

# Sub-trajectory Similarity Join with Obfuscation

Yanchuan Chang  
The University of Melbourne  
Australia  
yanchuanc@student.unimelb.edu.au

Jianzhong Qi  
The University of Melbourne  
Australia  
jianzhong.qi@unimelb.edu.au

Egemen Tanin  
The University of Melbourne  
Australia  
etanin@unimelb.edu.au

Xingjun Ma  
Deakin University  
Australia  
daniel.ma@deakin.edu.au

Hanan Samet  
University of Maryland  
USA  
hjs@cs.umd.edu

## ABSTRACT

User trajectory data is becoming increasingly accessible due to the prevalence of GPS-equipped devices such as smartphones. Many existing studies focus on querying trajectories that are similar to each other in their entirety. We observe that trajectories partially similar to each other contain useful information about users' travel patterns which should not be ignored. Such partially similar trajectories are critical in applications such as epidemic contact tracing. We thus propose to query trajectories that are within a given distance range from each other for a given period of time. We formulate this problem as a sub-trajectory similarity join query named as the *STS-Join*. We further propose a distributed index structure and a query algorithm for STS-Join, where users retain their raw location data and only send obfuscated trajectories to a server for query processing. This helps preserve user location privacy which is vital when dealing with such data. Theoretical analysis and experiments on real data confirm the effectiveness and the efficiency of our proposed index structure and query algorithm.

## CCS CONCEPTS

• **Information systems** → *Spatial-temporal systems*.

## KEYWORDS

Trajectory join, trajectory similarity, spatio-temporal indexing

### ACM Reference Format:

Yanchuan Chang, Jianzhong Qi, Egemen Tanin, Xingjun Ma, and Hanan Samet. 2021. Sub-trajectory Similarity Join with Obfuscation. In *33rd International Conference on Scientific and Statistical Database Management (SSDBM 2021)*, July 6–7, 2021, Tampa, FL, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3468791.3468822>

## 1 INTRODUCTION

Trajectory data is being captured by GPS-equipped devices such as smartphones. Such data can be used to query people's travel

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
SSDBM 2021, July 6–7, 2021, Tampa, FL, USA

© 2021 Association for Computing Machinery.  
ACM ISBN 978-1-4503-8413-1/21/07...\$15.00  
<https://doi.org/10.1145/3468791.3468822>

patterns. In this paper, we are interested in a type of query named as the *trajectory join* query, which returns all trajectory pairs from two trajectory sets that satisfy a given join predicate, e.g., finding people with similar commute routes for ride-sharing matches. Most existing trajectory join queries [26, 27, 33] compute trajectories that are similar in their entirety, i.e., their join predicates are defined on the full trajectories.

We observe that trajectories that are partially similar to each other also offer useful information and should not be ignored. Such partially similar trajectories are gaining importance in applications such as contact tracing for managing epidemics, e.g., to find people in close contact with confirmed cases of COVID-19 for a duration of over 15 minutes<sup>1</sup>. Another example is to compute partially similar trajectories to find matches to form *multi-hop* goods delivery or car-pooling arrangements that allow transits [31].

Motivated by these applications, we propose a trajectory join query that, given two sets of trajectories, computes every pair of trajectories that are within a distance threshold (e.g., 5 meters) lasting for a certain time span (e.g., 15 minutes). Our query is defined on sub-trajectory similarity and hence is named as the *STS-Join*.

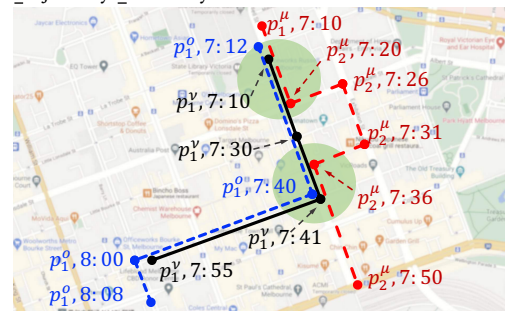


Figure 1: Trajectory join examples

Fig. 1 illustrates our join predicate. There are three trajectories in different colors. The sampled points on the trajectories are marked by the dots and are labelled with their time. Existing studies on full trajectory similarity may return the black and the blue trajectories as a similar pair, because they are very close in both space and time. The existing studies may lose sight of similarity of portions. Within the two green colored areas, the red trajectory is close to the blue and the black ones, respectively, even though their entire movements are quite dissimilar. STS-Join aims to find out both partially and fully similar trajectories.

<sup>1</sup>[www.dhhs.vic.gov.au/victorian-public-coronavirus-disease-covid-19](http://www.dhhs.vic.gov.au/victorian-public-coronavirus-disease-covid-19)

We assume a distributed (i.e., client-server) system architecture to process STS-Join queries. Each client device (e.g., a user’s mobile phone) stores a user’s own accurate trajectories, while the server only stores a modified version of the trajectories for privacy considerations. To showcase the feasibility of processing STS-Join queries over modified trajectories, we consider trajectory obfuscation. A user trajectory is obfuscated automatically where every sampled point on the trajectory is shifted with a bounded-distance noise before the trajectory is sent to the server. All users’ obfuscated trajectories are then maintained on the server in a spatio-temporal index for fast query retrieval. We note that there are limitations in using only obfuscation for trajectory privacy protection, and more advanced techniques exist in the literature such as *geo-indistinguishability* [4]. However, the core theme and contribution of our study is *not* to propose another privacy protection technique. Thus, we just use obfuscation for its simplicity and leave more advanced privacy protection techniques for future studies.

We propose a query algorithm for STS-Join based on a traversal over our index structure. To take advantage of the characteristic that segments on a trajectory are connected sequentially, we design a backtracking-based method to reduce node access of the traversal. It can avoid querying each segment individually. Note that this query method is applicable to any spatial indices that divide the space in a non-overlapping manner. Further, we derive an upper bound of the similarity between a query trajectory and an original user trajectory based on the similarity between the query trajectory and the corresponding obfuscated user trajectory. This enables additional pruning on the server, which reduces the number of query trajectories to be sent to the clients for final similarity checks and saves communication costs.

To sum up, we make the following contributions:

- (1) We define a new trajectory join predicate STS-Join that focuses on sub-trajectory similarity. Two similar trajectories can be very different, but they can contain parts that are related in both space and time. This similarity is especially applicable to bioinformatic datasets that are used for contact tracing and in computational transport science with shared-economy-based transportation systems.
- (2) We propose an efficient spatio-temporal index to manage trajectory data dynamically and a backtracking-based algorithm to process STS-Join queries. We further propose a similarity upper bound that is computed on obfuscated trajectories to enable data pruning and more efficient STS-Join processing. Trajectories in our index are not required to be accurate, but our join results are still correct.
- (3) We conduct experiments on real datasets. The proposed join algorithm outperforms adapted state-of-the-art methods by up to three orders of magnitude in running time.

## 2 RELATED WORK

Our study is related to studies on trajectory similarity measurements, trajectory join queries, and trajectory privacy.

### 2.1 Trajectory Similarity Measurement

Most trajectory similarity measurements are either spatial distance based or spatio-temporal distance based.

Spatial-based measurements [6, 8, 24, 25, 27] focus on the spatial distance between two trajectories. They aggregate distance between aligned point pairs from two trajectories, such as *dynamic time warping* (DTW)[27], *longest common subsequence* (LCSS)[6] and *edit distance on real sequence* (EDR)[8].

Spatio-temporal-based measurements [20, 26, 29, 30, 34] consider the distance in both space and time. For example, LCSS and EDR have been extended to incorporate temporal thresholds [29, 34]. Nanni and Pedreschi [20] assume a constant moving speed on each segment of a trajectory, which is computed as the length of the segment divided by its time span. They then compute the distance between two trajectories as the average distance between two users who travel along the two trajectories with the constant speed of each segment. Shang et al. [26] sum point-to-trajectory distances from one trajectory to the other as the distance between two trajectories, where the summed distance is a weighted sum of the spatial and temporal distances. Wu et al. [30] consider points from two trajectories to be compatible if their spatial distance over time difference is within a velocity threshold.

These measurements compute similarity based on full trajectories which are different from our sub-trajectory-based metric.

### 2.2 Trajectory Join

Studies leverage distributed structures to join trajectories, such as *DITA* [27] and *DISON* [33]. *DITA* supports a variety of trajectory distance functions based on full trajectories, while *DISON* measures trajectory similarity by counting the length of common road segments among the trajectories. Our trajectory join procedure supports distributed processing in two senses: (i) our join refinement procedure is distributed to client machines, and (ii) our index is based on non-overlapping space/time partitioning, which can be easily distributed.

There are also studies on sub-trajectory join [3, 5, 28]. *CSTJ* [5] joins trajectories online. It returns two trajectories once they stay as close-distance pairs for a given time threshold. Its distance measure is based on points in trajectories rather than segments as in our study. *DTJr* [28] also measures point distance. It finds the *longest* sub-trajectory pair satisfying both spatial and temporal distance thresholds, while we find *all* pairs that satisfy a sub-trajectory similarity metric.

*ALSTJ* [3] is the closest work to ours. Its similarity metric is based on the spatial span of sub-trajectory pairs that are within a given distance threshold. The main differences between *ALSTJ* and our STS-Join are in three aspects: (i) *ALSTJ* does not consider the temporal factor and may join trajectories generated at different times, while STS-Join requires the trajectories to be close in both space and time so as to be joined; (ii) *ALSTJ* may have false negatives in its query result due to its trajectory simplification procedure, while our STS-Join guarantees accurate query results; and (iii) *ALSTJ* does not consider user privacy while we do.

### 2.3 Trajectory Privacy

Trajectory privacy has been studied extensively in the last decade [4, 10, 17–19, 32]. For example, a study [32] introduces *position dummy* to hide a user’s location by mixing it with fake locations. Another study [9] extends *k-anonymity* to trajectory data. Studies [4, 17]

leverage *differential privacy* for release of trajectory data. *Geoindistinguishability* [4] generalizes differential privacy to user location data. *SDD* [17] applies the exponential-based randomized mechanism to trajectory data by sampling a rational distance and direction with noises between locations in the trajectory. There are also studies focusing on semantic privacy of trajectories [18, 19]. Monreale et al. [18] present a place taxonomy based method to preserve the trajectory semantic privacy. It guarantees that the probability of inferring the sensitive stops of a user is below a threshold. Naghizade et al. [19] propose an algorithm to protect the semantic information of a trajectory by substituting sensitive stops of a trajectory with less sensitive ones.

As mentioned earlier, our aim is *not* to propose a new privacy protection scheme but to show that it is feasible to process STS-Join queries with a client-server architecture where the server only stores a modified version of the user trajectories. We use obfuscation for its simplicity. Other privacy protection methods (e.g., position dummy) can be applied in our STS-Join if the modification on the trajectory points can be bounded by a distance threshold, which helps guarantee no false negative query results.

**Table 1: Frequently Used Symbols**

Symbol	Description
$\mathcal{D}_p$	An existing trajectory set
$\mathcal{D}_q$	A query trajectory set
$\delta_d, \delta_t$	Trajectory join distance and time thresholds
$\theta_{sp}$	Simplification threshold
$\theta_{ob}$	Maximum obfuscated shifting distance
$\mathcal{T}$	A trajectory
$p_i$	A point in trajectory
$s_i(p_i, p_{i+1})$	A segment in trajectory
$cdd(s_i^\mu, s_j^\nu)$	The close-distance duration of $s_i^\mu$ and $s_j^\nu$
$cdds(\mathcal{T}_\mu, \mathcal{T}_\nu)$	The CDD similarity between $\mathcal{T}_\mu$ and $\mathcal{T}_\nu$

### 3 PRELIMINARIES

Given two sets of trajectories  $\mathcal{D}_p$  (*known trajectory set*) and  $\mathcal{D}_q$  (*query trajectory set*), STS-Join returns pairs of trajectories  $(\mathcal{T}_\mu, \mathcal{T}_\nu) \in \mathcal{D}_p \times \mathcal{D}_q$  with sub-trajectories within  $\delta_d$  distance for at least  $\delta_t$  time, where  $\delta_d$  and  $\delta_t$  are query parameters. Below, we present a few basic concepts and formulate STS-Join. We list the frequently used symbols in Table 1.

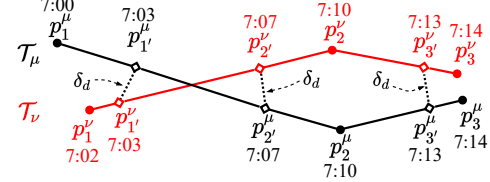
A trajectory  $\mathcal{T}$  is formed by a sequence of  $|\mathcal{T}|$  sampled points  $[p_1, p_2, \dots, p_{|\mathcal{T}|}]$ . A point  $p_i$  is a triple  $\langle x_i, y_i, t_i \rangle$ :  $p_i$  was generated by a user at location  $(x_i, y_i)$  (in Euclidean space) at time  $t_i$ . Two adjacent points  $p_i$  and  $p_{i+1}$  form a segment  $s_i = \overline{p_i, p_{i+1}} \in \mathcal{T}$ .

Following previous studies [3, 15], we consider a constant speed on each trajectory segment. This is valid as real-world trajectories have high sampling rates, e.g., 4.5 seconds in our experiments. The speed may not vary much in such short time frames. Such a constant-speed setting enables computing a user's location at any time  $t \in [t_i, t_{i+1}]$ , denoted as  $(\mathcal{X}_{\mathcal{T}}(t), \mathcal{Y}_{\mathcal{T}}(t))$ , given a trajectory  $\mathcal{T}$ , by linear interpolation:

$$\mathcal{X}_{\mathcal{T}}(t), \mathcal{Y}_{\mathcal{T}}(t) = x_i + \frac{t - t_i}{t_{i+1} - t_i} (x_{i+1} - x_i), y_i + \frac{t - t_i}{t_{i+1} - t_i} (y_{i+1} - y_i) \quad (1)$$

In Fig. 2, there are two trajectories  $\mathcal{T}_\mu = [p_1^\mu, p_2^\mu, p_3^\mu]$  (the black polyline) and  $\mathcal{T}_\nu = [p_1^\nu, p_2^\nu, p_3^\nu]$  (the red polyline). The solid points in the trajectories represent the sample points, and they are labeled with their timestamps, e.g.,  $p_1^\mu$  is recorded at 7:00. Using Equation 1, we can derive a user's location on  $\mathcal{T}_\mu$  (or  $\mathcal{T}_\nu$ ), e.g., at 7:03, the user should be at  $p_{1'}^\mu$ .

Now we can measure the distance between  $\mathcal{T}_\mu$  and  $\mathcal{T}_\nu$ , denoted by  $dist(\mathcal{T}_\mu, \mathcal{T}_\nu)$ , at any time  $t$  as the Euclidean distance. We call this distance the *point distance*. STS-Join computes the time duration where this distance is within a given threshold  $\delta_d$ .



**Figure 2: Example of trajectories**

Trajectories may be generated at different time spans and with different sampling rates. To compute the point distance between  $\mathcal{T}_\mu$  and  $\mathcal{T}_\nu$ , we need to first define the overlapping time span of segment pairs. Let  $ot(s_i^\mu, s_j^\nu)$  be the overlapping time span between  $s_i^\mu$  and  $s_j^\nu$ , i.e.,  $[t_i^\mu, t_{i+1}^\mu] \cap [t_j^\nu, t_{j+1}^\nu]$ , where  $(s_i^\mu, s_j^\nu) \in \mathcal{T}_\mu \times \mathcal{T}_\nu$ . Denote the length of the overlapping time span as  $|ot(s_i^\mu, s_j^\nu)|$ . Then, in Fig. 2,  $ot(s_1^\mu, s_1^\nu)$  is  $[7:02, 7:10]$ , and the length is 8 minutes.

**Close-distance duration.** When  $ot(s_i^\mu, s_j^\nu) < \theta$ , we call  $s_i^\mu$  and  $s_j^\nu$  two *time-overlapping* segments. Given such segments, we compute the time range  $[t_{i,j}^{\mu\nu}, t_{i,j+1}^{\mu\nu}]$  where  $dist(\mathcal{T}_\mu, \mathcal{T}_\nu) \leq \delta_d$ , i.e.,

$$\underbrace{(\mathcal{X}_{\mathcal{T}_\mu}(t) - \mathcal{X}_{\mathcal{T}_\nu}(t))^2 + (\mathcal{Y}_{\mathcal{T}_\mu}(t) - \mathcal{Y}_{\mathcal{T}_\nu}(t))^2}_{\leq \delta_d} \quad (2)$$

To solve this inequality, we expand Equation 2 with Equation 1 and compute the square of both sides of the inequality:

$$\begin{aligned} \delta_d^2 &> (\mathcal{X}_{\mathcal{T}_\mu}(t) - \mathcal{X}_{\mathcal{T}_\nu}(t))^2 + (\mathcal{Y}_{\mathcal{T}_\mu}(t) - \mathcal{Y}_{\mathcal{T}_\nu}(t))^2 \\ &= (k_x^\mu \cdot t + b_x^\mu - k_x^\nu \cdot t - b_x^\nu)^2 \\ &\quad + (k_y^\mu \cdot t + b_y^\mu - k_y^\nu \cdot t - b_y^\nu)^2 \\ &= [(k_x^\mu - k_x^\nu)^2 + (k_y^\mu - k_y^\nu)^2] t^2 \\ &\quad + 2[(k_x^\mu - k_x^\nu)(b_x^\mu - b_x^\nu) + (k_y^\mu - k_y^\nu)(b_y^\mu - b_y^\nu)] t \\ &\quad + (b_x^\mu - b_x^\nu)^2 + (b_y^\mu - b_y^\nu)^2, \text{ where } t \in ot(s_i^\mu, s_j^\nu) \text{ and} \\ k_x^\mu &= \frac{x_{i+1}^\mu - x_i^\mu}{t_{i+1}^\mu - t_i^\mu}, & b_x^\mu &= \frac{x_i^\mu t_{i+1}^\mu - x_{i+1}^\mu t_i^\mu}{t_{i+1}^\mu - t_i^\mu} \\ k_x^\nu &= \frac{x_{j+1}^\nu - x_j^\nu}{t_{j+1}^\nu - t_j^\nu}, & b_x^\nu &= \frac{x_j^\nu t_{j+1}^\nu - x_{j+1}^\nu t_j^\nu}{t_{j+1}^\nu - t_j^\nu} \\ k_y^\mu &= \frac{y_{i+1}^\mu - y_i^\mu}{t_{i+1}^\mu - t_i^\mu}, & b_y^\mu &= \frac{y_i^\mu t_{i+1}^\mu - y_{i+1}^\mu t_i^\mu}{t_{i+1}^\mu - t_i^\mu} \\ k_y^\nu &= \frac{y_{j+1}^\nu - y_j^\nu}{t_{j+1}^\nu - t_j^\nu}, & b_y^\nu &= \frac{y_j^\nu t_{j+1}^\nu - y_{j+1}^\nu t_j^\nu}{t_{j+1}^\nu - t_j^\nu} \end{aligned} \quad (3)$$

The resultant quadratic inequality has just one variable  $t$ . It can be solved straightforwardly by letting the inequality be equal and computing the roots for the equation with the quadratic formula. We omit the detailed computation for conciseness.

In Fig. 2, the distance threshold  $\delta_d$  is represented by the dotted lines. The distance between the first segments of the two trajectories first decreases and then increases. At time  $t = 7:03$  and  $7:07$ ,

$dist(\mathcal{T}_\mu, \mathcal{T}_v) = \delta_d$ , which yields the first time range [7:03, 7:07] that satisfies  $\delta_d$ . The distance between the second segments of the two trajectories keeps decreasing, which reaches  $\delta_d$  at 7:13. This yields the second time range [7:13, 7:14] that satisfies  $\delta_d$ .

We call the length of  $[t_{i,j}^{\mu\nu}, t_{i,j+1}^{\mu\nu}]$  the *close-distance duration* (CDD), denoted as  $cdd(s_i^\mu, s_j^\nu) = t_{i,j+1}^{\mu\nu} - t_{i,j}^{\mu\nu}$ . We define the similarity between two trajectories as their total CDD across all segments.

**Definition 3.1** (*Close-distance duration similarity*). Given a point distance threshold  $\delta_d$ , the *close-distance duration similarity* (CDDS) of two trajectories  $\mathcal{T}_\mu$  and  $\mathcal{T}_v$ ,  $cdds(\mathcal{T}_\mu, \mathcal{T}_v)$ , is the sum of  $cdd(s_i^\mu, s_j^\nu)$  of every time-overlapping segment pair  $(s_i^\mu, s_j^\nu) \in \mathcal{T}_\mu \times \mathcal{T}_v$ .

$$cdds(\mathcal{T}_\mu, \mathcal{T}_v) = \sum_{1 \leq i \leq |\mathcal{T}_\mu|, 1 \leq j \leq |\mathcal{T}_v|, ot(s_i^\mu, s_j^\nu) < \infty} cdd(s_i^\mu, s_j^\nu) \quad (4)$$

CDDS sums up the duration of all close segments including the partial ones. This differs from existing trajectory similarity metrics that require the full trajectories to be close. In Fig. 2,  $cdds(\mathcal{T}_\mu, \mathcal{T}_v)$  equals to the total length of the two time ranges [7:03, 7:07] and [7:13, 7:14], i.e., 5 minutes.

**Problem definition.** Now we can formulate our STS-Join.

**Definition 3.2** (*STS-Join query*). Given two trajectory datasets  $\mathcal{D}_p$  and  $\mathcal{D}_q$ , a point distance threshold  $\delta_d$ , and a close-distance duration similarity threshold  $\delta_t$ , *STS-Join* returns every trajectory pair  $(\mathcal{T}_\mu, \mathcal{T}_v) \in \mathcal{D}_p \times \mathcal{D}_q$  such that  $cdds(\mathcal{T}_\mu, \mathcal{T}_v) > \delta_t$ .

## 4 INDEX STRUCTURE

We assume set  $\mathcal{D}_p$  to be known (e.g., user trajectory dumps) and set  $\mathcal{D}_q$  to be given at query time (e.g., trajectories of newly confirmed COVID-19 cases). We build an index named *STS-Index* over  $\mathcal{D}_p$  such that  $\mathcal{D}_p$  can be STS-Joined with  $\mathcal{D}_q$  efficiently. We use a client-server architecture to protect location privacy. On a client, a user's trajectories are stored in their original form, which are obfuscated and sent to the server. The obfuscated trajectories from all clients together are indexed in a tree structure on the server that considers both their spatial and temporal features. Next, we detail the index structures on the clients and the server, respectively.

### 4.1 On Client Side

A user's original trajectories are stored on the client side. We *simplify* and *obfuscate* an original trajectory before sending it (together with the client ID and a local trajectory ID) to the server in order to reduce the communication and protect user's privacy. We index the trajectories by their local IDs (e.g., using a sorted array or a B-tree) for fast retrieval at the refinement stage of query processing.

**Trajectory simplification.** First, we simplify an original trajectory  $\mathcal{T}_\mu$  by reducing the number of sampled points. This reduces the storage space and improves the query efficiency later. Our simplification algorithm is adapted from the *Douglas-Peucker algorithm* (DP) [13]. The native DP algorithm ignores the temporal dimension. Consider two sampled points  $p_i^\mu$  and  $p_{i+k}^\mu$  ( $k > 1$ ) on  $\mathcal{T}_\mu$ . For all other sampled points between  $p_i^\mu$  and  $p_{i+k}^\mu$ , if their perpendicular distances to the segment between  $p_i^\mu$  and  $p_{i+k}^\mu$  are within a simplification threshold  $\theta_{sp}$  (an empirical parameter), then these points are all removed from  $\mathcal{T}_\mu$  by DP.

In our case, since we interpolate user locations on the trajectory segments, we require the user location on the simplified segment  $\overline{p_i^\mu, p_{i+k}^\mu}$  to be within  $\theta_{sp}$  distance from that on the original segments at every time point  $t \in [t_i^\mu, t_{i+k}^\mu]$ . This guarantees no false negatives in STS-Join over the simplified trajectories. Fig. 3 shows an example. The original (black) trajectory  $\mathcal{T}_\mu$  has three segments, which is simplified to just one (red) segment  $\overline{p_1^\mu, p_4^\mu}$ . We compute  $p_2^\mu$  and  $p_3^\mu$  on  $\overline{p_1^\mu, p_4^\mu}$  at times  $t_2^\mu$  and  $t_3^\mu$  (i.e., the time points of  $p_2^\mu$  and  $p_3^\mu$ ), respectively. To ensure valid simplification, the distance between  $p_2^\mu$  and  $p_3^\mu$  (i.e.,  $d_2$ ) and that between  $p_3^\mu$  and  $p_4^\mu$  (i.e.,  $d_3$ ) must both be within  $\theta_{sp}$ . This contrasts to the native DP that examines the perpendicular distances of  $p_2^\mu$  and  $p_3^\mu$  (i.e.,  $d_2^\perp$  and  $d_3^\perp$ ), which are shorter and may lead to false negatives at query processing.

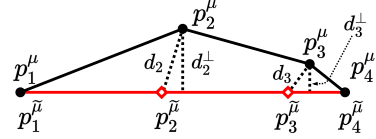


Figure 3: Example of trajectory simplification

**Trajectory obfuscation.** Our simplified trajectories retain a subset of the trajectory points. To protect privacy, we further adapt the *bounded Laplace mechanism* (i.e., an algorithm) to obfuscate the simplified trajectories as inspired by previous studies [14, 16].

The bounded Laplace mechanism adds a noise from the bounded Laplace distribution to the output of a (query) function. In our problem, we add a bounded noise to each dimension of a point on a simplified trajectory to protect users' location privacy. The probability distribution function (PDF) of the added noise is [16]:

$$Pr(x) = \lambda \cdot \frac{1}{2b} \cdot \exp\left(\frac{-|x|}{b}\right), x \in [-\theta_{ob}, \theta_{ob}] \quad (5)$$

where  $\lambda$  is a constant,  $b$  is the bias of the distribution,  $\theta_{ob} \in \mathbb{R}^+$  is the obfuscated distance threshold. Since the domain of the distribution is bounded in  $[-\theta_{ob}, \theta_{ob}]$ , the integral of PDF should be 1 for  $x \in [-\theta_{ob}, \theta_{ob}]$ . This yields the value of  $\lambda$ :

$$\lambda = \left( \int_{-\theta_{ob}}^{\theta_{ob}} \frac{1}{2b} \cdot \exp\left(\frac{-|x|}{b}\right) dx \right)^{-1} = (1 - \exp(-\frac{\theta_{ob}}{b}))^{-1} \quad (6)$$

We then leverage the *inverse cumulative distribution function* to generate random noises that satisfy the PDF. Firstly, we derive the CDF from the PDF by computing the integral of the PDF for  $x < 0$  and  $x > 0$  respectively:

$$F(x) = \begin{cases} \frac{\lambda}{2} \cdot (\exp(\frac{x}{b}) - \exp(\frac{-\theta_{ob}}{b})), & x < 0 \\ \frac{\lambda}{2} + \frac{\lambda}{2} \cdot (1 - \exp(\frac{-x}{b})), & x > 0 \end{cases} \quad (7)$$

Then, we can obtain the inverse CDF from Equation 7:

$$x = \begin{cases} -b \cdot \ln(1 + \lambda^{-1} - 2\lambda^{-1}y), & y \in (0, 0.5] \\ b \cdot \ln(1 - \lambda^{-1} + 2\lambda^{-1}y), & y \in (0.5, 1) \end{cases} \quad (8)$$

Here, variable  $y$  is a random number in  $(0, 1)$ . We generate  $y$  and use it to obtain noise  $x$ , which is then added to each point coordinate of the simplified trajectory. Parameter  $\lambda$  is determined by  $b$  and  $\theta_{ob}$  (Equation 6), while  $b$  is the distribution bias.

## 4.2 On Server Side

The obfuscated trajectories are stored in a tree structure on the server that indexes both the spatial and temporal features of the trajectories (in the form of segments). As Fig. 4 shows, the top levels of the structure together can be seen as a B-tree that indexes time intervals hierarchically, and the node capacity (which is a system parameter) of the example is 20. The time span of an entry in the leaf nodes is 3 minutes which is determined by the time span of the trajectory segments, and every entry points to a spatial index for the trajectory segments in that interval. If a segment spans across the intervals of multiple entries (which is rare as the segments are usually short even after simplification), we break the segment into multiple segments to suit the entry intervals. For example, if a segment starts at 8:29:50 and ends at 8:30:03, we break the segment at 8:30:00 into two segments that suit the time intervals  $[8:27:00, 8:30:00)$  (i.e.,  $N_5$ ) and  $[8:30:00, 8:33:00)$  of the leaf entries.

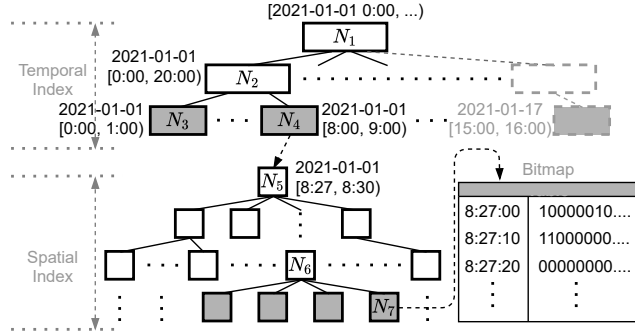


Figure 4: STS-Index server-side structure

We create a *quasi-quadtrees* as the spatial index to obtain nodes with non-overlapping *minimum bounding rectangles* (MBR). We use endpoints of the segments to build this tree. First, we insert the endpoints into the root. Once the root node capacity is reached, we partition the space into four quadrants each forming a child node. Every endpoint is moved to the child node enclosing it, while the root node now stores pointers to the child nodes. The process is done recursively until every node is within its capacity. In a leaf node, we further store the segments overlapped by the node MBR (a segment may be in multiple nodes), together with the IDs of the corresponding trajectories and clients.

When the data distribution is skewed, we may have many segments in the same leaf node. We add bitmaps to help query such segments. The bitmaps correspond to disjoint intervals that together cover the time interval of the segments in the node. The number of bitmaps can be empirically determined based on the sampling rate. Each bit value denotes whether a segment overlaps with a time interval. In Fig. 4, we build a bitmap for  $N_7$  with 10-second intervals, where both the first and the seventh segments in  $N_7$  overlap with the interval  $[8:27:00, 8:27:10)$ .

**Update handling.** When there are new obfuscated trajectories, their segments (and the endpoints of the segments) are added to the index by a top-down traversal to find the nodes to be inserted into (following the query procedure in the next section). Trajectory deletion can be also done by first a tree traversal to locate the trajectory segments to be deleted. Then, the segments are removed, together with any empty nodes.

## 5 STS-JOIN QUERY PROCESSING

We illustrate how to join a query trajectory set  $\mathcal{D}_q$  with a known trajectory set  $\mathcal{D}_p$  indexed in our STS-Index. Set  $\mathcal{D}_q$  is not obfuscated, e.g., the trajectories of confirmed COVID-19 cases are reported to the authority for contact tracing. Firstly, we present our STS-Join query algorithm and optimize it with a backtracking technique. Then, we illustrate how to prune dissimilar pairs on the server by bounding the actual trajectory similarity based on obfuscated trajectories. Finally, we analyze the algorithm cost.

### 5.1 The STS-Join Query Algorithm

A straightforward algorithm is to query every segment from  $\mathcal{D}_q$  independently over STS-Index, identify the data segments satisfying the query, and send the query trajectory to the corresponding clients for verification against the original trajectories.

We observe that, segments of a trajectory are connected end to end. Points of adjacent segments are likely to be in the leaf nodes that are in short-hop distance from each other in the index. Besides, our quasi-quadtrees divide the space without overlaps. By leveraging these features, we propose an efficient backtracking-based join algorithm that starts from the inner nodes instead of the root for querying each segment of a trajectory.

**MBR expansion for accurate query processing.** To ensure no false dismissals, we expand the MBRs of query segments to cover the simplified-and-obfuscated trajectories, as well as the query distance threshold  $\delta_d$ . Overall, we expand the MBR of a query segment by  $\delta_d + \theta_{sp} + \theta_{ob}$  on each side. This ensures no false negatives, because by definition, the distance between a point on an original trajectory segment and the corresponding point on its simplified-and-obfuscated version cannot be greater than  $\delta_d + \theta_{sp} + \theta_{ob}$  in any dimension. We use  $[mbr]_{s_i}$  to denote the expanded MBR of  $s_i$ . In Fig. 5,  $[mbr]_{s_1^v}$  is the expanded MBR of the segment  $s_1^v$ .

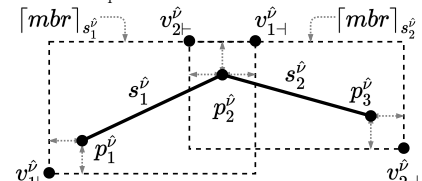


Figure 5: Example of expanded MBRs

**Pivots.** We use *pivots* to help locate the adjacent segments in our algorithm. Pivots are vertices of the expanded MBR of a query segment that are close to the endpoints of the segment. Using a pivot instead of an endpoint helps avoid false negative query results, since the query MBR is expanded. In Fig. 5,  $v_{1-}^v$  and  $v_{1+}^v$  are the pivots of query segment  $s_1^v$ .

**On server side.** As summarized in Algorithm 1, our server-side algorithm iteratively searches for similar segments for the query trajectories (lines 1 to 20). For every query trajectory  $\mathcal{T}_v$ , we split its segments according to the time intervals of the root nodes of the quasi-quadtrees (line 2), like we did in index construction. Then, we search for data segments similar to every segment in the split query trajectory  $\mathcal{T}_v$  (lines 3 to 20). We first update *pivot* and its corresponding leaf node  $N$ . This is only needed for the first segment in  $\mathcal{T}_v$  or when we move on to a segment in a new time interval (lines 4 to 7), which rarely happens. Function  $FindNode(N, pivot)$

**Algorithm 1:** STS-Join (server side)

---

**Input:**  $\mathcal{D}_q$ : query trajectory set;  $\delta_d$ : query distance threshold

**Output:**  $\mathcal{P}$ : the set of STS-Joined trajectory pairs

```

1 for  $\mathcal{T}_v$  in  $\mathcal{D}_q$  do
2    $\mathcal{T}_v \leftarrow$  split  $\mathcal{T}_v$  to suit node intervals of the temporal
   index;
3   for  $s_i^v$  in  $\mathcal{T}_v$  do
4     if  $i$  is 1 or  $s_i^v$  and  $s_{i-1}^v$  are in different
       quasi-quadtrees then
5       pivot  $\leftarrow v_{i-1}^v$ ;
6        $N \leftarrow$  the root of the quasi-quadtrees whose time
       interval overlaps with that of  $s_i^v$ ;
7        $N \leftarrow FindNode(N, pivot)$ ;
8       pivot  $\leftarrow v_{i+1}^v$ ;
9        $N_p \leftarrow Backtrack(N, v_{i+1}^v)$ ;
10       $Q.enqueue(N_p)$ ;
11      while  $Q < \emptyset$  do
12         $N_p \leftarrow Q.dequeue()$ ;
13        for each entry in  $N_p$  do
14          if  $overlap(entry.mbr, [mbr]_{s_i^v})$  then
15            if  $N_p$  is not a leaf node then
16               $Q.enqueue(entry.child)$ ;
17            else
18              add  $\langle s^{\#}, s_i^v \rangle$  into  $\mathcal{S}$ , where  $s^{\#}$  is the
              segment indexed at  $entry$ ;
19          if  $N_p$  is a leaf node and its  $mbr$  covers  $pivot$ 
20            then
               $N \leftarrow N_p$ ;
21 Update every pair  $\langle s^{\#}, s_i^v \rangle \in \mathcal{S}$  to its corresponding
   non-time-interval-split segment pair  $\langle s^{\#}, s^v \rangle$ ;
22 Group the pairs in  $\mathcal{S}$  by the client ID of  $s^{\#}$ , and send
   corresponding pairs to the clients for further verification;
23  $\mathcal{P} \leftarrow$  the set of trajectory pairs that are returned from
   clients;
24 return  $\mathcal{P}$ ;

```

---

locates node  $N$  whose MBR covers the pivot in the quasi-quadtrees by a point query (line 7). Then, we can leverage node  $N$  to backtrack with function  $Backtrack(N, p)$  that starts from  $N$  and recursively traces back to the ancestor node whose MBR covers  $p$ . For each query segment, we stop the backtracking at the ancestor node  $N_p$  that covers the upper pivot  $v_{i+1}^v$  of the expanded MBR of the current query segment (lines 8 and 9). Then, we search for similar segment pairs from  $N_p$  for the current query segment (lines 10 to 20). For pruning, only tree nodes whose MBRs overlap with the expanded MBR  $[mbr]_{s_i^v}$  of the query segment are visited, which are stored in a queue  $Q$  to support the traversal (lines 14 to 16). When the traversal reaches a simplified-and-obfuscated segment  $s^{\#}$ , we add  $\langle s^{\#}, s_i^v \rangle$  to the result set  $\mathcal{S}$  (line 18). Meanwhile, we verify

each leaf node for whether it contains the lower pivot of the next query segment which will be used at the stage of backtracking in the next segment query (line 19). After all query trajectories have been processed, we update the query segments in  $\mathcal{S}$  to their corresponding non-time-interval-split segments (line 21). Then, we group the pairs in  $\mathcal{S}$  by the client that generated  $s^{\#}$ , and send the pairs to the clients based on the client IDs of the segments (line 22). STS-Join verifies the trajectory similarity on the clients, since the server only stores simplified-and-obfuscated trajectories. After each client has computed the trajectory similarity, it returns the result set to the server.

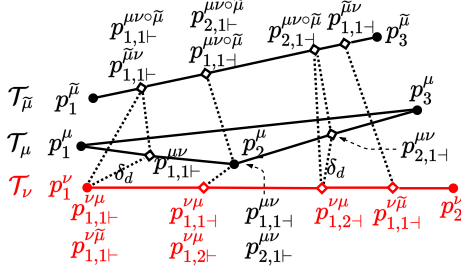
**On client side.** On each client, the trajectories corresponding to the segments  $s^{\#}$  received from the server are retrieved (by ID lookups using the local trajectory IDs of the segments). Then, we compute the CDD of the segment pairs from the server and add up the CDDs for every trajectory pair formed by the segment pairs. The pairs satisfying the time threshold  $\delta_t$  are returned as the query result to the server. This guarantees no false positives.

**Discussion.** Backtracking is not limited to quasi-quadtrees. It is applicable to all space-partitioning indices in which the process can stop at a common parent node. A query starting from such parent nodes will not have false negatives, since there are no overlaps among MBRs on the same level in the tree. That means one location can be covered only by one MBR at each level. In addition, backtracking can be applied in index construction, insertion, and deletion, because segments are inserted or deleted sequentially, and we can leverage the common parent node approach to reduce the node accesses when operating on the next segment.

## 5.2 CDDS-Based Pruning

Algorithm 1 sends all segment pairs that may satisfy the query distance threshold  $\delta_d$  to the clients for further verification and CDDS computation. In this subsection, we compute an upper bound of the actual CDDS between a query trajectory  $\mathcal{T}_v$  and a known trajectory  $\mathcal{T}_\mu$  using  $\mathcal{T}_v$  and the simplified-and-obfuscated segments of  $\mathcal{T}_\mu$  stored on the server. We only send the segment pairs to the corresponding client when the upper bound exceeds  $\delta_t$ , thus saving both communication costs between the server and the clients and computation costs on the clients. This essentially adds a subprocedure to prune the segment pairs before Line 22 of Algorithm 1 using an upper bound of  $cdds(\mathcal{T}_\mu, \mathcal{T}_v)$ . For simplicity, we only describe the CDDS-based pruning procedure below but do not repeat the full pseudocode of the STS-Join algorithm powered by it.

**CDDS-based pruning procedure.** We group the segment pairs that satisfy the query distance threshold (i.e., the segment pairs in set  $\mathcal{S}$  at Line 21 of Algorithm 1) by their client IDs, local trajectory IDs, and query trajectory IDs. The segment pairs that share the same client ID, local trajectory ID, and query trajectory ID all come from the same pair of simplified-and-obfuscated known and query trajectories ( $\mathcal{T}_\mu, \mathcal{T}_v$ ) for CDDS upper bound computation. Let the set of such segment pairs be  $\mathcal{S}_{\mu, v}$  and the original known trajectory corresponding to  $\mathcal{T}_\mu$  be  $\mathcal{T}_\mu$ , respectively. Then, we compute an upper bound of the close-distance duration (i.e., CDD) for every pair of segments in  $\mathcal{S}_{\mu, v}$ . Summing up these upper bounds for all segment pairs in  $\mathcal{S}_{\mu, v}$  yields our upper bound of  $cdds(\mathcal{T}_\mu, \mathcal{T}_v)$ , since these pairs are the only ones in  $\mathcal{T}_\mu \times \mathcal{T}_v$  (and hence  $\mathcal{T}_\mu \times \mathcal{T}_v$ ) that satisfy the query distance threshold by definition. The set  $\mathcal{S}_{\mu, v}$  is pruned



**Figure 6: Example of trajectory similarity bounding**

from being sent to a client if the CDDs upper bound is less than  $\delta_t$ . Next, we detail our CDD upper bound computation.

**CDD upper bound.** Recall that the CDD between a known segment  $s_i^\mu$  of  $\mathcal{T}_\mu$  and a query segment  $s_j^\nu$  of  $\mathcal{T}_\nu$ ,  $cdd(s_i^\mu, s_j^\nu)$ , is the time duration when their spatial distance is within  $\delta_d$ . We further define the CDD between a simplified-and-obfuscated segment  $s_i^\mu$  and a query segment  $s_j^\nu$  as the time duration where their spatial distance is within  $\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob}$ , denoted as  $\overline{cdd}(s_i^\mu, s_j^\nu)$ . Then, for all known segments  $s_i^\mu, s_{i+1}^\mu, \dots, s_{i+\Delta}^\mu$  corresponding to  $s_i^\mu, \dots, s_{i+\Delta}^\mu$ ,  $cdd(s_{i+k}^\mu, s_j^\nu) \leq \overline{cdd}(s_i^\mu, s_j^\nu)$ . This is because, given a point on a known segment, its corresponding point on the simplified-and-obfuscated segment is within a distance of  $\theta_{sp} + \sqrt{2}\theta_{ob}$  by definition. Thus, if the distance between (points on)  $s_{i+k}^\mu$  and  $s_j^\nu$  is within  $\delta_d$ , the distance between (points on)  $s_i^\mu$  and  $s_j^\nu$  must be within  $\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob}$ . Therefore, we use  $\overline{cdd}(s_i^\mu, s_j^\nu)$  as our CDD upper bound of the original known trajectories.

We formulate the CDD upper bound and show its correctness with the following lemma.

**LEMMA 5.1.** *Given a query segment  $s_j^\nu \in \mathcal{T}_\nu$ , a sequence of known segments  $s_i^\mu, s_{i+1}^\mu, \dots, s_{i+\Delta}^\mu \in \mathcal{T}_\mu$ , and their corresponding simplified-and-obfuscated segment  $s_i^\mu \in \mathcal{T}_\mu$ , we have:*

$$\bigcap_{k=0}^{\Delta} cdd(s_{i+k}^\mu, s_j^\nu) \leq \overline{cdd}(s_i^\mu, s_j^\nu) \quad (9)$$

where  $\overline{cdd}$  is the CDD with distance threshold  $\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob}$ .

**PROOF.** We use Fig. 6 to help illustrate the proof, where known segments  $s_i^\mu$  can be  $p_1^\mu, p_2^\mu$  and  $p_2^\mu, p_3^\mu$  correspond to a simplified-and-obfuscated segment  $s_i^\mu = p_1^\mu, p_3^\mu$ , and the query segment is shown as  $s_j^\nu = p_1^\nu, p_2^\nu$ . Besides, any two points connected by a dash line have the same timestamp.

Given a known segment  $s_{i+k}^\mu$  and a query segment  $s_j^\nu$ , their CDD can be non-zero only if they overlap in their time span, i.e.,  $ot(s_{i+k}^\mu, s_j^\nu) > \emptyset$ . Further, if the CDD is non-zero, it is defined by the two roots (i.e., two time points) of the quadratic equation  $dist^2(s_{i+k}^\mu, s_j^\nu) = \delta_d^2$  (cf. Equation 3). Let the two roots be  $t_{j,i+k}^{\nu\mu}$  and  $t_{j,i+k}^{\mu\nu}$ . Then,  $cdd(s_{i+k}^\mu, s_j^\nu) = t_{j,i+k}^{\nu\mu} - t_{j,i+k}^{\mu\nu}$ . Note that, if either  $t_{j,i+k}^{\nu\mu}$  or  $t_{j,i+k}^{\mu\nu}$  is outside the range of  $ot(s_{i+k}^\mu, s_j^\nu)$ , we replace it with the corresponding boundary value of  $ot(s_{i+k}^\mu, s_j^\nu)$  to meet the overlapping time span requirement of the segments. Let the points

corresponding to  $t_{j,i+k}^{\nu\mu}$  on  $s_{i+k}^\mu$  and  $s_j^\nu$  be  $p_{i+k,j}^{\mu\nu}$  and  $p_{j,i+k}^{\nu\mu}$ , respectively. Similarly, let the points corresponding to  $t_{j,i+k}^{\mu\nu}$  on the two segments be  $p_{i+k,j}^{\nu\mu}$  and  $p_{i+k,j}^{\mu\nu}$ , respectively. In Fig. 6, these four points on  $s_{i+k}^\mu$  and  $s_j^\nu$  are  $p_{1,1}^\mu, p_{1,1}^{\nu\mu}, p_{1,1}^{\mu\nu}$ , and  $p_{1,1}^{\nu\mu}$  ( $i = 1, j = 1$ , and  $k = 0$ ).

Then, we analyze the distance between a known segment  $s_{i+k}^\mu$  and its corresponding simplified-and-obfuscated segment  $s_i^\mu$ . Note that a sequence of known segments can be simplified into a single segment with a simplification threshold  $\theta_{sp}$ , while the simplified segment is further obfuscated by shifting the endpoints with a maximum shifting distance  $\theta_{ob}$  along each dimension. Thus, the distance between any point  $p$  on  $s_{i+k}^\mu$  and its corresponding point (i.e., the point at the same time  $t$  as that of  $p$ ) on  $s_i^\mu$ , is bounded within  $\theta_{sp} + \sqrt{2}\theta_{ob}$ . This applies to the distance between points (at any given time  $t$ ) on  $p_1^\mu, p_2^\mu$  ( $p_2^\mu, p_3^\mu$ ) and  $p_1^\mu, p_3^\mu$  in Fig. 6.

Next, we derive the distance between the query segment and the simplified-and-obfuscated segment by leveraging the distance relationship above. By definition, the distance between the known segment  $s_{i+k}^\mu$  and the query segment  $s_j^\nu$  does not exceed  $\delta_d$  when  $t \in [t_{j,i+k}^{\nu\mu}, t_{j,i+k}^{\mu\nu}]$ , while the distance between  $s_{i+k}^\mu$  and its corresponding simplified-and-obfuscated segment  $s_i^\mu$  does not exceed  $\theta_{sp} + \sqrt{2}\theta_{ob}$ . Therefore, the distance  $dist(s_{i+k}^\mu, s_j^\nu)$  between the query segment and the simplified-and-obfuscated known segment does not exceed  $\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob}$ , when  $t \in [t_{j,i+k}^{\nu\mu}, t_{j,i+k}^{\mu\nu}]$ . In Fig. 6, we locate a point  $p_{2,1}^{\mu\nu}$  on  $p_1^\mu, p_3^\mu$  whose timestamp is  $t_{1,2}^{\nu\mu}$  ( $i = 1, j = 1$ , and  $k = 1$ ). Then, the distance between  $p_{2,1}^{\mu\nu}$  and  $p_{2,1}^{\nu\mu}$  does not exceed  $\theta_{sp} + \sqrt{2}\theta_{ob}$ , while the distance between  $p_{2,1}^{\mu\nu}$  and  $p_{1,2}^{\nu\mu}$  is  $\delta_d$  as shown above. Thus, the distance between  $p_{2,1}^{\mu\nu}$  and  $p_{1,2}^{\nu\mu}$  is within  $\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob}$ .

Since there exists a time range  $[t_{j,i+k}^{\nu\mu}, t_{j,i+k}^{\mu\nu}]$  where the distance between  $s_{i+k}^\mu$  and  $s_j^\nu$  does not exceed  $\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob}$ , the quadratic distance equation  $dist^2(s_{i+k}^\mu, s_j^\nu) = (\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob})^2$  must have two roots, which are denoted as  $t_{j,i+k}^{\nu\mu}$  and  $t_{j,i+k}^{\mu\nu}$ , respectively. Also, since this quadratic function has a non-negative quadratic coefficient (cf. Equation 3), and the distance does not exceed  $\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob}$  when  $t \in [t_{j,i+k}^{\nu\mu}, t_{j,i+k}^{\mu\nu}]$ , we can derive that  $t_{j,i+k}^{\nu\mu} \leq t_{j,i+k}^{\mu\nu} < t_{j,i+k}^{\nu\mu} \leq t_{j,i+k}^{\mu\nu}$ . Since  $cdd(s_{i+k}^\mu, s_j^\nu) = t_{j,i+k}^{\nu\mu} - t_{j,i+k}^{\mu\nu}$  and  $\overline{cdd}(s_{i+k}^\mu, s_j^\nu) = t_{j,i+k}^{\nu\mu} - t_{j,i+k}^{\mu\nu}$ , we have  $cdd(s_{i+k}^\mu, s_j^\nu) \leq \overline{cdd}(s_{i+k}^\mu, s_j^\nu)$  when  $t \in ot(s_{i+k}^\mu, s_j^\nu)$ . Recall that  $\overline{cdd}$  is CDD computed with distance threshold  $\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob}$ . In Fig. 6,  $t_{1,1}^{\nu\mu}$  is the upper boundary of  $\overline{cdd}(p_1^\mu, p_3^\mu, p_1^\nu, p_2^\nu)$ , and the corresponding points on the query segment and the simplified-and-obfuscated segment are  $p_{1,1}^{\nu\mu}$  and  $p_{1,1}^{\mu\nu}$  ( $i = 1, j = 1$ , and  $k = 1$ ) respectively. Meanwhile, when  $t = t_{1,2}^{\nu\mu}$ , the distance between  $p_{1,2}^{\nu\mu}$  and  $p_{1,2}^{\mu\nu}$

does not exceed  $\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob}$ . Thus, we have  $t_{1,2+}^{v\mu} \leq t_{1,1+}^{v\mu}$ . Such an inequality is also applicable to other CDD boundaries.

Finally, we sum up the CDD of each segment in  $s_i^\mu, s_{i+1}^\mu, \dots, s_{i+\Delta}^\mu$  with  $s_j^v$ , and we sum up  $\overline{cdd}(s_i^\mu, s_j^v)$  that corresponds to different time ranges in  $ot(s_{i+k}^\mu, s_j^v)$  where  $0 \leq k \leq \Delta$ . Since the inequality is satisfied on each separate time range, such inequality is also satisfied for the sums, i.e.,  $\overline{cdd}(p_1^\mu, p_2^\mu, p_1^v, p_2^v) + \overline{cdd}(p_2^\mu, p_3^\mu, p_1^v, p_2^v) \leq \overline{cdd}(p_1^\mu, p_3^\mu, p_1^v, p_2^v)$  in Fig. 6. This completes the proof.

Given Lemma 5.1, we have the following lemma to bound the CDDS for the original known trajectories.

LEMMA 5.2. *Given a query trajectory  $\mathcal{T}_v$ , a known trajectory  $\mathcal{T}_\mu$ , and its corresponding simplified-and-obfuscated version  $\mathcal{T}_{\mathcal{P}}$ , we have:*

$$\overline{cdds}(\mathcal{T}_\mu, \mathcal{T}_v) \leq \overline{cdds}(\mathcal{T}_{\mathcal{P}}, \mathcal{T}_v) \quad (10)$$

where  $\overline{cdds}$  is CDDS with distance threshold  $\delta_d + \theta_{sp} + \sqrt{2}\theta_{ob}$ .

PROOF. The proof is straightforward and hence is omitted.

### 5.3 Cost Analysis

We analyze the cost of Algorithm 1 in terms of the number of node accesses. If a query segment spans the whole spatio-temporal space (worst case), all index nodes may be visited, with or without backtracking. However, this rarely happens. In many cities, points of interest are distributed in a polycentric structure [7, 12]. Trajectories are likely to gather near local centers, where backtracking helps query the sub-spatial index of the centers.

The algorithm iterates through the query segments. In each iteration, the costs are spent on finding the common parent node of two adjacent points (and hence adjacent segments) in the query trajectory and on traversing the index to reach leaf nodes.

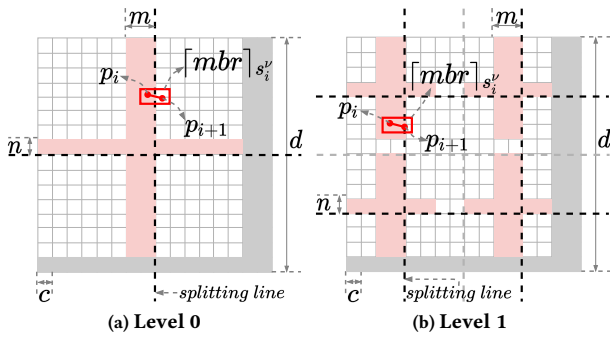


Figure 7: Space division in the index

First, we derive the backtracking distance  $bd$  (i.e., the number of tree levels) between the leaf nodes and the common parent node of two adjacent points. We use the expectation  $E(bd)$  to measure the average cost, which is obtained by summing up different backtracking lengths weighted by their probabilities.

We use Fig. 7 to illustrate how to derive  $E(bd)$  with a space division in the quasi-quadtrees. The width of the expanded MBR covering the (red) query segment is  $m$  and the height is  $n$ , where  $m > n$ . Let the side length of the space be  $d$  and that of a cell be  $c$ . Here, each cell corresponds to a leaf node. In our quasi-quadtrees, once the node capacity is exceeded, a node will be divided into

four sub-nodes. The space covered by the node is split into four quadrants. We call the vertical and horizontal boundaries (the black dash lines) between the sub-spaces the *splitting lines*.

Different splitting lines divide the space at different index levels. The backtracking distance can be derived by the levels of the splitting lines that intersect with the expanded MBR. To derive whether an expanded MBR intersects with a splitting line, we compare the location of the upper left vertex of the MBR with the splitting line. In Fig. 7a, when the upper left vertex of the MBR falls in the pink area, the query segment intersects with a splitting line at level 0. The common parent node is the root of the quasi-quadtrees, and  $bd$  is the height  $h$  of this tree, i.e.,  $\lfloor \log_2 \frac{d}{c} \rfloor$ . Fig. 7b shows the level-1 splitting lines, where the backtracking distance is  $h - 1$  if the expanded MBR intersects with them.

Next, we consider the stop condition of segment intersection. When  $m$  and  $n$  are large, a query segment may intersect with splitting lines at multiple lower levels, while common parent nodes cannot be at such levels. Let the lowest feasible level of the common parent node be  $I$ ,  $I = \lfloor \log_2 \frac{d}{2m} \rfloor$ . This is because, given a pivot of a query segment in a node, we need at most a common parent node with an MBR side length of  $2m$  to also cover the other pivot of the segment. In Fig. 7b, if  $m > 4c$  and  $n > 4c$ , the query segment intersects with splitting lines at levels 1 to 4. We only need to consider the case where the query segment intersects with the splitting lines at level 0.

The probabilities of different backtracking lengths are derived by the ratio of the data space occupied by the pink area of different levels. We cut off the grey area in Fig. 7, as the query segment is outside the space when its expanded MBR is in this area.

Then, we compute  $E(bd)$  by summing up the product of the probability and the backtracking distance at every level:

$$E(bd) = \sum_{i=0}^{h-1} \frac{[m(\frac{d}{2^i} - n) + n(\frac{d}{2^i} - m) - mn] \cdot 4^i}{(b-m)(b-n)} \cdot (h-i) \quad (11)$$

where the fraction part is the probability determined by the size of the pink area, and  $h - i$  is the distance between a leaf node and the common parent node. We let  $\Lambda = \frac{1}{(b-m)(b-n)}$  and then expand Equation 11:

$$\begin{aligned} E(bd) &= \Lambda \cdot \sum_{i=0}^{h-1} [m(\frac{d}{2^i} - n) + n(\frac{d}{2^i} - m) - mn] \cdot 4^i \cdot (h-i) \\ &= \Lambda \cdot h \sum_{i=0}^{h-1} (2^i md + 2^i nd - 3 \cdot 4^i mn) - \sum_{i=0}^{h-1} (2^i mdi + 2^i ndi \\ &\quad - 3 \cdot 4^i mni) \end{aligned} \quad (12)$$

Then, we expand the two terms in square brackets in Equation 12 separately. For the first term, by summing up the geometric series, we can have:

$$\begin{aligned} &h \sum_{i=0}^{h-1} (2^i md + 2^i nd - 3 \cdot 4^i mn) \\ &= h[md(2^{I+1} - 1) + nd(2^{I+1} - 1) - mn(4^{I+1} - 1)] \\ &= h[md(\frac{d}{m} - 1) + nd(\frac{d}{m} - 1) - mn(\frac{d^2}{m^2} - 1)] \\ &= h(d-m)(d-n) \end{aligned} \quad (13)$$

The second term can be expanded as:



$$\begin{aligned}
& \sum_{i=0}^{I-1} (2^i m d i + 2^i n d i - 3 \cdot 4^i m n i) \\
&= 2I(d-m)(d-n) - (I+1) \sum_{i=1}^I (2^i m d + 2^i n d - 3 \cdot 4^i m n) \quad (14) \\
&= 2I(d-m)(d-n) - (I+1)(d-2m)(d-2n)
\end{aligned}$$

Lastly, we plug Equation 13, 14 into Equation 12:

$$\begin{aligned}
E(bd) &= \Lambda h(d-m)(d-n) - 2I(d-m)(d-n) \\
&\quad + (I+1)(d-2m)(d-2n) \\
&= h - 2I + (I+1) \frac{(d-2m)(d-2n)}{(d-m)(d-n)} \quad (15) \\
&6 h - I + 1
\end{aligned}$$

So far, we have that  $E(bd)$  is up to  $h - I + 1$ .

**Total cost.** To query a trajectory, STS-Join iteratively takes  $O((h-I)|\tilde{T}|)$  node accesses for backtracking, since only one node is visited at each level, where  $|\tilde{T}|$  is the number of segments in a query trajectory. Besides, it visits  $O(4^{h-I}|\tilde{T}|)$  nodes while searching for similar segments of a query segment. In total, the cost of STS-Join with backtracking is  $O(4^{h-I}|\tilde{T}||\mathcal{D}_q|)$ .

## 6 EXPERIMENTS

### 6.1 Experimental Setup

All experiments are conducted on a 64-bit machine running Ubuntu 20.04 LTS with a 6-core AMD Ryzen 5 CPU, 32 GB RAM, and a 500GB SSD. All algorithms are implemented in C++ and run in main memory.

**Datasets.** To the best of our knowledge, there is no public contact-tracing datasets. We use two real trajectory datasets instead: **DiDi** [2] and **GeoLife** [1]. DiDi contains vehicle trajectories in Chengdu, a capital city in China. We randomly sample trajectories recorded in the first week of November, 2016 to generate a dataset with 500k trajectories and 67.7 million segments (this is limited by our memory size). The trajectories come from 201k users (each is considered a client in the experiments), i.e., there are 2.5 trajectories per client on average. There are 135 segments per trajectory on average. The average length and time span of a segment are 30.0 meters and 4.5 seconds, respectively.

GeoLife contains trajectories of different transport modes (e.g., walking, driving, and cycling) recorded mainly in Beijing, China. We only keep the trajectories within the Fifth Ring Road of Beijing. There are some 11k trajectories in this area and very few outside. We randomly sample 10k trajectories from them to form the GeoLife dataset used in the experiments. We shift the trajectory times into the same week without changing their time span to increase the density in time. The trajectories come from 182 users (clients), i.e., there are 55 trajectories per client on average. There are 9.0 million segments in these trajectories and 896 segments per trajectory on average. The average length and time span of a segment are 12.2 meters and 9.2 seconds, respectively. The two datasets are summarized in Table 2. We further generate their random subsets for scalability tests. See Table 3 for the sizes of the subsets.

**Algorithms.** There is no existing algorithm for STS-Join. We adapt two techniques which result in three baseline algorithms. Meanwhile, we replace the join predicates in these algorithms with ours, so *all algorithms return the same accurate query results*. We test three variants of our proposed STS-Join algorithm to confirm the

**Table 2: Datasets**

Dataset	Configuration			
	# traj.	# seg.	Avg. seg. length	Avg. seg. time span
DiDi	500k	67,654k	30.0 meters	4.5 seconds
GeoLife	10k	8,981k	12.2 meters	9.2 seconds

effectiveness of the proposed backtracking-based query algorithm and the CDDS-based pruning strategy.

- **3DR:** The 3DR-tree [21] extends the R-tree by adding time ranges as an extra dimension. We use it to index segments in  $\mathcal{D}_p$  and run an R-tree join like algorithm for queries.
- **ALSTJ:** As described in Section 2, ALSTJ [3] (an R-tree variant) is the closest study to ours, which also considers sub-trajectory join. We replace the join predicate in its query algorithm with ours to process STS-Join.
- **ALSTJ-T:** The original ALSTJ index does not consider the time dimension. We improve it by adding a B-tree-like temporal index at its top levels like our STS-Index.
- **STS:** STS denotes the naive STS-Join without backtracking or CDDS-based pruning as described at the start of Section 5.1.
- **STS-BT:** STS-BT is the backtracking-based STS-Join (no pruning), described in Algorithm 1.
- **STS-BTB:** STS-BTB further prunes matched trajectory segment pairs from being sent to the clients by a CDDS upper bound as described in Section 5.2.

We set the node capacity of the baselines and the B-trees in STS-Index to 100. We set the B-tree and the bitmap time interval lengths in STS-Index to 1,800 and 30 seconds, respectively. The client-server model of STS-Join is simulated, since it aims to enhance privacy but not computation capacity.

**Parameter settings.** Table 3 summarizes parameters tested in the experiments, where the default values are in bold.

We follow the closest work [3] to compare the algorithm response times, and we also measure the I/O cost that is a widely used indicator in spatial indexing studies [11, 22, 23] if using external memory based implementation. Besides, we study the communication cost by recording the number of segments sent to the clients for final verification and the number of clients receiving such segments. *All algorithms return the same result set under the same query setting, which verifies the correctness of our methods.*

**Table 3: Parameters and Their Values**

Parameter	Settings
$ \mathcal{D}_p $ (DiDi)	100k, <b>200k</b> , 300k, 400k, 500k
$ \mathcal{D}_p $ (GeoLife)	2k, <b>4k</b> , 6k, 8k, 10k
$ \mathcal{D}_q $ (DiDi)	100, <b>200</b> , 300, 400, 500
$ \mathcal{D}_q $ (GeoLife)	2, <b>4</b> , 6, 8, 10
$\delta_d$	10, <b>20</b> , 30, 40, 50 (meters)
$\delta_t$	0, 10, <b>20</b> , 30, 40, 50 (seconds)
$\theta_{sp}$	0, 10, <b>20</b> , 30, 40, 50 (meters)
$\theta_{ob}$	0, 10, <b>20</b> , 30, 40, 50 (meters)

### 6.2 Join Performance

We set the default sizes of the known DiDi and GeoLife datasets  $\mathcal{D}_p$  to 200k and 4k, respectively. The query datasets  $\mathcal{D}_q$  are randomly sampled from the original DiDi and GeoLife datasets with the same time ranges as the known datasets  $\mathcal{D}_p$ . Their sizes are set to be

proportional to those of the respective known DiDi and GeoLife datasets. The default sizes of DiDi and GeoLife query datasets are 200 and 4 (i.e., 0.1% of  $|\mathcal{D}_p|$ ), respectively. The GeoLife query datasets might seem small, but there are some 900 segments per trajectory on average, which are sufficient to show the performance difference of the algorithms.

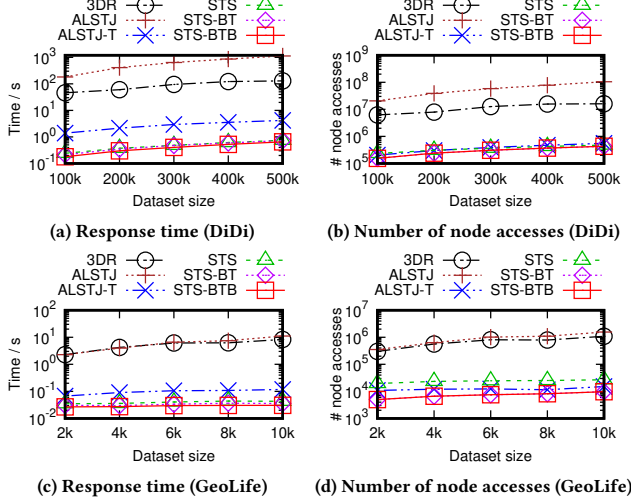


Figure 8: Query costs vs. known dataset size

**Varying the known dataset size  $|\mathcal{D}_p|$ .** Fig. 8 shows the query costs, which increase with  $|\mathcal{D}_p|$ . Our STS-BTB algorithm with both backtracking and CDDS-based pruning outperforms all competitors consistently in response time. It also takes the fewest node accesses among all algorithms except our STS-BT, which has the same node accesses as STS-BTB since the CDDS-based pruning does not impact node accesses. The advantage is up to three orders of magnitude. Comparing STS-BTB with STS-BT and STS, it achieves  $\sim 14\%$  and  $\sim 18\%$  lower running times, which confirms the effectiveness of our CDDS-based pruning to reduce the computational costs. STS has more node accesses ( $\sim 24\%$  on average) than STS-BT, confirming the effectiveness of our backtracking strategy to reduce node accesses. STS has very similar running times to those of STS-BT, as the algorithms are running in memory, where the time saved on node accesses is lost in the more complex backtracking procedure of STS-BT. STS also outperforms the baselines in running time, while it has slightly more node accesses than ALSTJ-T when  $|\mathcal{D}_p| \geq 200k$  on DiDi datasets. This is because ALSTJ-T has a much larger node capacity in its index, i.e., 100, while we use 4 in our quasi-quadtrees. The response time of ALSTJ-T is still up to 6.7 times larger than those of STS and STS-BT. This is because of its worse pruning power – its index allows overlapping node MBRs while ours does not. ALSTJ is even worse, i.e., up to three orders of magnitude slower and two orders of magnitude more node accesses than STS-BTB, because of its lack of pruning power in the time dimension. 3DR considers the time dimension but is not optimized for trajectories. It runs up to 240 times slower and has up to 43 times more node accesses than STS-BTB.

**Varying the query dataset size  $|\mathcal{D}_q|$ .** Next, we vary the query dataset size  $|\mathcal{D}_q|$ . Fig. 9 shows that the query costs increase with  $|\mathcal{D}_q|$ , which is expected. The relative performance among the algorithms is very similar to that when varying  $|\mathcal{D}_p|$ . STS, STS-BT, and

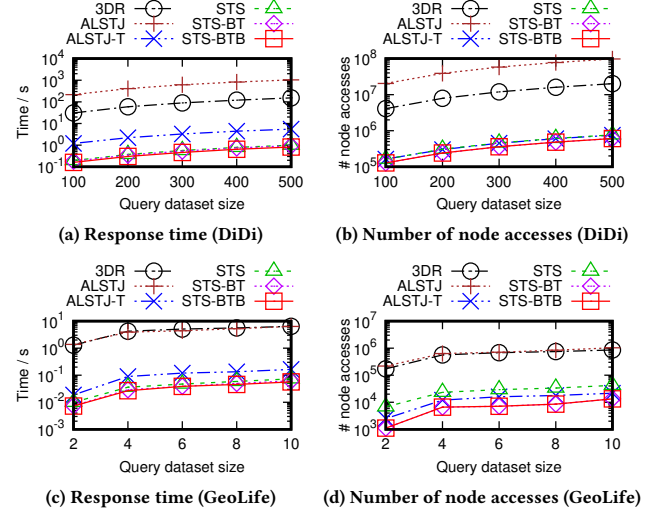


Figure 9: Query costs vs. query dataset size

STS-BTB outperform all competitors on both datasets in response time (up to three orders of magnitude), while ALSTJ-T has a slightly smaller number of node accesses than that of STS. These confirm the superiority of the proposed algorithms and the effectiveness of the proposed backtracking and CDDS-based pruning techniques to reduce the algorithm costs.

Since the query performance patterns are similar on both DiDi and GeoLife data, we omit the figures for GeoLife in the following query experiments for conciseness.

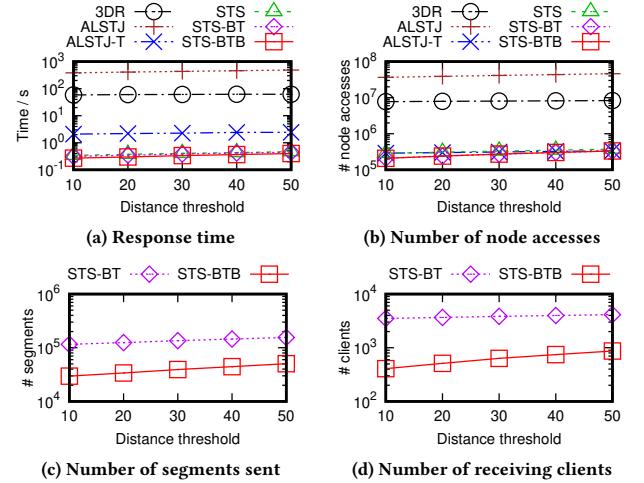


Figure 10: Query costs vs. query distance threshold

**Varying the query distance threshold  $\delta_d$ .** Fig. 10 shows the algorithm performance as  $\delta_d$  increases from 10 to 50 meters. The query costs of all algorithms increase with  $\delta_d$ , because more trajectory pairs satisfy a larger  $\delta_d$  and need to be returned. The relative performance between the proposed STS algorithms and the competitors are similar to those observed above. Like before, STS-BTB is up to 22% and 14% faster than STS and STS-BT, respectively, while STS-BT and STS-BTB both have the fewest node accesses. We also study the impact of  $\delta_d$  on the communication cost, which is reflected by the number of segments that are sent to the clients and the number of clients receiving at least one segment, as reported in

Figs. 10c and 10d, respectively. We see that, while the communication costs of STS-BTB increase with  $\delta_d$  which is expected, they are at least 3.1 and 4.8 times lower than those of STS-BT (for  $\delta_d \leq 50$ ) in terms of the number of segments sent and the number of clients receiving the segments, respectively. Here, STS has been omitted since it has the same communication costs as STS-BT, and the same applies to the figures below. The results further confirm the effectiveness of our CDDS-based pruning.

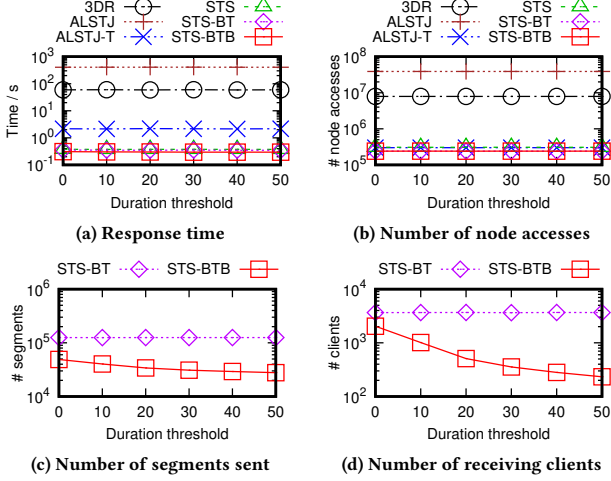


Figure 11: Query costs vs. query duration threshold

**Varying the query duration threshold  $\delta_t$ .** Figs. 11a and 11b show that the join query costs are relatively stable as the query duration threshold  $\delta_t$  increases from 0 to 50 seconds. This is because all algorithms compare the CDDS of trajectories with  $\delta_t$  at almost their final stages. The value of  $\delta_t$  has little impact on the time and node access costs. The relative performance among the algorithms is again very similar to that in the previous figures. From Figs. 11c and 11d, we see that the communication costs of STS-BT is relatively stable, while those of STS-BTB decrease quickly as  $\delta_t$  grows. This is because a larger  $\delta_t$  makes the query predicate more difficult to be satisfied, even for the CDDS used in pruning which has been inflated to guarantee no false negatives. This indicates that our CDDS-based pruning strategy may gain more pruning power in application scenarios with larger  $\delta_t$  values.

**Varying the simplification threshold  $\theta_{sp}$ .** We also study the impact of the simplification threshold on STS-Join. As shown in Figs. 12a and 12b, when  $\theta_{sp}$  increases from 0 to 50 meters, the query costs first drop and then increase for all STS algorithm variants. The drop is because introducing simplification helps reduce the number of segments in STS-Index and hence the query costs. However, as  $\theta_{sp}$  grows, the data segments become longer which leads to larger node MBRs (and greater MBR expansion to ensure no false negatives) and reduced the index pruning power. This explains for the increase in query costs. STS-BTB and STS-BT have the same node accesses since the CDDS-based pruning has no impact on node accesses. Meanwhile, STS-BTB again outperforms STS-BT in the communication costs, thus showing the robustness of STS-BTB (cf. Figs. 12c and 12d).

**Varying the shifting distance  $\theta_{ob}$ .** Fig. 13 shows that the query costs of the STS algorithm variants increase with  $\theta_{ob}$ , as

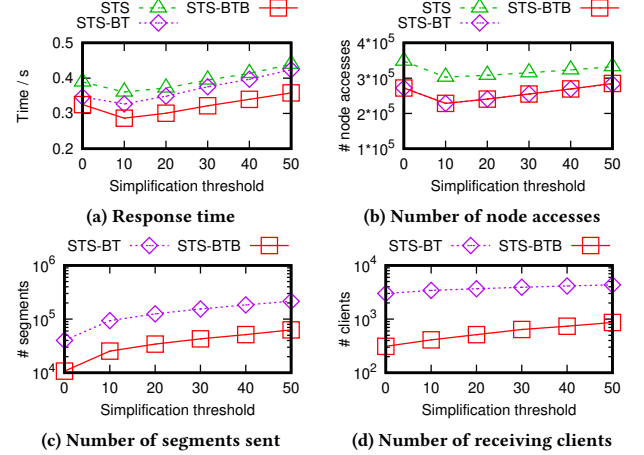


Figure 12: Query costs vs. simplification threshold

expected. A larger  $\theta_{ob}$  may enlarge the node MBRs and hence leads to higher query costs. STS-BTB still takes the smallest time costs, number of node accesses, and communication costs, which have similar trends to those when varying  $\theta_{sp}$ , further confirming the robustness of the algorithm and the optimization techniques.

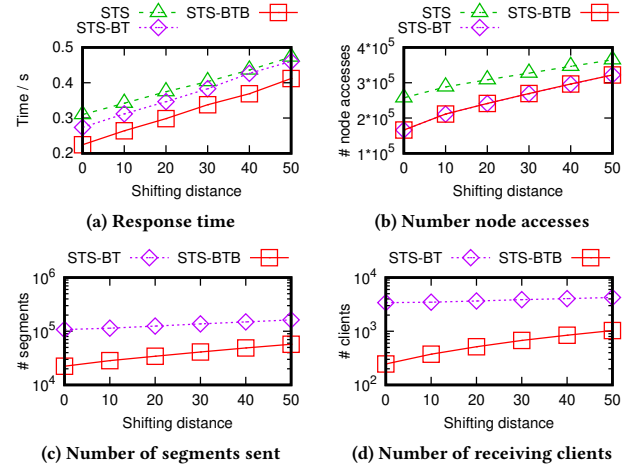


Figure 13: Query costs vs. shifting distance

When other STS-Index parameters are varied (the B-tree, and the bitmap time interval lengths), STS-BTB consistently outperforms STS-BT and STS. We omit the figures for conciseness.

### 6.3 Index Construction Performance

We show the costs of index construction in Fig. 14. Since STS-BT and STS-BTB have the same process for index construction, we omit STS-BTB from the figure. STS-BT index is the fastest to build, which takes up to 15% less time than ALSTJ-T. STS-BT index also takes fewer node accesses to build than ALSTJ-T index does except when  $|\mathcal{D}_p| > 400k$  (recall that our index has smaller node capacities). STS takes more time and node accesses to build without backtracking, while ALSTJ suffers in its lack of pruning power in the time dimension for data insertion. 3DR takes the largest costs to build, for that it is not optimized for trajectories.

We report the index sizes (RAM space) in Fig. 14c. STS and STS-BT indices have the same size, so we only report for STS. STS,

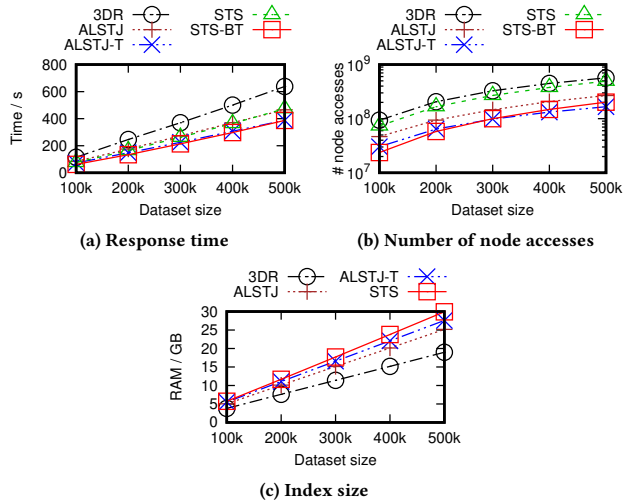


Figure 14: Index construction costs vs. dataset size

ALSTJ, and ALSTJ-T indices are larger than those of 3DR, because they store simplified and original trajectories while 3DR only stores original trajectories. STS indices take slightly (6.85%) more space than ALSTJ-T does. This extra cost comes from the bitmap structure and more index nodes, which is justified by our much faster query performance. ALSTJ indices are smaller than ALSTJ-T's as they do not index the time dimension.

**Effectiveness of the join predicate.** To show the differences between join queries, we run the original ALSTJ algorithm (with its own join predicate [3]) and collect the resultant trajectory pairs. We find that this misses 36% of the trajectory pairs returned by STS-Join, confirming the effectiveness of STS-Join in identifying trajectory pairs that may be missed by existing queries.

## 7 CONCLUSIONS

We propose a trajectory join query named STS-Join that returns pairs of trajectories within a given distance threshold for a time longer than a given time threshold. Unlike existing join queries that require the full trajectories to be similar, STS-Join can find trajectories that are similar for a time period while being quite different overall. We propose a client-server index named STS-Index and an efficient backtracking-based algorithm and a similarity time period based pruning strategy to support queries while supporting privacy protection. Our cost analysis and experiments on real data confirm the superiority of the proposed algorithm – it outperforms adapted state-of-the-art join algorithms consistently and by up to three orders of magnitude in the query time.

For future work, we plan to extend STS-Join to the road network, and allow temporal shifts in the similarity metric. We also intend to further study the privacy protection on trajectory similarity queries from a theoretical perspective.

## ACKNOWLEDGMENTS

This work was sponsored in part by the Australian Research Council under the Discovery Project Grant DP180103332 and by the NSF under Grants IIS-18-16889 and IIS-20-41415.

## REFERENCES

- [1] 2012. GeoLife. [www.microsoft.com/en-us/download/details.aspx?id=52367](http://www.microsoft.com/en-us/download/details.aspx?id=52367).
- [2] 2016. DiDi. [outreach.didichuxing.com/research/opendata/en/](http://outreach.didichuxing.com/research/opendata/en/).
- [3] Omid Isfahani Alamdari, Mirco Nanni, Roberto Trasarti, and Dino Pedreschi. 2020. Towards in-memory sub-trajectory similarity search. In *EDBT/ICDT Workshops*.

- [4] Miguel E. Andrés, Nicolás E. Bordenabe, Konstantinos Chatzikokolakis, and Catuscia Palamidessi. 2013. Geo-indistinguishability: Differential privacy for location-based systems. In *CCS*. 901–914.
- [5] Petko Bakalov and Vassilis J Tsotras. 2006. Continuous spatiotemporal trajectory joins. In *International Conference on GeoSensor Networks*. 109–128.
- [6] Dan Buzan, Stan Sclaroff, and George Kollios. 2004. Extraction and clustering of motion trajectories in video. In *International Conference on Pattern Recognition*. 521–524.
- [7] Jixuan Cai, Bo Huang, and Yimeng Song. 2017. Using multi-source geospatial big data to identify the structure of polycentric cities. *Remote Sensing of Environment* 202 (2017), 210–221.
- [8] Lei Chen, M Tamer Özsu, and Vincent Oria. 2005. Robust and fast similarity search for moving object trajectories. In *SIGMOD*. 491–502.
- [9] Chi-Yin Chow and Mohamed F Mokbel. 2007. Enabling private continuous queries for revealed user locations. In *SSTD*. 258–275.
- [10] Chi-Yin Chow and Mohamed F. Mokbel. 2011. Trajectory privacy in location-based services and data publication. *SIGKDD Explorations Newsletter* 13, 1 (2011), 19–29.
- [11] Douglas Comer. 1979. Ubiquitous B-tree. *Comput. Surveys* 11, 2 (1979), 121–137.
- [12] Yue Deng, Jiping Liu, Yang Liu, and An Luo. 2019. Detecting urban polycentric structure from POI data. *ISPRS International Journal of Geo-Information* 8, 6 (2019), 283.
- [13] David H. Douglas and Thomas K. Peucker. 1973. Algorithms for the reduction of the number of points required to represent a digitized line or its caricature. *Cartographica: The International Journal for Geographic Information and Geovisualization* 10, 2 (1973), 112–122.
- [14] Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. 2006. Our data, ourselves: Privacy via distributed noise generation. In *EUROCRYPT*. 486–503.
- [15] Elias Frenntzos, Kostas Gratsias, and Yannis Theodoridis. 2007. Index-based most similar trajectory search. In *ICDE*. 816–825.
- [16] Naosie Holohan, Spiros Antonatos, Stefano Braghin, and Pól Mac Aonghusa. 2020. The bounded Laplace mechanism in differential privacy. *Journal of Privacy and Confidentiality* 10, 1 (2020).
- [17] Kaifeng Jiang, Dongxu Shao, Stéphane Bressan, Thomas Kister, and Kian-Lee Tan. 2013. Publishing trajectories with differential privacy guarantees. In *SSDBM*. 1–12.
- [18] Anna Monreale, Roberto Trasarti, Chiara Renso, Dino Pedreschi, and Vania Bogorny. 2010. Preserving privacy in semantic-rich trajectories of human mobility. In *SIGSPATIAL*. 47–54.
- [19] Elham Naghizade, Lars Kulik, and Egemen Tanin. 2014. Protection of sensitive trajectory datasets through spatial and temporal exchange. In *SSDBM*. 1–4.
- [20] Mirco Nanni and Dino Pedreschi. 2006. Time-focused clustering of trajectories of moving objects. *Journal of Intelligent Information Systems* 27, 3 (2006), 267–289.
- [21] Dieter Pfoser, Christian S. Jensen, Yannis Theodoridis, et al. 2000. Novel approaches to the indexing of moving object trajectories. In *VLDB*. 395–406.
- [22] Jianzhong Qi, Yufei Tao, Yanchuan Chang, and Rui Zhang. 2018. Theoretically optimal and empirically efficient r-trees with strong parallelizability. *PVLDB* 11, 5 (2018), 621–634.
- [23] Jianzhong Qi, Yufei Tao, Yanchuan Chang, and Rui Zhang. 2020. Packing R-trees with space-filling curves: Theoretical optimality, empirical efficiency, and bulk-loading parallelizability. *TODS* 45, 3 (2020), 1–47.
- [24] Yasushi Sakurai, Masatoshi Yoshikawa, and Christos Faloutsos. 2005. FTW: Fast similarity search under the time warping distance. In *PODS*. 326–337.
- [25] Stan Salvador and Philip Chan. 2007. Toward accurate dynamic time warping in linear time and space. *Intelligent Data Analysis* 11, 5 (2007), 561–580.
- [26] Shuo Shang, Lisi Chen, Zhewei Wei, Christian S. Jensen, Kai Zheng, and Panos Kalnis. 2017. Trajectory similarity join in spatial networks. *PVLDB* 10, 11 (2017), 1178–1189.
- [27] Zeyuan Shang, Guoliang Li, and Zhifeng Bao. 2018. DITA: A distributed in-memory trajectory analytics system. In *SIGMOD*. 1681–1684.
- [28] Panagiotis Tampakis, Christos Doukeridis, Nikos Pelekis, and Yannis Theodoridis. 2020. Distributed subtrajectory join on massive datasets. *ACM Transactions on Spatial Algorithms and Systems* 6, 2 (2020), 1–29.
- [29] Michail Vlachos, George Kollios, and Dimitrios Gunopoulos. 2002. Discovering similar multidimensional trajectories. In *ICDE*. 673–684.
- [30] Huayu Wu, Mingqiang Xue, Jianneng Cao, Panagiotis Karras, Wee Siong Ng, and Kee Kiat Koo. 2016. Fuzzy trajectory linking. In *ICDE*. 859–870.
- [31] Yixin Xu, Lars Kulik, Renata Borovica-Gajic, Abdullah Aldwyish, and Jianzhong Qi. 2020. Highly efficient and scalable multi-hop ride-sharing. In *SIGSPATIAL*. 215–226.
- [32] Tun-Hao You, Wen-Chih Peng, and Wang-Chien Lee. 2007. Protecting moving trajectories with dummies. In *MDM*. 278–282.
- [33] Haitao Yuan and Guoliang Li. 2019. Distributed in-memory trajectory similarity search and join on road network. In *ICDE*. 1262–1273.
- [34] Yihong Yuan and Martin Raubal. 2014. Measuring similarity of mobile phone user trajectories—a Spatio-temporal Edit Distance method. *International Journal of Geographical Information Science* 28, 3 (2014), 496–520.