, $MR=1$ indicates that the algorithm finds the optimal solution. (4) Relative Rank (RR). It is the percentage of all subtrajectories of the data trajectory that is better than the result returned by the algorithm.

Evaluation Platform. The methods are implemented in C++14. The experiments are conducted on a Linux server with 48-cores of Intel(R) Xeon(R) 2.20GHz CPUs and 128.00 GB RAM.

6.2 Experimental Results

Effectiveness compared with other algorithms We used different algorithms for each distance function in different datasets to find the subtrajectories of the data trajectories with the smallest distance from the query trajectory. The experimental results are shown in Table 2. The approximation algorithms have substantial uncertainty in terms of effectiveness. Although the subtrajectory found by these approximation algorithms when using ERP as the distance function is close to the optimal subtrajectory, the subtrajectory found by approximate algorithms when using DTW as the distance function is far from the optimal subtrajectory. POS and PSS tend to select the trajectories with the same length as the query trajectory due to the higher cost of deleting a point when ERP is used as the distance function. In contrast, the length of the optimal subtrajectory tends to vary when DTW is used as the distance function. In addition, the subtrajectories found by the RLS and RLS-Skip algorithms learned based on reinforcement learning are also far from the optimal subtrajectories. CMA can find the exact optimal solution in all cases.

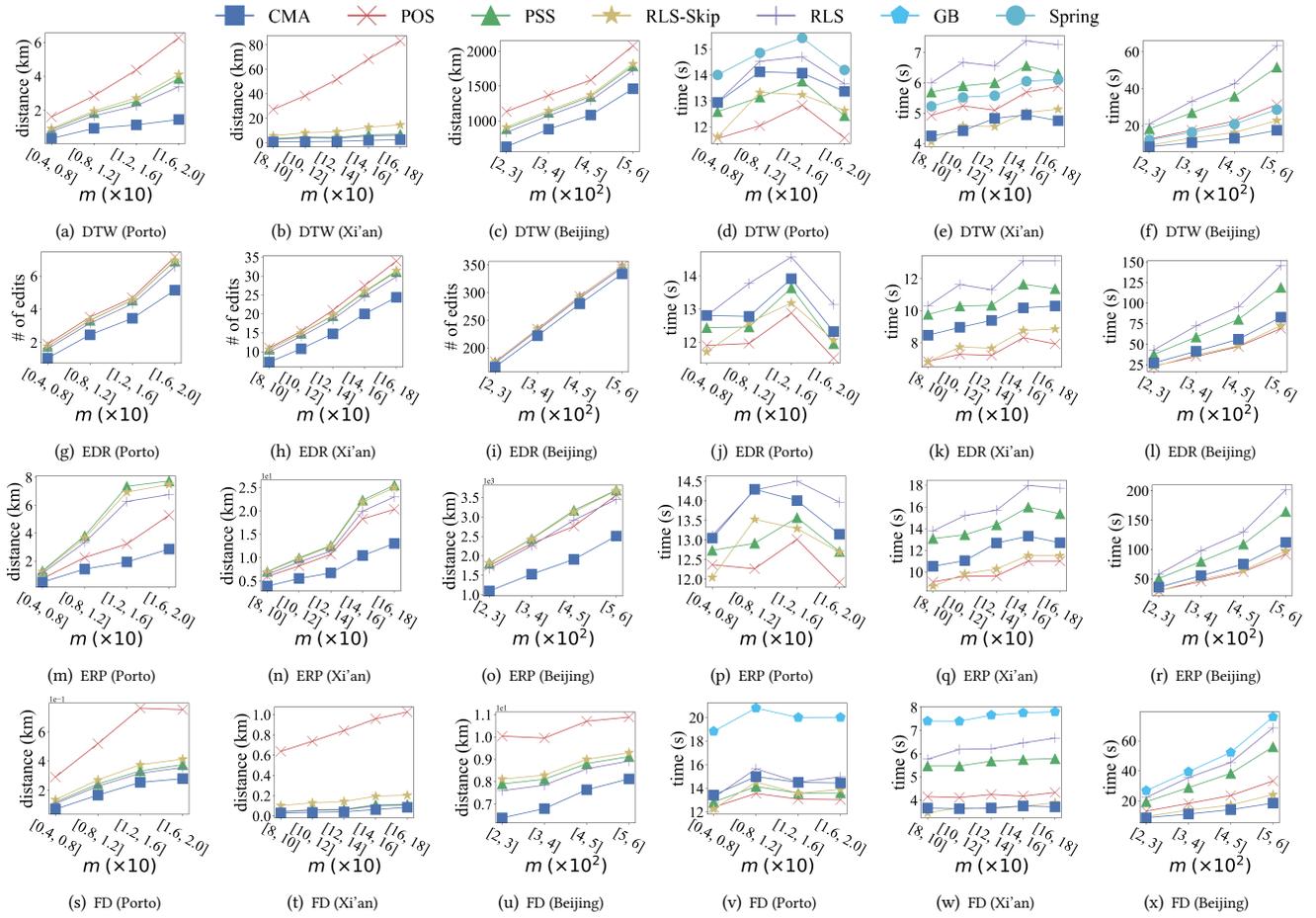


Figure 6: Effectiveness and efficiency with varying query lengths

Efficiency compared with other algorithms With the pruning algorithm, we can find the optimal subtrajectory from many data trajectories faster. Compared with ExactS, the efficiency of CMA has improved nearly 200 times on Xi'an datasets and nearly 50 times on the Porto dataset according to the Table 3. The longer the length of the trajectory, the more the improvement of CMA over ExactS. CMA can find the optimal subtrajectory relatively quickly regardless of the distance function. POS and RLS-Skip are the fastest, but they are approximate algorithms. The experimental results in Table 3 indicate that CMA exhibits superior efficiency compared to other precise algorithms. Compared to CMA, Spring requires many additional computations. In addition to finding the optimal subtrajectory, Spring can identify all subtrajectories whose distances to the query trajectory are less than a given threshold (without overlaps between these subtrajectories). To achieve this, Spring continuously checks the DP matrix for subtrajectories that satisfy the criteria and outputs them, resulting in some additional computations. CMA, on the other hand, performs only one check after completing the calculation of the DP matrix. Spring is specifically designed for large-scale streaming data. The search space of GB is $O(mn)$. However, during the algorithm's execution, backtracking is required repeatedly, which can result in some nodes

being searched multiple times. In contrast, each cell in the DP matrix of CMA is computed only once, making its efficiency slightly higher than that of GB.

Effectiveness of the length of query trajectory. For the Beijing dataset, we select the length ranges m as [200, 300], [300, 400], [400, 500], and [500, 600]. For the Xi'an dataset, we choose the length ranges [80, 100], [100, 120], [120, 140], [140, 160], and [160, 180]. For the Porto dataset, we use the length ranges [4, 8], [8, 12], [12, 16], and [16, 20] for the query trajectories. Figure 6 shows that the execution time of the algorithm increases with the length of the trajectory regardless of the dataset and distance function because the search algorithm takes less time to find the optimal subtrajectory for each query trajectory when the trajectory length is small. However, in the dataset of Porto, the execution time increases and then decreases with the length of the query trajectory, which may be attributed to the fact that there are fewer trajectories similar to the query trajectory in the dataset when the size of the query trajectory becomes longer. Thus, most trajectories are screened out in the filtering phase, resulting in a decrease in the final search time. RLS takes more time than other algorithms in almost all cases. All algorithms except CMA have poor performance when DTW is used as the distance function, regardless of the length of the query trajectory. It is because that DTW allows different points in query

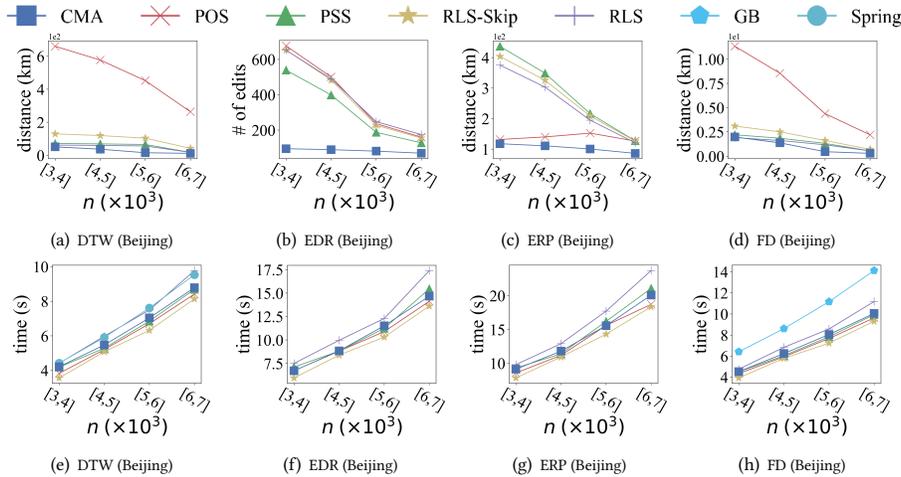


Figure 7: Effectiveness and efficiency with varying data lengths

Table 3: Efficiency of Algorithms.

| Dataset | Algorithm | Time Cost (s) | | | |
|---------|-----------|---------------|----------|----------|----------|
| | | DTW | EDR | ERP | FD |
| Porto | POS | 16.32 | 17.75 | 16.91 | 18.42 |
| | PSS | 18.06 | 16.90 | 17.14 | 18.05 |
| | RLS | 17.84 | 19.87 | 19.39 | 19.62 |
| | RLS-Skip | 16.62 | 15.28 | 17.92 | 18.90 |
| | CMA | 18.78 | 14.64 | 19.26 | 18.78 |
| | ExactS | 7794.59 | 6731.42 | 7225.32 | 8334.16 |
| | Spring | 20.04 | - | - | - |
| | GB | - | - | - | 29.01 |
| Xi'an | POS | 6.69 | 9.69 | 13.12 | 5.48 |
| | PSS | 8.03 | 12.47 | 16.21 | 7.12 |
| | RLS | 7.93 | 14.66 | 18.33 | 7.74 |
| | RLS-Skip | 5.79 | 9.30 | 13.45 | 4.91 |
| | CMA | 5.65 | 9.79 | 14.08 | 4.31 |
| | ExactS | 1625.58 | 2789.93 | 3429.52 | 1312.26 |
| | Spring | 7.37 | - | - | - |
| | GB | - | - | - | 10.76 |
| Beijing | POS | 17.53 | 35.15 | 45.54 | 18.30 |
| | PSS | 26.95 | 58.79 | 79.31 | 28.81 |
| | RLS | 33.17 | 72.37 | 97.63 | 35.46 |
| | RLS-Skip | 13.35 | 36.94 | 48.57 | 13.94 |
| | CMA | 10.81 | 41.53 | 55.18 | 11.29 |
| | ExactS | overtime | overtime | overtime | overtime |
| | Spring | 16.46 | - | - | - |
| | GB | - | - | - | 75.86 |

trajectory matches the same point in data trajectory. Furthermore, the effectiveness of the approximation algorithm is improved as the length of the query trajectory increases when EDR is used as the distance function. As the query trajectory length increases, the number of eligible data trajectories decreases, thus the approximation algorithm has a higher probability of finding the optimal solution. Both algorithms, RLS and RLS-Skip, also find much worse subtrajectories than CMA, where RLS has a worse execution time than CMA in almost all cases.

Effectiveness of the length of data trajectories We varied the length of data trajectories on the Beijing city dataset. In the experiment, we selected 1000 trajectories, each with lengths in the intervals [3000,4000], [4000,5000], [5000,6000], and [6000,7000], from all

trajectories in Beijing city. The experimental results are presented in Figure 7. The figure shows that the time to find the optimal solution increases linearly with the length of the data trajectory for all algorithms. Additionally, the distance of the subtrajectories found by the CMA, Spring, and GB algorithms decreases as the length of the data trajectory increases, indicating that longer data trajectories are more likely to contain subtrajectories that are more similar to the query trajectory. We also observed that longer trajectory lengths make it easier for approximation algorithms to find better solutions. There are more subtrajectories similar to the query trajectory with the increase of the length of data trajectories, which enables the approximation algorithms to find better solutions.

Performance of Spring and GB In this paper, we also explore the performance of Spring and GB; the experimental results of Spring are shown in Figure 6(b) ~ 6(d), while the results of GB are shown in Figure 6(t) ~ 6(v). The experimental results show that the AR of Spring and GB is 1 in all cases, which means that both algorithms can find the optimal solution. However, the execution time of Spring is similar to that of CMA, while GB is less efficient. **Summary of Results.** We verify that CMA can accurately find the nearest subtrajectory from the data trajectory to the query trajectory. Meanwhile, the execution time of CMA is about the same as the two approximation methods (i.e., PSS and POS), and is much smaller than ExactS. Therefore, the proposed algorithm can quickly and accurately find the subtrajectories of the closest data trajectory for each query trajectory.

7 RELATED WORK

Trajectory Distance Function. Many studies have proposed metrics to measure the distance between two trajectories [2, 4, 5, 13, 22, 24, 29–31]. We divided them into two categories: order-insensitive and order-sensitive. The order-insensitive distance functions are independent of the position of the point in the trajectory; the order-sensitive distance functions are just the opposite. For example, the order-insensitive functions, DTW [30] and Fréchet distance (FD) [2], define the distance between trajectories as the cost of turning one trajectory into another through substitution operations. DTW allows different points in one trajectory to be mapped to the same point in another trajectory, enabling DTW to deal well with the

Table 4: Summary of subtrajectory similarity search algorithms.

| Algorithms | Accuracy | order-insensitive | | | | | | | order-sensitive | |
|------------------------------|----------|-------------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------------|-----------|
| | | DTW | ERP | EDR | FD | NetERP | NetEDR | SURS | LCSS | LCRS |
| CMA (Ours) | exact | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | - | - |
| ExactS [27] | exact | $O(mn^2)$ | $O(mn^2)$ | $O(mn^2)$ | $O(mn^2)$ | $O(mn^2)$ | $O(mn^2)$ | $O(mn^2)$ | $O(mn^2)$ | $O(mn^2)$ |
| Spring [20] | exact | $O(mn)$ | - | - | - | - | - | - | - | - |
| Greedy Backtracking (GB) [8] | exact | - | - | - | $O(mn)$ | - | - | - | - | - |
| POS [27] | approx. | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ |
| PSS [27] | approx. | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ |
| RLS [27] | approx. | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ |
| RLS-Skip [27] | approx. | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ | $O(mn)$ |

case where two trajectories are sampled at different frequencies. Edit distance with real penalty (ERP) [4] introduces the insert and delete operations. The cost of inserting a point and deleting a point equals replacing it with a pre-defined default point. However, when the position of the default point is not set reasonably, the cost of deleting and inserting a point can be much greater than replacing it. Then, edit distance on real sequences (EDR) [5] fixes this issue by introducing an upper bound. Specifically, when the distance between a point in the trajectory and its replacement is greater than this upper bound, the replacement cost equals the deletion cost. WED is a generic distance function that allows users to customize the cost of deletion, insertion, and replacement. The order-sensitive distance functions (e.g., longest common subsequence (LCSS) [22], longest overlapping road segments (LORS) [24], and longest common road segments (LCRS) [31]) calculate the distance of a point in a trajectory from another trajectory considering the point positions in the trajectories.

Subtrajectory Search. The previous work [13] divides the subtrajectory search into two stages: filtering and verification. In the filtering phase, most of the trajectories whose distance from the query trajectory exceeds a given threshold are filtered out to reduce the number of validations [7, 27]; in the validation phase, the execution time of the validation phase is simplified with the help of indexes. Unfortunately, this work invokes the trajectory distance function calculation method for all candidate subtrajectories within a trajectory during the validation phase, which makes the validation phase take much time. Another work [27] focuses on how to find the subtrajectory with the minimum distance from the query trajectory in the data trajectory given a query trajectory of length m and a data trajectory of length n . ExactS [27] is proposed to find the optimal subtrajectory in time complexity of $O(mn^2)$. Meanwhile, this work also proposes approximate algorithms (e.g., POS and PSS [27]) with $O(mn)$ time complexity. In addition to these traditional methods, this work proposes two reinforcement learning-based approximate methods (RLS, RLS-Skip [27]) to find the optimal subtrajectory. Furthermore, RLS and RLS-Skip can adaptively select appropriate split points to improve the efficiency of the search. With DTW as the distance function, Spring [20] can find the optimal subtrajectory exactly in $O(mn)$ time complexity. Besides, GB [8] can find the exact optimal similar subtrajectory with $O(mn)$ time complexity on FD. However, Spring and GB do not apply to other distance functions. In contrast, our CMA can be applied to most order-insensitive distance functions. Table 4 summarizes the existing subtrajectory search methods. Experiments on NetERP, NetEDR and SURS can be found in Appendix D of our report [11].

Applications of Subtrajectory search. Some previous studies [25, 28] implement the travel time estimation of a segment of the trajectory by a similar subtrajectory search. One specific process is to search the most similar subtrajectory from the database and then use its time as an estimate of the current trajectory’s communication time. The advantage of subtrajectory search is that it can solve the sparsity of trajectories in the database and thus find more similar trajectories. Another common application is to analyze the movement and behavioral performance of players on the sports ground through subtrajectory search [26]. In addition, subtrajectory search can be used to count the frequency of a given road section in the database for better road planning [6, 12, 16].

8 CONCLUSION

This paper focuses on a similar subtrajectory search problem, i.e., finding the subtrajectory of the data trajectory with the minimum distance for the query trajectory. We convert the problem of finding the optimal subtrajectory to finding the optimal matching sequence. For a given query trajectory of length m and a data trajectory of length n , we propose the CMA algorithm to find the subtrajectory with the minimum distance to the query trajectory from the data trajectory in the time complexity of $O(mn)$. Finally, we conduct sufficient experiments on the datasets of Xi’an, Beijing and Porto, and the experimental results show that our CMA algorithm can find efficiently the exact optimal subtrajectory for each query trajectory.

9 ACKNOWLEDGMENT

Peng Cheng’s work is supported by the National Natural Science Foundation of China under Grant No. 62102149 and Open Foundation of Key Laboratory of Transport Industry of Big Data Application Technologies for Comprehensive Transport. Lei Chen’s work is supported by National Key Research and Development Program of China (2022YFE0200500), National Science Foundation of China (NSFC) under Grant No. U22B2060, the Hong Kong RGC GRF Project 16209519, CRF Project C6030-18G, C2004-21GF, AOE Project AoE/E-603/18, RIF Project R6020-19, Theme-based project TRS T41-603/20R, China NSFC No. 61729201, Guangdong Basic and Applied Basic Research Foundation 2019B151530001, Hong Kong ITC ITF grants MHX/078/21 and PRP/004/22FX, Microsoft Research Asia Collaborative Research Grant, HKUST-Webank joint research lab grant and HKUST Global Strategic Partnership Fund (2021 SJTU-HKUST). Xuemin Lin’s work is supported by NSFC U2241211 and U20B2046. Wenjie Zhang’s work is supported by the Australian Research Council FT210100303 and DP230101445. Corresponding author: Peng Cheng.

REFERENCES

- [1] Pankaj K Agarwal, Kyle Fox, Kamesh Munagala, Abhinandan Nath, Jiangwei Pan, and Erin Taylor. 2018. Subtrajectory clustering: Models and algorithms. In *Proceedings of the 37th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*. 75–87.
- [2] Helmut Alt and Michael Godau. 1995. Computing the Fréchet distance between two polygonal curves. *Int. J. Comput. Geom. Appl.* 5 (1995), 75–91.
- [3] Kevin Buchin, Maike Buchin, Joachim Gudmundsson, Maarten Löffler, and Jun Luo. 2011. Detecting commuting patterns by clustering subtrajectories. *International Journal of Computational Geometry & Applications* 21, 03 (2011), 253–282.
- [4] Lei Chen and Raymond T. Ng. 2004. On The Marriage of Lp-norms and Edit Distance. In *VLDB*. Morgan Kaufmann, 792–803.
- [5] Lei Chen, M. Tamer Özsu, and Vincent Oria. 2005. Robust and Fast Similarity Search for Moving Object Trajectories. In *SIGMOD Conference*. ACM, 491–502.
- [6] Zaiben Chen, Heng Tao Shen, and Xiaofang Zhou. 2011. Discovering popular routes from trajectories. In *ICDE*. IEEE Computer Society, 900–911.
- [7] Christos Faloutsos, Mudumbai Ranganathan, and Yannis Manolopoulos. 1994. Fast subsequence matching in time-series databases. *Acm Sigmod Record* 23, 2 (1994), 419–429.
- [8] Joachim Gudmundsson, Martin P. Seybold, and John Pfeifer. 2021. On Practical Nearest Sub-Trajectory Queries under the Fréchet Distance. In *SIGSPATIAL/GIS*. ACM, 596–605.
- [9] Bo Hui, Da Yan, Haiquan Chen, and Wei-Shinn Ku. 2021. TrajNet: A Trajectory-Based Deep Learning Model for Traffic Prediction. In *KDD*. ACM, 716–724.
- [10] Bo Hui, Da Yan, Haiquan Chen, and Wei-Shinn Ku. 2021. Trajectory WaveNet: A Trajectory-Based Model for Traffic Forecasting. In *2021 IEEE International Conference on Data Mining (ICDM)*. IEEE. <https://doi.org/10.1109/icdm51629.2021.00131>
- [11] Jiabao Jin, Peng Cheng, Lei Chen, Xuemin Lin, and Wenjie Zhang. 2023. Efficient Non-Learning Similar Subtrajectory Search (Technical Report). arXiv:2307.10082 [cs.DB]
- [12] Satoshi Koide, Yukihiro Tadokoro, Takayoshi Yoshimura, Chuan Xiao, and Yoshiharu Ishikawa. 2018. Enhanced Indexing and Querying of Trajectories in Road Networks via String Algorithms. *ACM Trans. Spatial Algorithms Syst.* 4, 1 (2018), 3:1–3:41.
- [13] Satoshi Koide, Chuan Xiao, and Yoshiharu Ishikawa. 2020. Fast Subtrajectory Similarity Search in Road Networks under Weighted Edit Distance Constraints. *Proc. VLDB Endow.* 13, 11 (2020), 2188–2201.
- [14] Jae-Gil Lee, Jiawei Han, and Kyu-Young Whang. 2007. Trajectory clustering: a partition-and-group framework. In *Proceedings of the 2007 ACM SIGMOD international conference on Management of data*. 593–604.
- [15] Xiucheng Li, Kaiqi Zhao, Gao Cong, Christian S Jensen, and Wei Wei. 2018. Deep representation learning for trajectory similarity computation. In *2018 IEEE 34th international conference on data engineering (ICDE)*. IEEE, 617–628.
- [16] Wuman Luo, Haoyu Tan, Lei Chen, and Lionel M. Ni. 2013. Finding time period-based most frequent path in big trajectory data. In *SIGMOD Conference*. ACM, 713–724.
- [17] online. 2016. GAIA Open Dataset. <https://outreach.didichuxing.com/appEnvue/ChengDuOct2016?id=7>.
- [18] online. 2023. Porto Dataset. <https://www.kaggle.com/c/pkdd-15-predict-taxi-service-trajectory-i/data>.
- [19] Sayan Ranu, Padmanabhan Deepak, Aditya D Telang, Prasad Deshpande, and Sriram Raghavan. 2015. Indexing and matching trajectories under inconsistent sampling rates. In *2015 IEEE 31st International conference on data engineering*. IEEE, 999–1010.
- [20] Yasushi Sakurai, Christos Faloutsos, and Masashi Yamamuro. 2007. Stream Monitoring under the Time Warping Distance. In *ICDE*. IEEE Computer Society, 1046–1055.
- [21] Panagiotis Tampakis, Christos Doulkeridis, Nikos Pelekis, and Yannis Theodoridis. 2020. Distributed subtrajectory join on massive datasets. *ACM Transactions on Spatial Algorithms and Systems (TSAS)* 6, 2 (2020), 1–29.
- [22] Michail Vlachos, Dimitrios Gunopulos, and George Kollios. 2002. Discovering Similar Multidimensional Trajectories. In *ICDE*. IEEE Computer Society, 673–684.
- [23] Jiachuan Wang, Peng Cheng, Libin Zheng, Chao Feng, Lei Chen, Xuemin Lin, and Zheng Wang. 2020. Demand-Aware Route Planning for Shared Mobility Services. *Proc. VLDB Endow.* 13, 7 (2020), 979–991.
- [24] Sheng Wang, Zhifeng Bao, J. Shane Culpepper, Zizhe Xie, Qizhi Liu, and Xiaolin Qin. 2018. Torch: A Search Engine for Trajectory Data. In *SIGIR*. ACM, 535–544.
- [25] Yilun Wang, Yu Zheng, and Yexiang Xue. 2014. Travel time estimation of a path using sparse trajectories. In *KDD*. ACM, 25–34.
- [26] Zheng Wang, Cheng Long, Gao Cong, and Ce Ju. 2019. Effective and Efficient Sports Play Retrieval with Deep Representation Learning. In *KDD*. ACM, 499–509.
- [27] Zheng Wang, Cheng Long, Gao Cong, and Yiding Liu. 2020. Efficient and Effective Similar Subtrajectory Search with Deep Reinforcement Learning. *Proc. VLDB Endow.* 13, 11 (2020), 2312–2325.
- [28] Robert Waury, Christian S. Jensen, Satoshi Koide, Yoshiharu Ishikawa, and Chuan Xiao. 2019. Indexing Trajectories for Travel-Time Histogram Retrieval. In *EDBT*. OpenProceedings.org, 157–168.
- [29] Min Xie. 2014. EDS: a segment-based distance measure for sub-trajectory similarity search. In *SIGMOD Conference*. ACM, 1609–1610.
- [30] Byoung-Kee Yi, H. V. Jagadish, and Christos Faloutsos. 1998. Efficient Retrieval of Similar Time Sequences Under Time Warping. In *ICDE*. IEEE Computer Society, 201–208.
- [31] Haitao Yuan and Guoliang Li. 2019. Distributed In-memory Trajectory Similarity Search and Join on Road Network. In *ICDE*. IEEE, 1262–1273.
- [32] Jing Yuan, Yu Zheng, Xing Xie, and Guangzhong Sun. 2011. Driving with knowledge from the physical world. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. 316–324.
- [33] Jing Yuan, Yu Zheng, Chengyang Zhang, Wenlei Xie, Xing Xie, Guangzhong Sun, and Yan Huang. 2010. T-drive: driving directions based on taxi trajectories. In *Proceedings of the 18th SIGSPATIAL International conference on advances in geographic information systems*. 99–108.