

Parallel Trajectory-to-Location Join

Shuo Shang¹, Lisi Chen, Kai Zheng, Christian S. Jensen, *Fellow, IEEE*, Zhewei Wei, and Panos Kalnis

Abstract—The matching between trajectories and locations, called Trajectory-to-Location join (TL-Join), is fundamental functionality in spatiotemporal data management. Given a set of trajectories, a set of locations, and a threshold θ , the TL-Join finds all (trajectory, location) pairs from the two sets with spatiotemporal correlation above θ . This join targets diverse applications, including location recommendation, event tracking, and trajectory activity analyses. We address three challenges in relation to the TL-Join: how to define the spatiotemporal correlation between trajectories and locations, how to prune the search space effectively when computing the join, and how to perform the computation in parallel. Specifically, we define new metrics to measure the spatiotemporal correlation between trajectories and locations. We develop a novel parallel collaborative (PCol) search method based on a divide-and-conquer strategy. For each location o , we retrieve the trajectories with high spatiotemporal correlation to o , and then we merge the results. An upper bound on the spatiotemporal correlation and a heuristic scheduling strategy are developed to prune the search space. The trajectory searches from different locations are independent and are performed in parallel, and the result merging cost is independent of the degree of parallelism. Studies of the performance of the developed algorithms using large spatiotemporal data sets are reported.

Index Terms—Trajectory-to-location join, parallel processing, spatial networks, spatial databases

1 INTRODUCTION

WITH the continuous proliferation of GPS-enabled mobile devices (e.g., vehicle navigation systems and smart phones) and the rapid development of online map-based services (e.g., Google Maps,¹ and MapQuest²), it is easy to collect and share trajectories, e.g., at specialized sites such as Bikely,³ GPS-way-points,⁴ Share-my-routes,⁵ and Microsoft Geolife.⁶ Also, more and more social networking sites, including Twitter,⁷ Facebook,⁸ and Foursquare,⁹ are starting to support trajectory collection and sharing [1], [2]. The availability of massive trajectory data motivates new studies in spatiotemporal data management. The matching between trajectories and locations, called Trajectory-to-Location Join (TL-Join), is

fundamental functionality. Given a set T of trajectories, a set O of locations, and a threshold θ , the TL-Join finds all (trajectory, location) pairs from T and O with a spatiotemporal correlation above θ .

The TL-Join may benefit diverse applications, including location recommendation [3], event tracking [4], and trajectory activity analyses [5], [6], [7]. For example, people may want to place new facilities (e.g., shopping malls, banks, and petrol stations) in a city according to available trajectories of the potential customers. They may use the TL-Join to find the locations that join with the most trajectories. Such locations have high visibility to trajectories and may be most attractive to customers. These locations may then maximize the commercial value of new facilities. As another example, when events occur (e.g., accidents or terrorist attacks), the police may want to find eyewitnesses of the events. The TL-Join can find such people by matching their trajectories to the events' locations. In addition, we can use the TL-Join to analyze the activities of trajectories. Depending on the points of interest (e.g., restaurants, shopping malls, and sightseeing places) that a trajectory joins with, we can infer activities related to the trajectory (e.g., dinner, shopping, and tourism).

To the best of our knowledge, this is the first trajectory-to-location matching study that takes into account both the spatial and temporal ranges when computing spatial and temporal correlations. We use a linear method [1], [2], [8] to combine the spatial and temporal correlations into a spatiotemporal correlation metric. In contrast, existing studies typically perform (i) the matching solely in the spatial domain [3], [4], [5], [6], [8] or (ii) using point-to-point matching in the spatial domain or the temporal domain [2], [3], [4], [8]. As a result, they may fail to support time-aware applications. For example, they may match a morning trajectory to an evening activity (e.g., drinking at a bar), or they may match a midnight trajectory to a facility open only during the

1. <https://maps.google.com/>
2. <https://www.mapquest.com/>
3. <https://www.bikely.com/>
4. <https://www.gps-waypoints.net>
5. <https://www.sharemyroutes.com/>
6. <https://research.microsoft.com/en-us/projects/geolife/>
7. <https://www.twitter.com/>
8. <https://www.Facebook.com/>
9. <https://www.Foursquare.com/>

- S. Shang and P. Kalnis are with the King Abdullah University of Science and Technology, Thuwal 23955, Saudi Arabia. E-mail: jedi.shang@gmail.com, panos.kalnis@kaust.edu.sa.
- L. Chen is with the University of Wollongong, Wollongong, NSW 2522, Australia. E-mail: LCHEN012@e.ntu.edu.sg.
- K. Zheng is with the University of Electronic Science and Technology of China, Chengdu, Sichuan 610054, China. E-mail: zhengkai@uestc.edu.cn.
- C. S. Jensen is with Aalborg University, Aalborg 9100, Denmark. E-mail: csj@cs.aau.dk.
- Z. Wei is with the Renmin University of China, Beijing 100872, P.R. China. E-mail: zhewei@ruc.edu.cn.

Manuscript received 25 Nov. 2017; revised 30 June 2018; accepted 5 July 2018. Date of publication 10 July 2018; date of current version 26 Apr. 2019.

(Corresponding author: Kai Zheng.)

Recommended for acceptance by M. A. Cheema.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below.

Digital Object Identifier no. 10.1109/TKDE.2018.2854705

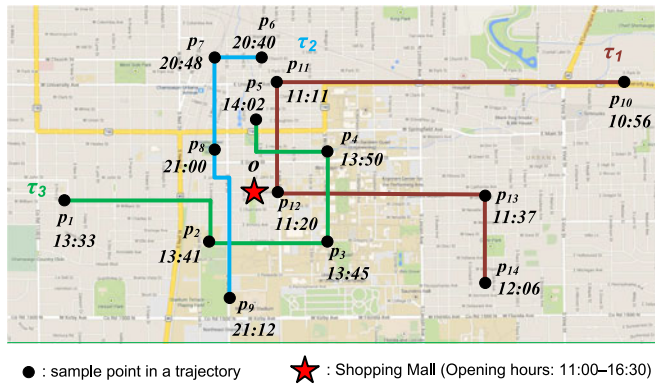


Fig. 1. TL-join example.

day (e.g., a bank or a shopping mall). Further, the matched pairs cannot guarantee a long-term and continuous correlation between locations and trajectories in the spatial and temporal domains. For example, a trajectory may have a single sample point and a short duration matching to a location (its other sample points are too far away from the location); but the result caused by this matching may be of little use in activity analyses because its matching duration is too short to denote a significant relationship. Furthermore, the so-called Semantic Enrichment [7] approach utilizes the stay-time at a location to infer a traveler's activity. This type of matching is not attractive in our intended applications because it is not flexible, i.e., it relies on a fixed visit position (e.g., the intended location) when defining a matching. In contrast, the matching in the TL-Join allows travelers to visit multiple positions close to the intended location within a matching duration. Such flexible matching is appropriate in applications such as location recommendation, event tracking, and activity analyses. The Semantic Enrichment can be viewed as a special case of the matching in the TL-Join.

An example of the TL-Join is shown in Fig. 1, where τ_1 , τ_2 , and τ_3 are trajectories and o is a location (a shopping mall) with opening hours from 11:00 to 16:30. For the matchings that solely consider the spatial domain [3], [4], [5], [6], [8], trajectory-location pairs (τ_1, o) , (τ_2, o) , and (τ_3, o) are returned because all trajectories are spatially close to o . However, the timestamps of τ_2 do not match the time range (opening hours) of o , so pair (τ_2, o) has little meaning in this scenario. For point-to-point matching [2] in the spatial and temporal domains, pairs (τ_1, o) and (τ_3, o) are returned because (p_{12}, o) and (p_4, o) are matched point pairs, where $p_{12} \in \tau_1$ and $p_4 \in \tau_3$. But there is not a long-term and continuous correlation between p_{12} and o , so this matching result is of little use in applications such as location recommendation, event tracking, and activity analyses. The TL-Join returns the pair (τ_3, o) because it has a long-term and continuous correlation in both the spatial and temporal domains (e.g., matched point pairs (p_2, o) , (p_3, o) , (p_4, o) , and (p_5, o) and a duration of around 20 minutes) and because its spatiotemporal correlation exceeds threshold θ . Notice that the Semantic Enrichment approach [7] considers the stay duration at a location, which means that multiple trajectory sample points are at the corresponding location for the stay duration (the number of sample points depends on the trajectory sampling rate and the stay duration). The TL-Join can support this special case.

TABLE 1
Trajectory-to-Location Matchings

| Studies | Query Space | Spatiotemporal Matching | Data |
|-------------------------|-------------|--|-------|
| RPNN [3] | Network | Spatial only (point-to-point) | 1.6 K |
| ATSQ [4] | euclidean | Spatial only (point-to-point) | 49 K |
| UOTS [8] | Network | Spatial only (point-to-point) | 30 K |
| PTM [2] | Network | Spatial (point-to-point) and Temporal (point-to-point) | 30 K |
| Semantic Enrichment [7] | euclidean | Spatial (point-to-point) and Temporal (range) | 13 K |
| VID Join [5], [6] | euclidean | Spatial only (range) | 12 K |
| TL-Join (our proposal) | Network | Spatial (range) and Temporal (range) | 10 M |

The TL-Join is applied in a spatial network because objects move in a spatial network rather than in euclidean space in many practical scenarios. In a spatial network, network distance is the relevant distance between two objects, and using euclidean distance [4], [5], [6], [7] may lead to errors.

An overview of a comparison to existing trajectory-location matching studies is shown in Table 1. Existing methods [2], [3], [4], [5], [6], [7], [8] cannot process the TL-Join due to four reasons. (i) Different query types: most trajectory-to-location matching studies [2], [3], [4], [7], [8] are not related to the join operation. For example, RPNN [3] concerns reverse path nearest neighbor querying; ATSQ [4], UOTS [8], and PTM [2] concern trajectory search by locations; and Semantic Enrichment [7] concerns the use of trajectories to infer travelers' activities. Their solutions cannot be used in the TL-Join because the solutions are for different query types. (ii) Different matching functions: existing studies are based on point-to-point matching [2], [3], [4], [8] or spatial-only matching [3], [4], [5], [6], [8], and their solutions are inapplicable to spatiotemporal range matching. (iii) Different query spaces: the VID join [5], [6] is conducted in euclidean space, and its spatial index and accompanying pruning techniques are not competitive in spatial networks. (iv) Parallel processing requirement: existing centralized trajectory-to-location joins (VID Join) cannot process large trajectory data sets. Based on the experiments reported in the literature [5], [6], the VID join can process at most 12 K trajectories. In contrast, our implementation of the TL-Join can process 10 M trajectories with a reasonable runtime (the PCol solution can process $10 \text{ M} \times 0.5 \text{ M}$ (trajectory, location) pairs with 120 threads in 651 seconds). Table 1 offers further details on the scale of the data considered, indicating that we consider several orders of magnitude more data than do previous studies.

Next, the algorithm used for computing the TS-Join [1] cannot process the TL-Join because the query arguments are different (two trajectory sets versus a trajectory and a location sets) and because the matching functions are different (point-to-point matching versus range matching). The TL-Join needs its own specific solutions.

We propose two baseline solutions to the TL-Join, called parallel temporal-first search (PTF) and parallel spatial-first search (PSF). For PTF, we improve the equal-partition grid index in the TS-Join [1], and we propose a new balanced grid index in the temporal domain (i.e., each leaf node has similar numbers of trajectories and locations so that the parallel computation load is balanced). We define spatiotemporal correlation upper and lower bounds to prune the search space, and we perform the refinement of (trajectory, location) pairs from the leaf nodes towards the root. The computations at each index level occur in parallel. The main drawback of PTF is threefold: (i) weak spatial pruning power (temporal driven pruning), (ii) high merging cost (having more leaf nodes enables more parallel processing, but also higher merging cost), and (iii) additional computation cost to acquire network distances when computing spatial correlations.

Next, PSF is based on a divide-and-conquer strategy and performs better than PTF. For each location o , PSF explores the spatial domain to find trajectories with high spatial correlation to o . In the temporal domain, it checks whether the corresponding timestamps are within the time range of o (temporal correlation). We define upper bounds on the spatial correlation to prune the search space. Each trajectory search is independent and is performed in parallel, and the merging cost is independent of the degree of parallelism. The network distances needed for spatial correlation computations can be derived directly during trajectory searches from locations. The limitation of PSF lies in its weak pruning power in the temporal domain.

To process the TL-Join more efficiently, we propose a novel parallel collaborative search (PCol) approach. PCol uses the parallel mechanism of PSF. For each location o , PCol explores the spatial and temporal domains concurrently to find trajectories with high spatiotemporal correlation to o . We define upper bounds on the spatiotemporal correlation and a heuristic scheduling strategy that result in strong pruning power in the two domains.

To sum up, we make the following contributions:

- We propose a new trajectory-to-location join, called TL-Join, targeting applications such as location recommendation, event tracking, and trajectory activity analyses.
- The TL-Join takes both spatial and temporal range matching into account to compute spatiotemporal correlation. No other proposal provides this functionality.
- We develop two baseline algorithms for computing the TL-Join called parallel temporal-first (PTF) search and parallel spatial-first (PSF) search.
- We develop a parallel collaborative algorithm (PCol) with effective pruning techniques and a heuristic scheduling strategy in the spatial and temporal domains.
- We conduct extensive experiments on large trajectory data sets to study the performance of the developed algorithms. We can handle about 3 orders of magnitude more trajectories than the state-of-the-art VID join.

The rest of the paper is organized as follows. Section 2 defines the setting, including spatial networks, trajectories, locations, and the spatiotemporal correlation metrics considered in the paper; it ends by defining the problem.

Parallel temporal-first (PTF) search and parallel spatial-first (PSF) search are covered in Sections 3 and 4, while parallel collaborative (PCol) search is covered in Section 5. Experimental results are presented in Section 6. Related work is covered in Section 7, and conclusions and future directions are presented in Section 8.

2 PRELIMINARIES AND PROBLEM DEFINITION

2.1 Spatial Networks

A spatial network is modeled as a connected, undirected, and weighted graph $G = (V, E, F, W)$, where V is a vertex set and $E \subseteq \{\{v_i, v_j\} | v_i, v_j \in V \wedge v_i \neq v_j\}$ is an edge set. A vertex $v_i \in V$ represents a road intersection or an end of a road, and an edge $e_k = \{v_i, v_j\} \in E$ represents a road segment that enables travel between vertices v_i and v_j . Function $F : V \cup E \rightarrow \text{Geometries}$ maps a vertex to the point location of the corresponding road intersection and maps an edge to a polyline representing the corresponding road segment. Function $W : E \rightarrow \mathbb{R}$ assigns a real-valued weight $W(e)$ to an edge e that represents the corresponding road segment's length.

The shortest path between two vertices v_i and v_j is a sequence of edges linking v_i and v_j such that the sum of the edge weights is minimal. Such a path is denoted by $SP(v_i, v_j)$, and its length is denoted by $sd(v_i, v_j)$. euclidean-space based spatial indices (e.g., the R-tree [9]) and accompanying techniques are ineffective in network environments due to loose lower bounds.

For simplicity, we assume that the data points considered (e.g., trajectory sample points) are located at vertices. It is straightforward to also support data points on edges. Assume a data point p is on an edge e with given network distances to the two end vertices e_a and e_b . Then, a new vertex is created for p with the appropriate geometry, and edge e is replaced by edges (e_a, p) and (p, e_b) with the appropriate weights and geometries.

2.2 Trajectories and Locations

Raw trajectory samples obtained from GPS devices are typically of the form (longitude, latitude, time), and trajectory sample points are captured periodically at some sampling rate. We assume that all sample points have already been map matched onto the spatial network using a map-matching algorithm (e.g., [10], [11]) and that an object always follows the shortest path when moving between two adjacent sample points p_a and p_b . A trajectory is defined as follows.

Definition: Trajectory

A trajectory τ of a moving object is a finite, time-ordered sequence $\langle v_1, v_2, \dots, v_n \rangle$, where $v_i = (p_i, t_i)$, $i \in [1, n]$, with p_i being a sample point (equal to some vertex in $G.V$) and t_i being a timestamp.

Assuming that $\tau.sr$ is the sampling rate of trajectory τ , we have that $t_{i+1} - t_i = \tau.sr$, $i \in [1, n-1]$.

The above modeling of spatial networks and trajectories aligns with previous studies [1], [2].

Definition: Location

A location o contains a spatial attribute $o.p$ and a temporal attribute $o.R$, where $o.p$ is a vertex in $G.V$ and $o.R$ is a time range. Time range $o.R$ describes the valid duration of o (e.g., opening hours of facilities, or time ranges of events).

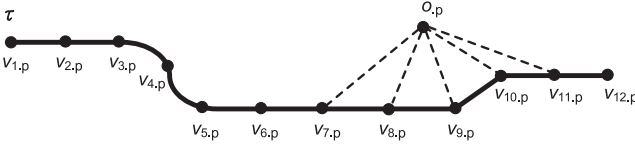


Fig. 2. Spatiotemporal correlation.

The values of timestamps and time ranges are set to be within the range of 24 hours, and the date is not taken into account because in many practical scenarios, like urban transportation, movements are often studied within the range of a day [1], [2].

2.3 Spatiotemporal Correlation

Given a location o and a trajectory τ , K_τ is the set of trajectory sample points in τ that are spatially closest to $o.p$: $\forall v \in K_\tau (\forall v' \in \tau \setminus K_\tau (sd(o.p, v.p) \leq sd(o.p, v'.p)))$. The cardinality of K_τ is set as follows.

$$k_\tau = |K_\tau| = \left\lceil \frac{d_c}{\tau.sr} \right\rceil + 1. \quad (1)$$

Here, d_c controls the coupling duration (to describe the term of correlations) between o and τ . Its value is user-defined. We assume that trajectories are sampled uniformly. As different locations may have different coupling duration d_c and different trajectories may have different sampling rates $\tau.sr$, the value of k_τ may be different for different trajectories. The following algorithms support this.

The spatial correlation $C_S(o, \tau)$ and the temporal correlation $C_T(o, \tau)$ between o and τ are defined as follows.

$$C_S(o, \tau) = \frac{\sum_{v_i \in K_\tau} e^{-sd(o.p, v_i.p)}}{k_\tau} \quad (2)$$

$$C_T(o, \tau) = \frac{|\{v_j.t | v_j \in K_\tau \wedge v_j.t \in o.R\}|}{k_\tau}. \quad (3)$$

In the spatial domain, we count the sum of the spatial distances between location o and trajectory sample points in K_τ , while in the temporal domain, we check the validity of the sample points in K_τ by matching their timestamps to the time range $o.R$.

An example that illustrates these definitions is shown in Fig. 2, where o is an object and $\tau = \langle v_1, v_2, \dots, v_{12} \rangle$ is a trajectory. The coupling duration d_c is 8 minutes, and the sampling rate of τ is 2 minutes, so $k_\tau = \lfloor \frac{8 \text{ minutes}}{2 \text{ minutes}} \rfloor + 1 = 4 + 1 = 5$. Points $v_7.p, v_8.p, \dots, v_{11}.p$ are the top-5 trajectory sample points spatially closest to $o.p$, so $K_\tau = \{v_7.p, v_8.p, \dots, v_{11}.p\}$. The value of $C_S(o, \tau)$ is computed by substituting K_τ into Equation (2). Next, assuming that $v_7.t=12:50$, $v_8.t=12:52$, $v_9.t=12:54$, $v_{10}.t=12:56$, $v_{11}.t=12:58$, and $o.R = [12:55, 13:00]$, we have that $v_7.t \notin o.R$, $v_8.t \notin o.R$, $v_9.t \notin o.R$, $v_{10}.t \in o.R$, and $v_{11}.t \in o.R$. According to Equation (3), $C_T(o, \tau) = |\{v_{10}.t, v_{11}.t\}|/5 = 2/5 = 0.4$.

The spatial and temporal correlations of τ are both in the range $[0, 1]$. We use a linear method [1], [2] to combine the spatial and temporal correlations (Equations (2) and (3)), and the spatiotemporal correlation is defined as follows.

$$C_{ST}(o, \tau) = \lambda \cdot C_S(o, \tau) + (1 - \lambda) \cdot C_T(o, \tau). \quad (4)$$

Here, parameter $\lambda \in [0, 1]$ controls the relative importance of the spatial and temporal correlations. The value of λ can be adjusted at query time.

2.4 Problem Definition

Given a set T of trajectories, a set O of locations, and a threshold θ , the trajectory-to-location join (TL-Join) returns the set A of all (trajectory, location) pairs from the two sets whose spatiotemporal correlations are at least θ , i.e., $\forall (\tau_i, o_j) \in A (C_{ST}(\tau_i, o_j) \geq \theta) \wedge \forall (\tau'_i, o'_j) \in ((T \times O) \setminus A) (C_{ST}(\tau'_i, o'_j) < \theta)$.

| A List of Notions | |
|--------------------|---|
| Notion | Description |
| $G.V$ | the set of vertices in graph G |
| $G.E$ | the set of edges in graph G |
| $sd(p_i, p_j)$ | shortest path distance between vertices p_i and p_j |
| $\tau.sr$ | sampling rate of trajectory τ |
| $o.R$ | time range of location o |
| d_c | coupling duration |
| K_τ | the set of top- k_τ sample points in τ that are spatially closest to location o |
| C_S, C_T, C_{ST} | spatial, temporal, and spatiotemporal correlation |
| λ | the relative importance of the spatial and temporal correlations |
| UB, LB | global upper and lower bounds |

3 PARALLEL TEMPORAL-FIRST SEARCH

3.1 Basic Idea

Parallel temporal-first (PTF) search is a baseline approach to TL-Join processing. We improve the equal-partition grid index used in the TS-Join [1], and we propose a new balanced hierarchical grid index in the temporal domain (Section 3.2). We also define upper and lower bounds to prune the search space in the spatial and temporal domains. PTF refines the (trajectory, location) pairs in the same leaf node and merges the results from the leaf nodes towards the root. The join result is then obtained from the root. The computations at the nodes at the same level occur in parallel (Section 3.3). The pseudocode of PTF and its time complexity are given in Section 3.4.

3.2 Balanced Grid Index

In the TS-Join [1], the temporal domain is partitioned into m equal-sized time slots, each of which is assigned to a leaf node. The drawback of this approach is that the distributions of trajectories and locations are imbalanced, and different leaf nodes may have quite different numbers of trajectories (e.g., peak hours may have more, off-peak hours may have fewer, and midnight may have none). Such imbalance yields poor performance in parallel processing. To address this issue, we propose a new balanced grid index in the temporal domain for PTF. Here, each leaf node n has a matching-times upper bound $M \geq |n_\tau| \times |n_o|$, where n_τ and n_o are the sets of trajectories and locations contained in n . The optimal value of M that achieves the highest performance is determined through extensive experiments. Notice that the balanced grid index is a temporal index, which indexes the time ranges of trajectories and locations. Other trajectory indexes (e.g., [12], [13], [14]) are spatial index, and they are not suitable for this scenario.

The balanced grid index is constructed as follows. Given a value of M and a slot $s = [0, 24:00]$, we recursively partition s into two equal-sized nodes if $|s_\tau| \times |s_o| > M$, where s_τ and s_o are sets of trajectories and locations in slot s . For example, given a trajectory $\tau = \langle v_1, v_2, \dots, v_i \rangle$, its temporal

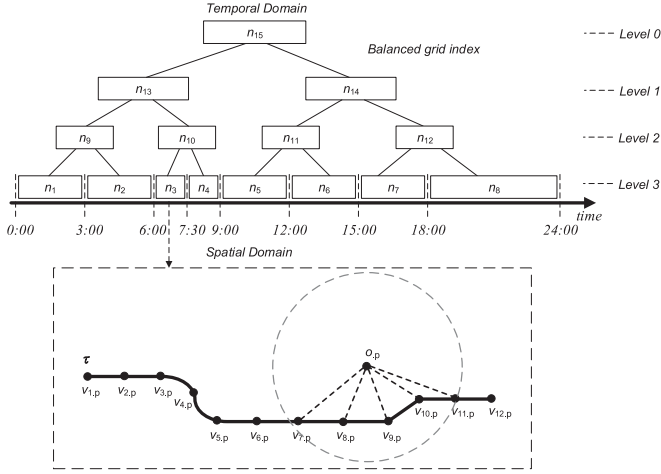


Fig. 3. Example of PTF.

range $\text{range}(\tau) = [v_1.t, v_i.t]$. If $\text{range}(\tau) \subseteq \text{range}(s)$, τ is contained in s . Similarly, given a location o , if $o.R \subseteq \text{range}(s)$, o is contained in s . For example, given $\text{range}(s) = [9:00, 12:00]$, $\text{range}(\tau) = [10:00, 11:00]$, and $o.R = [9:30, 11:30]$, τ and o are contained in s .

Once the partitioning terminates, each slot corresponds to a leaf node. We build a tree structure bottom-up. Assume that there are m nodes at the leaf level. Then we build $\lceil \frac{m}{2} \rceil$ parent nodes. We do this recursively until there is one parent, which is the root. The height of the tree is $\lceil \log(m) \rceil + 1$. An example is shown in Fig. 3, where n_1, n_2, \dots, n_8 are leaf nodes and n_{15} is the root. Each trajectory τ and each location o are stored in the lowest node n that fully covers its temporal range, i.e., $\text{range}(\tau) \subseteq \text{range}(n)$ and $o.R \subseteq \text{range}(n)$ and $\text{range}(\tau)$ and $o.R$ are not contained in the range of any child node of n . For example, given $o'.R = [9:30, 17:30]$, o' is stored in n_{14} ($\text{range}(n_{14}) = [9:00, 24:00]$) because $o'.R \subseteq \text{range}(n_{14})$ and $o'.R \not\subseteq \text{range}(n_{11})$ and $o'.R \not\subseteq \text{range}(n_{12})$ (n_{11} and n_{12} are child nodes of n_{14}).

3.3 Filtering, Refinement, and Merging

In the example in Fig. 3, a trajectory τ and a location o are stored in node n_3 . As they are temporally close to each other, we estimate the upper bound on their temporal correlation $C_T(o, \tau)$ (cf. Equation (3)) as follows.

$$\begin{aligned} |\{v_j | v_j \in K_\tau \wedge v_j.t \in o.R\}| &\leq k_\tau \\ \Rightarrow C_T(o, \tau).ub = 1 &\geq C_T(o, \tau). \end{aligned} \quad (5)$$

By substituting Equation (5) into Equation (4), we have that

$$\begin{aligned} C_{ST}(o, \tau) &= \lambda \cdot C_S(o, \tau) + (1 - \lambda) \cdot C_T(o, \tau) \geq \theta \\ \Rightarrow C_S(o, \tau) &\geq \frac{\theta - (1 - \lambda) \cdot C_T(o, \tau).ub}{\lambda} = \frac{\theta - 1 + \lambda}{\lambda}. \end{aligned}$$

For each “qualified” (trajectory, location) pair (o, τ) (i.e., $C_{ST}(o, \tau) \geq \theta$), its spatial correlation exceeds the value of $\frac{\theta - 1 + \lambda}{\lambda}$. We define a global lower bound LB_S of the spatial correlation between (trajectory, location) pairs in the same leaf node as follows.

$$LB_S = \frac{\theta - 1 + \lambda}{\lambda}. \quad (6)$$

We use network expansion to compute the spatial correlation $C_S(o, \tau)$ (Equation (2)). The network expansion is performed from location o using Dijkstra’s algorithm [15]. Dijkstra’s algorithm always selects the vertex with the minimum distance label for expansion. Hence, the first k_τ sample points in τ scanned by the expansion are just the top- k_τ sample points closest to o . For example, in Fig. 3, assuming $k_\tau = 5$ and $v_7.p, v_8.p, \dots, v_{11}.p$ are top-5 first scanned sample points in τ . According to Equation (2), $C_S(o, \tau) = \frac{1}{5}(e^{-d(v_7.p, o.p)} + e^{-d(v_8.p, o.p)} + \dots + e^{-d(v_{11}.p, o.p)})$. If $C_S(\tau, o) < LB_S$, then $C_{ST}(o, \tau) < \theta$, and the (trajectory, location) pair (o, τ) can be pruned safely. Otherwise, we compute the exact spatiotemporal correlation $C_{ST}(o, \tau)$ (Equation (4)) and compare to θ to check the pair’s validity. The computations in different leaf nodes are independent and occur in parallel.

Having computed the spatiotemporal correlations of the (trajectory, location) pairs in the leaf nodes, we merge the results from the leaf level to the root level (bottom-up). At each level, when two nodes n and n' have the same parent n'' , we merge their results and assign this to the parent (e.g., merge n_3, n_4 , and n_{10} to obtain the result for n_{10} in Fig. 3). In addition to these qualified results ($C_{ST}(o, \tau) \geq \theta$), we also need to consider the (trajectory, location) pairs (o, τ) in the following three cases: (i) one item is stored in n or n' and another item is stored in n'' (e.g., $\text{range}(\tau) \subseteq \text{range}(n)$ and $o.R \subseteq \text{range}(n'')$); (ii) two items are stored in n'' (e.g., $\text{range}(\tau) \subseteq \text{range}(n'')$ and $o.R \subseteq \text{range}(n'')$); (iii) one item is stored in n and another item is stored in n' (e.g., $\text{range}(\tau) \subseteq \text{range}(n)$ and $o.R \subseteq \text{range}(n')$).

For the first and the second cases, we use the same lower and upper bounds (Equations (5) and (6)) and pruning techniques as we use for the (trajectory, location) pairs in the same node. The qualified pairs are stored in n'' . For the third case, as trajectory τ and location o are stored in different nodes, we have that $C_T(\tau, o) = 0$. By substituting this into Equation (4), we have that

$$C_{ST}(\tau, o) \geq \theta \Leftrightarrow C_S(\tau, o) \geq \frac{\theta}{\lambda}. \quad (7)$$

As the value of $C_S(\tau, o)$ is in the range $[0, 1]$, if $\theta > \lambda$, (trajectory, location) pair (τ, o) is pruned directly. Otherwise, we compute the spatiotemporal correlation $C_{ST}(\tau, o)$ and compare to θ to check the pair’s validity.

The merging processes of adjacent node pairs (e.g., merge n_1 and n_2 to n_9 , n_3 and n_4 to n_{10}) at the same level of the tree are independent. Thus they again occur in parallel. Having merged the computation results from the leaf nodes all the way to the root node, the join result in $[0:00, 24:00]$ is found.

3.4 Algorithm and Time Complexity

The pseudocode of PTF is shown in Algorithm 1. The computation is bottom-up, and h is the current level of computation. Initially, we compute the global spatial lower bound LB_S (Equation (6)) for leaf nodes (lines 1–2). For each (trajectory, location) pair (τ, o) in n (i.e., $\text{range}(\tau) \subseteq \text{range}(n)$ and $o.R \subseteq \text{range}(n)$), we compute its spatial correlation $C_S(\tau, o)$ (Equation (2)), and if $C_S(\tau, o)$ is less than LB_S , pair (τ, o) is pruned (lines 3–7). Otherwise, we compute the exact spatiotemporal correlation $C_{ST}(\tau, o)$ (Equation (4)), and if it is no less than θ , we store (τ, o) in P_n (lines 8–10). Having refined all leaf nodes, we merge the results from the leaf level

towards the root. If two nodes n_a and n_b are at the same level and they have the same parent node n_c , we merge the results for n_a , n_b , and n_c (e.g., n_1 , n_2 , and n_9 in Fig. 3) and store the qualified (trajectory, location) pairs in P_{n_c} (lines 11–14). If $h = 1$, the root node n_c is reached, and all (trajectory, location) pairs stored in P_{n_c} are returned. Otherwise, we repeat the procedure for the next level of the tree (line 15–17).

Algorithm 1. PTF Search

Data: a balanced grid index tree T_r , a trajectory set T , a location set O , and a threshold θ
Result: $\{(\tau, o) | C_{ST}(\tau, o) \geq \theta, \forall \tau \in T, \forall o \in O\}$

```

1  $h \leftarrow T_r.height - 1;$ 
2 compute  $LB_S$ ;
3 for each leaf node  $n$  in  $T_g$  do
4   for each (trajectory, location) pair  $(\tau, o)$  in  $n$  do
5     compute  $C_S(\tau, o)$ ;
6     if  $C_S(\tau, o) < LB_S$  then
7       prune  $(\tau, o)$ ;
8     compute  $C_{ST}(\tau, o)$ ;
9     if  $C_{ST}(\tau, o) \geq \theta$  then
10       $P_n.add(\tau, o)$ ;
11 while true do
12   if  $n_a, n_b$  in level  $h$ ,  $n_a.parent = n_b.parent = n_c$  then
13     merge  $n_a, n_b$ , and  $n_c$ ;
14     compute and store qualified (trajectory, location) pairs in  $P_{n_c}$ ;
15   if  $h = 1$  then
16     return  $P_{n_c}$ ;
17    $h \leftarrow h - 1$ ;
```

Let $|T|$ and $|O|$ denote the cardinalities of trajectory set T and location set O . We use $|V|$ and $|E|$ to denote the numbers of vertices and edges in G . Then $O(|V|\log|V| + |E|)$ is the time complexity of computing the spatial correlation between a trajectory and a location by using Dijkstra's algorithm. PTF follows the filter-and-refine paradigm, and the time complexity of the filtering phase is $O((|V|\log|V| + |E|)|T||O|)$.

The time complexity to verify candidates by computing their exact spatiotemporal correlations is $O(k_\tau|C|)$ (the spatial correlations are computed in the filtering phase, so in the refinement phase we only need to compute the temporal correlations), where $|C|$ is the cardinality of the candidate set and $C \subseteq P \times O$. The total time complexity is $O((|V|\log|V| + |E|)|T||O| + k_\tau|C|) = O((|V|\log|V| + |E|)|T||O|)$, which does not depend on the candidate set size.

The computations for nodes at the same level of the tree occur in parallel. If we have multiple cores and threads (each leaf node corresponds to a thread), it is possible to accelerate the computation at the leaf level by generating many leaf nodes and processing them in parallel. However, more leaf nodes also leads to more tree levels (m is the number of leaf number, and the height of the tree is $\lceil \log(m) \rceil + 1$), which increases the merging cost.

4 PARALLEL SPATIAL-FIRST SEARCH

4.1 Basic Idea

PTF has three weaknesses. (i) Weak spatial pruning power: the pruning is driven by the temporal domain so it has low effectiveness in the spatial domain. (ii) High merging cost:

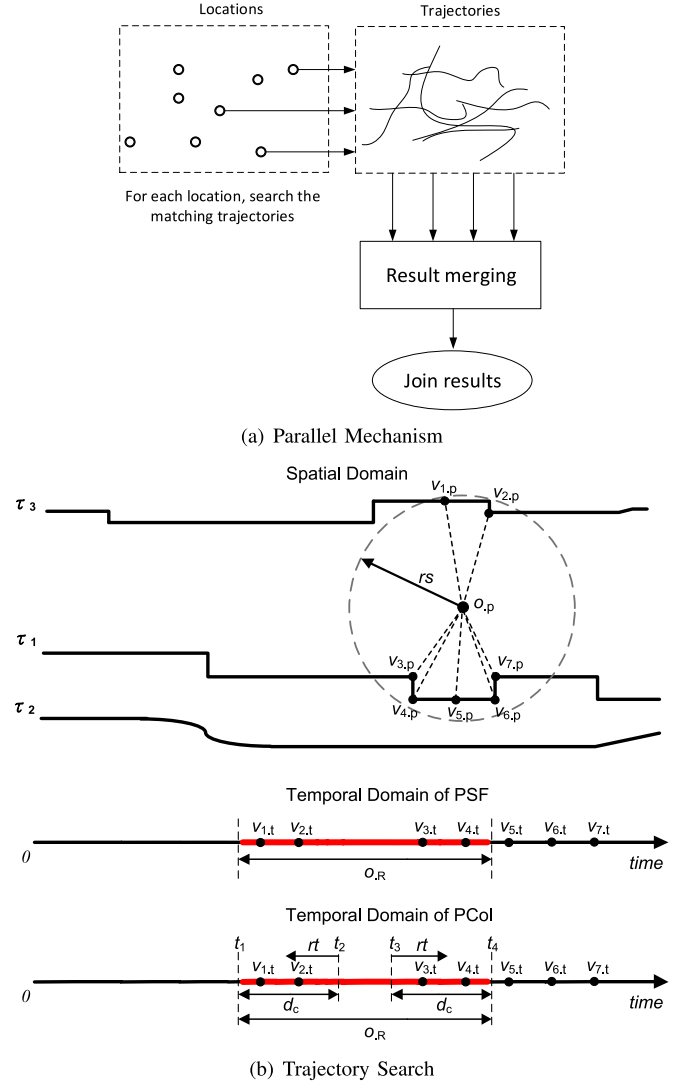


Fig. 4. Examples of PSF and PCOL.

more leaf nodes (each leaf node corresponds to a thread) lead to a higher merging costs, which decreases performance. (iii) Additional network distance computations are needed to compute the spatial correlations (Equation (2)), which again yields poor performance.

Parallel spatial first (PSF) search is another baseline for TL-Join computation. Its parallel mechanism is shown in Fig. 4a. For each location $o \in O$, we search the trajectories with high spatiotemporal correlations to o . The trajectory-search processes at different locations are performed in parallel. In the spatial domain, we use network expansion [15] to explore the spatial network and to find trajectories spatially close to o (spatial correlation). In the temporal domain, we check whether the corresponding timestamps are within the time range of o (temporal correlation). Upper and lower bounds on the spatiotemporal correlations are defined to prune the search space. By merging the search results from each location, the solution of the TL-join is found. Compared to PTF, PSF has two advantages. First, its result merging cost is independent of the degree of parallelism. We can simply combine the trajectory-search results of all locations to get the solution. Second, the network distances for the spatial correlation computation can be acquired during the

trajectory-search processes. PSF has better time complexity than PTF.

4.2 Filtering, Refinement, and Merging

An example of PSF is given in Fig. 4b, where o is a location and $o.R$ is its time range; τ_1 , τ_2 , and τ_3 are trajectories; $v_1.p$ and $v_2.p$ are the top-2 vertices in τ_3 spatially closest to o ; $v_3.p, v_4.p, \dots, v_7.p$ are the top-5 vertices in τ_1 spatially closest to o , and $v_1.t, \dots, v_7.t$ are the corresponding timestamps. Assuming $k_{\tau_1} = k_{\tau_2} = k_{\tau_3} = 5$.

In the spatial domain, network expansion is performed from o according to Dijkstra's algorithm [15]. The explored space is a circular region (o, rs) with center o and radius rs . As Dijkstra's algorithm always selects the vertex with the minimum distance label for expansion, the top- k first scanned vertices in τ are the top- k vertices spatially closest to o . For example, in Fig. 4b, $v_1.p$ and $v_2.p$ are the top-2 first scanned vertices in τ_3 , and $v_3.p, v_4.p, \dots, v_7.p$ are the top-5 first scanned vertices in τ_1 .

Assuming a trajectory τ has $\tau.k$ vertices that have been scanned by the expansion from o . If $\tau.k \geq k_\tau$, trajectory τ is called "fully scanned" (e.g., τ_1 in Fig. 4b). If $k_\tau > \tau.k > 0$, τ is called "partly scanned" (e.g., τ_3 in Fig. 4b). If $\tau.k = 0$, τ is called "unscanned" (e.g., τ_2 in Fig. 4b).

For a partly scanned trajectory τ' , we estimate an upper bound on its spatial correlation as follows. Assuming that $v_j.p \in \tau'$ is an unscanned vertex in the spatial domain, we have:

$$rs < sd(o.p, v_j.p) \Rightarrow e^{-sd(o.p, v_j.p)} < e^{-rs}.$$

By substitution into Equation (2), the spatial correlation upper bound $C_S(\tau, o).ub$ is defined as follows.

$$C_S(\tau, o).ub = \frac{\sum_{v_i.p \in V_s} e^{-sd(o.p, v_i.p)} + \sum_{v_j.p \in V_u} e^{-rs}}{k_\tau}. \quad (8)$$

Here, V_s is a set of scanned vertices in τ ($V_s \subset K_\tau$), and V_u is the set of unscanned top- k_τ vertices in τ ($V_s \cup V_u = K_\tau$ and $|V_s \cup V_u| = k_\tau$). Among all partly scanned trajectories in the spatial domain, we define a global upper bound on spatial correlation as

$$UB_S = \max_{\tau \in T_p} \{C_S(\tau, o).ub\}, \quad (9)$$

where T_p is a set of partly scanned trajectories in the spatial domain, and the value of UB_S changes dynamically during the query processing.

Filter-and-Refine: if $UB_S < \theta$, we prune all partly scanned and unscanned trajectories. For fully scanned trajectories, we compute the exact spatial (Equation (2)) and temporal (Equation (3)) correlations. The spatiotemporal correlation $C_{ST}(\tau, o)$ is derived by combining them (Equation (4)). For example, in Fig. 4b, τ_1 is fully scanned, so we compute $C_S(\tau_1, o.p) = \frac{1}{5}(e^{-d(v_3.p, o.p)} + e^{-d(v_4.p, o.p)} + e^{-d(v_5.p, o.p)} + e^{-d(v_6.p, o.p)} + e^{-d(v_7.p, o.p)})$, $C_T(\tau_1, o) = \frac{\{v_3.t, v_4.t\}}{5} = \frac{2}{5} = 0.4$, and $C_{ST}(\tau_1, o.p) = \lambda \cdot C_S(\tau_1, o.p) + (1 - \lambda) \cdot C_T(\tau_1, o)$. If the spatiotemporal correlation $C_{ST}(\tau, o)$ does not exceed θ , we prune trajectory τ . Otherwise, (trajectory, location) pair (τ, o) is stored in M_o (the set of matched (trajectory, location) pairs of o). By combining M_o of all objects in O , the result $\bigcup_{o \in O} M_o$ of the TL-Join is found.

Notice that we do not maintain upper bounds on the unscanned trajectories to reduce the computation and storage cost. Given a partly scanned trajectory τ and an unscanned trajectory τ' (e.g., τ_2 in Fig. 4b), according to Equations (8) and 9, we have:

$$C_S(\tau', o).ub = \frac{\sum_{v_j.p \in V_u} e^{-rs}}{k'_\tau} = e^{-rs} < C_S(\tau, o).ub \leq UB_S. \quad (10)$$

Here, $V_u = K_{\tau'}$ and $|V_u| = k'_\tau$. If $UB_S < \theta$, we have that $C_S(\tau', o).ub < \theta$. So the unscanned trajectories can be pruned safely, and it is not necessary to maintain spatial upper bounds.

4.3 Algorithm and Time Complexity

PSF adopts a divide-and-conquer strategy. For each location o in set O , we retrieve the trajectories with high spatiotemporal correlation to o . The search processes for different locations are independent so they are performed in parallel. Unlike for PTF, PSF has a constant merging cost (its merging cost is independent of the degree of parallelism, and we simply combine the search result of each location to achieve join result). The pseudocode of PSF is shown in Algorithm 2.

Initially, for each location $o \in O$, the set of its matched (trajectory, location) pairs M_o is set to \emptyset . The global spatial upper bound UB_S is set to 0. For each trajectory $\tau \in T$, the number of its scanned vertices $\tau.k$ is set to 0. We perform network expansion from each location o to explore the spatial network (lines 1–4). For each newly scanned vertex p , all trajectories passing P have one more scanned vertex (lines 5–7). If the number of scanned vertices of τ is equal to k_τ (τ is fully scanned), we compute its spatiotemporal correlation $C_{ST}(\tau, o)$ (Equation (4)). If the value of $C_{ST}(\tau, o)$ exceeds that of θ , we store (trajectory, location) pair (τ, o) in M_o . Then, we remove τ from the partly scanned trajectory set T_p and update the value of UB_S (lines 8–13). If τ is partly scanned ($0 < \tau.k < k_\tau$), we compute its spatial correlation upper bound $C_S(\tau, o).ub$, and we update the value of the global spatial upper bound UB_S (lines 14–17). If the value of UB_S does not exceed that of θ , the expansion from o terminates (lines 18–20). Having searched all locations, we combine their results and get the result $\bigcup_{o \in O} M_o$ of the TL-Join (line 21).

Let $|O|$ denote the cardinality of location set O and let T_θ denote the scanned trajectory set for the search process from each location, which includes the partly and fully scanned trajectories ($T_\theta = T_p \cup T_f$). According to Equations (8) and (9), the maximum spatial expansion radiuses rs is inversely proportional to θ . Assuming the trajectories are uniformly distributed in the spatial domain, it follows that $|T_\theta|$ is inversely proportional to θ . Thus, $|T_\theta|$ is sensitive to the value of threshold θ and the pruning effectiveness. We use $|V|$ and $|E|$ to denote the numbers of vertices and edges in G . Then $O(|V| \log |V| + |E|)$ is the time complexity of network expansion using Dijkstra's algorithm. The time complexity of PSF is $O(|T_\theta| |O| (|V| \log |V| + |E|))$. If the value of θ is sufficiently large, the time complexity is close to $O(|O| (|V| \log |V| + |E|))$.

5 PARALLEL COLLABORATIVE SEARCH

5.1 Basic Idea

The main weakness of PSF lies in its weak temporal pruning power since its pruning is driven by the spatial domain. To

overcome that weakness and to process the TL-Join more efficiently, we propose a parallel collaborative (PCol) search algorithm that improves PSF. In contrast to PSF, PCol performs trajectory search in the spatial and temporal domains concurrently. An upper bound on the spatiotemporal correlation and a heuristic search strategy are proposed to prune the search space. PCol follows the same parallel mechanism as PSF (cf. Fig. 4a). Compared to PSF, PCol has stronger pruning power, which should translate into higher performance.

Algorithm 2. PSF Search

Data: a set O of locations, a set T of trajectories, and a threshold θ

Result: $\bigcup_{o \in O} M_o$

```

1  $\forall o \in O (M_o \leftarrow \emptyset);$ 
2 for each location  $o$  in  $O$  do
3    $UB_S \leftarrow 0;$ 
4    $\forall \tau \in T (\tau.k \leftarrow 0);$ 
5    $p \leftarrow \text{expand}(o);$ 
6   for each trajectory  $\tau$  passing  $p$  do
7      $\tau.k \leftarrow \tau.k + 1;$ 
8     if  $\tau.k = k_\tau$  then
9       compute  $C_{ST}(\tau, o);$ 
10      if  $C_{ST}(\tau, o) \geq \theta$  then
11         $M_o.add(\tau, o);$ 
12         $T_p.remove(\tau);$ 
13        update  $UB_S;$ 
14      if  $0 < \tau.k < k_\tau$  then
15        update  $C_S(\tau, o).ub;$ 
16        if  $C_S(\tau, o).ub > UB_S$  then
17           $UB_S \leftarrow C_S(\tau, o).ub;$ 
18      if  $UB_S < \theta$  then
19        store  $M_o;$ 
20        break;
21 return  $\bigcup_{o \in O} M_o;$ 

```

5.2 Upper Bound

In the spatial domain, PCol, like PSF, adopts network expansion [15] to explore the spatial network and to find trajectories with high spatial correlation to the query location o . In the temporal domain, we partition time range $o.R$ into three parts (if $|o.R| > 2d_c$). An example is shown in Fig. 4b, where $|\text{range}(t_1, t_2)| = |\text{range}(t_3, t_4)| = d_c$ and d_c is the coupling duration between τ and o . Initially we search the trajectory timestamps in $\text{range}(t_2, t_3)$, and then we expand the search from t_2 and t_3 concurrently towards the boundaries of $o.R$, and rt is the radius of the search space. If $|o.R| \leq 2d_c$, we only partition $o.R$ into two parts from the middle point (i.e., merging t_2 and t_3 in Fig. 4b to the middle point), and then we expand the search from the middle point towards the boundaries.

We estimate the upper bound on the temporal correlation of an unscanned trajectory τ as follows.

$$|\text{range}(t_1, t_2 - rt)| = |\text{range}(t_3 + rt, t_4)| = d_c - rt$$

$$\Rightarrow C_T(\tau, o).ub = \frac{\left\lfloor \frac{|\text{range}(t_1, t_2 - rt)|}{\tau.sr} \right\rfloor + 1}{k_\tau} = \frac{\left\lfloor \frac{d_c - rt}{\tau.sr} \right\rfloor + 1}{k_\tau}. \quad (11)$$

Here, $\tau.sr$ is the sampling rate of trajectory τ , and $\text{range}(t_1, t_2 - rt)$ and $\text{range}(t_3 + rt, t_4)$ are the unscanned spaces in $o.R$. Because trajectories are sampled continuously

and uniformly and because $\text{range}(t_2 - rt, t_3 + rt)$ has been scanned in the current step, it is impossible for an unscanned trajectory to appear in both $\text{range}(t_1, t_2 - rt)$ and $\text{range}(t_3 + rt, t_4)$. Notice that for trajectories with non-uniform sampling rate, we simply need to count the number n of sample points in the corresponding time range, or to use the minimum sampling rate of τ to compute the bounds.

By combining the upper bounds on the spatial (Equation (10)) and temporal (Equation (11)) correlation according to Equation (4), we obtain an upper bound $C_{ST}(\tau, o).ub$ of the spatiotemporal correlation. The value of $C_{ST}(\tau, o).ub$ is used as a global upper bound UB for all unscanned trajectories in both domains, and it changes dynamically during query processing.

$$C_{ST}(\tau, o).ub = \lambda \cdot C_S(\tau, o).ub + (1 - \lambda) \cdot C_T(\tau, o).ub$$

$$\Rightarrow UB = C_{ST}(\tau, o).ub = \lambda \cdot e^{-rs} + (1 - \lambda) \cdot \frac{\left\lfloor \frac{d_c - rt}{\tau.sr} \right\rfloor + 1}{k_\tau}. \quad (12)$$

5.3 Filtering and Refinement

If the value of spatiotemporal upper bound $C_{ST}(\tau, o).ub$ is less than θ , the search in the spatial and temporal domains terminate and all unscanned trajectories are pruned. Then we refine the fully and partly scanned trajectories in the two domains. If a trajectory τ is fully scanned in the spatial domain, we compute its exact spatial, temporal, and spatiotemporal correlations according to Equations (2), (3), and (4). If $C_{ST}(\tau, o) \geq \theta$, we store (trajectory, location) pair (τ, o) in the set of the matched pairs of o . Otherwise, trajectory τ is pruned.

If a trajectory $\tau' = \langle v_1, v_2, \dots, v_n \rangle$ is partly scanned in the spatial domain and is unscanned in the temporal domain, we estimate the temporal correlation upper bound $C_T(\tau', o).ub$ as follows.

$$C_T(\tau', o).ub = \frac{1}{k'_\tau} \left(\left\lfloor \frac{|\text{range}(v_1.t, t_n.t) \cap \text{range}(t_1.t, t_2.t - rt)|}{\tau'.sr} \right\rfloor + \left\lfloor \frac{|\text{range}(v_1.t, t_n.t) \cap \text{range}(t_3.t + rt, t_4.t)|}{\tau'.sr} \right\rfloor + 1 \right). \quad (13)$$

If a trajectory τ' is scanned in the temporal domain, its temporal correlation upper bound is defined as follows.

$$C_T(\tau', o).ub = \frac{1}{k'_\tau} \left(\left\lfloor \frac{|\text{range}(v_1.t, t_n.t) \cap o.R|}{\tau'.sr} \right\rfloor + 1 \right). \quad (14)$$

By combining the spatial correlation upper bound (Equation (8)) and temporal correlation upper bound (Equations (13) and (14)) according to Equation (4), we obtain a spatiotemporal correlation upper bound $C_{ST}(\tau', o).ub$ as follows.

$$C_{ST}(\tau', o).ub = \begin{cases} \lambda \cdot C_S(\tau', o).ub \\ + (1 - \lambda) \cdot C_T(\tau', o).ub & \text{if Case 1} \\ (1 - \lambda) \cdot C_T(\tau', o).ub & \text{if Case 2.} \end{cases} \quad (15)$$

Case 1: τ' is partly scanned in the spatial domain.

Case 2: τ' is unscanned in the spatial domain.

If the value of $C_{ST}(\tau', o).ub$ is less than that of θ , we prune trajectory τ' . Otherwise, we refine the trajectory in the spatial domain until it is fully scanned. Then we compute its exact spatiotemporal correlation and compare to θ .

5.4 Heuristic Scheduling

We propose a heuristic method to schedule the two query sources in the spatial and temporal domains (i.e., expansion center $o.p$ in the spatial domain, and expansion centers t_2 and t_3 in the temporal domain). Our target is to let more trajectories be scanned in the both domains, which is helpful to (i) reduce the number of scanned trajectories to be refined and to (ii) improve the pruning power of Equation (16).

For example, T_S is the set of scanned trajectories in the spatial domain, and T_T is the set of scanned trajectories in the temporal domain. We refine $|T_S \cup T_T|$ trajectories in total. If we are able to increase the intersection between T_S and T_T (the trajectories that have been scanned in both domains), we can reduce the total number of trajectories to be refined and can improve the query efficiency correspondingly. Moreover, we can use the spatial and temporal correlation upper bounds (Equation (16)) of the trajectories in $T_S \cap T_T$ to prune the search space, which yields better pruning than using only the spatial or the temporal upper bound (if trajectories only have been scanned in one domain).

Priority labels of the query sources in the two domains are then defined as follows. At each time, we only search the top-ranked query source (the query source has a larger value of its label) until a new query source takes its place.

$$q.l = \begin{cases} \lambda \cdot |(T_S \cup T_T) \setminus T_S| & \text{if Case 3} \\ (1 - \lambda) \cdot |(T_S \cup T_T) \setminus T_T| & \text{if Case 4} \end{cases} \quad (16)$$

Case 3: q is in the spatial domain ($q = o.p$).

Case 4: q is in the temporal domain (q is for t_2 and t_3).

Here λ and $(1 - \lambda)$ control the relative importance of the spatial and the temporal domains (Equation (4)).

5.5 Algorithm and Time Complexity

The PCol search procedure is detailed in Algorithm 3. The query arguments include a location o , a trajectory τ , and a threshold θ . The query result is returned in $\bigcup_{o \in O} M_o$. Initially, UB and the priority labels are set to 0 and M_o is set to \emptyset . If the value of $2d_c$ is less than that of $|o.R|$, we scan the timestamps in $\text{range}(t_2, t_3)$ (cf. Fig. 4b)(lines 1–5). We select the top-ranked query source q from heap H as the current-search query source, and we expand the search from q . We update the value of UB (Equation (12)). If the value of UB is less than θ , we prune all unscanned trajectories in the two domains (lines 6–11). Then we refine all scanned trajectories in the two domains. If a trajectory τ is fully scanned in the spatial domain, we compute its exact spatiotemporal correlation $C_{ST}(\tau.o)$ (Equation (4)) and compare it to θ . If $C_{ST}(\tau.o) \geq \theta$, (trajectory, location) pair (τ, o) is added in M_o . Otherwise, τ is pruned (lines 12–17). If a trajectory τ' is partly scanned, we compute its spatiotemporal upper bound $C_{ST}(\tau'.o).ub$ (Equation (16)). If $C_{ST}(\tau'.o).ub < \theta$, τ' is pruned. Otherwise, we further refine trajectory τ' and compute $C_{ST}(\tau'.o)$. If $C_{ST}(\tau'.o) \geq \theta$, (τ', o) is added to M_o . Otherwise, τ' is pruned (lines 18–27). Set M_o is stored. If q is not at the top of heap H , we update q to be the top-ranked query source (lines 28–31).

By combining the matching sets of all locations, the solution $\bigcup_{o \in O} M_o$ of the TL-Join is found (line 32).

Algorithm 3. PCol Search

Data: a set O of locations, a set T of trajectories, and a threshold θ

Result: $\bigcup_{o \in O} M_o$

```

1  $\forall o \in O (M_o \leftarrow \emptyset);$ 
2 for each location  $o$  in  $O$  do
3    $UB \leftarrow 0; \forall q \in H (q.l \leftarrow 0);$ 
4   if  $2d_c < |o.R|$  then
5     scan timestamps in  $\text{range}(t_2, t_3);$ 
6    $q \leftarrow H.top;$ 
7   while true do
8     expand( $q$ );
9     update  $UB$ ;
10    if  $UB < \theta$  then
11      prune all unscanned trajectories;
12    for each spatially fully scanned trajectory  $\tau$  do
13      compute  $C_{ST}(\tau.o);$ 
14      if  $C_{ST}(\tau.o) \geq \theta$  then
15         $M_o.add(\tau, o);$ 
16      else
17        prune  $\tau$ ;
18    for each partly scanned trajectory  $\tau'$  do
19      compute  $C_{ST}(\tau'.o).ub;$ 
20      if  $C_{ST}(\tau'.o).ub < \theta$  then
21        prune  $\tau'$ ;
22      else
23        refine  $\tau'$  and compute  $C_{ST}(\tau'.o);$ 
24        if  $C_{ST}(\tau'.o) \geq \theta$  then
25           $M_o.add(\tau', o);$ 
26        else
27          prune  $\tau'$ ;
28    store  $M_o$ ;
29    break;
30  if  $q \neq H.top$  then
31     $q \leftarrow H.top;$ 
32 return  $\bigcup_{o \in O} M_o;$ 

```

Let T'_θ denote the scanned trajectory set for the search process from each location. In the spatial domain, the time complexity is $O(|T'_\theta| |O| (|V| \log |V| + |E|))$ (the same as PSF), while in the temporal domain, the time complexity is $O(|T'_\theta| |O|)$. The time complexity of PCol is $O(|T'_\theta| |O| (|V| \log |V| + |E|)) + O(|T'_\theta| |O|) = O(|T'_\theta| |O| (|V| \log |V| + |E|))$. If the value of θ is sufficiently large, the time complexity is close to $O(|O| (|V| \log |V| + |E|))$. The time complexity of PCol is the same as that of PSF, and the advantage of PCol lies in that it has a higher pruning power and defines a smaller candidate set T'_θ .

6 EXPERIMENTAL RESULTS

We report on experiments with real and synthetic spatial data that offer insight into the properties of the developed algorithms.

6.1 Settings

We use two spatial networks, namely the Beijing Road Network (BRN) and the New York Road Network (NRN),¹⁰

10. <https://publish.illinois.edu/dbwork/open-data/>

TABLE 2
Parameter Settings

| | NRN | BRN |
|-----------------------------------|--|------------------------------------|
| Trajectory cardinality $ T $ | 1,000,000–10,000,000 / default 1,000,000 | 50,000–200,000 / default 100,000 |
| Location cardinality $ O $ | 500,000–2,000,000 / default 500,000 | 25,000–100,000 / default 50,000 |
| Average location time range $o.R$ | 1–7 hours / default 1 hour | 1–7 hours / default 1 hour |
| Coupling duration d_c | 20–40 minutes / default 25 minutes | 20–40 minutes / default 25 minutes |
| Threshold θ | 0.9–0.98 / default 0.96 | 0.9–0.98 / default 0.96 |
| Preference parameter λ | 0.1–0.9 / default 0.5 | 0.1–0.9 / default 0.5 |
| Thread count m | 24–120 / default 24 | 24–120 / default 24 |

which contain 28,342 vertices and 27,690 edges, and 95,581 vertices and 260,855 edges, respectively. The graphs are stored using adjacency lists. In BRN, we use a real taxi trajectory data set collected by the T-drive project [16], while in NRN, we use a real taxi trajectory data set from New York¹⁰. The time range of a trajectory is 1–2 hours. We use real location data in BRN, i.e., for POIs (e.g., restaurants and shopping malls), we use real locations and real time ranges (opening hours, e.g., 3–5 hours for a restaurant, 7–10 hours for a shopping mall), and for accidents, we use real locations and synthetic time ranges (e.g., 0.5–2 hours), and we use synthetic location data in NRN.

In the experiments, the index structure of PTF (cf. Section 3) and the spatial networks of PSF and PCol (when running Dijkstra’s expansion [15], cf. Sections 4 and 5) are memory resident, as the memory occupied is 34 MB and 44 MB for BRN and 42 MB and 55 MB for NRN. Trajectories and locations are also memory resident for all algorithms, and they occupy 279 MB for BRN and 2.2 GB for NRN. All algorithms are implemented in Java and run on a cluster with 10 data nodes. Each node is equipped with two Intel® Xeon® Processors E5-2620 v3 (2.4GHz) and 128GB RAM. Unless stated otherwise, experimental results are averaged over 10 independent trials using different query inputs. The main performance metrics are runtime and the number of location-trajectory pair visits. The number of location-trajectory visits is used as a metric because it reflects the number of data accesses. In multi-threaded executions, the total runtime is the maximum runtime among all individual threads.

Trajectories in T are selected randomly from the real data sets. The parameter settings are listed in Table 2. For PTF (Section 3), the best performance is achieved when the index contains 56 leaf nodes for BRN and 545 leaf nodes for NRN, and when each leaf node contains at most 8,192 (trajectory, location) pairs ($M = 8,192$) in BRN and at most 16,384 (trajectory, location) pairs ($M = 16,384$) in NRN. Compared to the equal-partition grid index [1], the performance of the balanced grid index is improved by around 20 percent. Because computing network distances online is time-consuming, we pre-compute the all-pairs shortest paths distances in the graphs (for PTF only, not for PSF and PCol). PTF,

TABLE 3
Pruning Effectiveness for TL-Join

| | PTF | PSF | PCol-w/o-h | PCol |
|-----------------------|--------|--------|------------|--------|
| Candidate ratio (BRN) | 0.98% | 0.51% | 0.07% | 0.06% |
| Pruning ratio (BRN) | 99.02% | 99.49% | 99.93% | 99.94% |
| Candidate ratio (NRN) | 0.26% | 0.08% | 0.04% | 0.03% |
| Pruning ratio (NRN) | 99.74% | 99.92% | 99.96% | 99.97% |

PSF (Section 4), and PCol (Section 5) are denoted by “PTF,” “PSF,” and “PCol” in subsequent figures. The PCol algorithm without the heuristic scheduling strategy is denoted by “PCol-w/o-h.”

6.2 Pruning Effectiveness

First, we study the pruning effectiveness of the algorithms using the default settings. The results are shown in Table 3, where the reported candidate and pruning ratios are defined as follows.

$$\text{Candidate ratio} = \frac{|C|}{|T||O|}$$

$$\text{Pruning ratio} = 1 - \text{Candidate ratio},$$

where $|C|$ is the size of the candidate set. The pruning ratio shows how many trajectory-location pairs are pruned, while the candidate ratio shows how many trajectory-location pair remains (to be processed in the next step). The candidate ratio is directly proportional to the running time. We see that the candidate ratio of PCol is only 6.1–11.5 percent of that of PTF and 12–37.5 percent of that of PSF. Further, the heuristic scheduling strategy reduces the candidate ratio by 14–25 percent.

6.3 Effect of Trajectory Cardinality $|T|$

Fig. 5 shows the effect of trajectory cardinality $|T|$ on the performance of the algorithms. Intuitively, a larger $|T|$ causes more (trajectory, location) pairs to be processed (cf. the complexity analysis in Sections 3.4, 4.3, and 5.5), meaning that the runtime and the number of (trajectory, location) pair visits are expected to increase for all algorithms. We see that

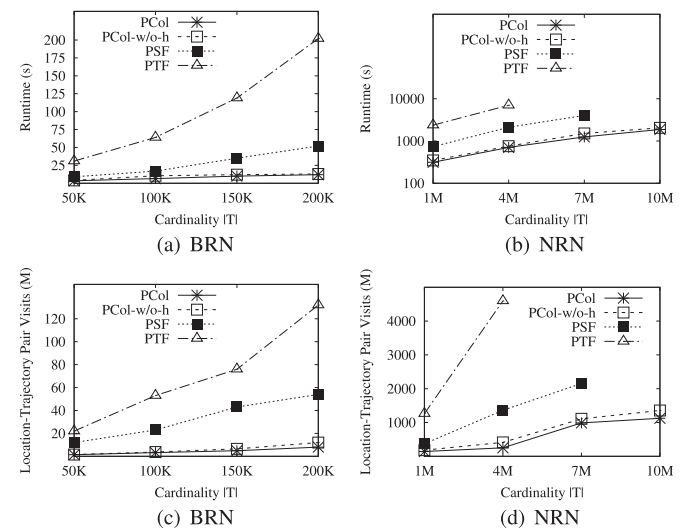
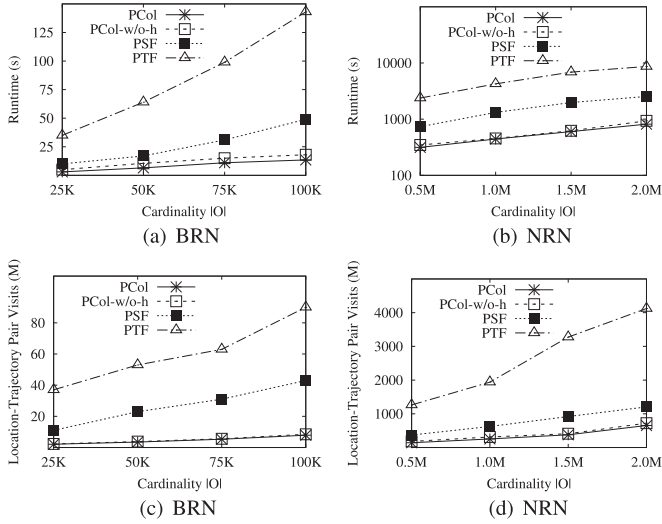


Fig. 5. Effect of trajectory cardinality $|T|$.

Fig. 6. Effect of location cardinality $|O|$.

PCol outperforms PTF by almost an order of magnitude and that it outperforms PSF by 230–300 percent in terms of both runtime and (trajectory, location) pair visits; and we see that the heuristic scheduling strategy can further improve PCol by 15–33 percent in terms of both runtime and (trajectory, location) pair visits. PCol is able to process 1 M trajectories ($|T| = 1$ M and $|O| = 0.5$ M) in 314 seconds and 10 M trajectories ($|T| = 10$ M and $|O| = 0.5$ M) in 1,874 seconds with the default 24 threads (see Fig. 5b). These results demonstrate the importance of balancing the pruning power in the spatial and temporal domains (Section 5.2) and the benefit of the heuristic scheduling strategy (Section 5.4).

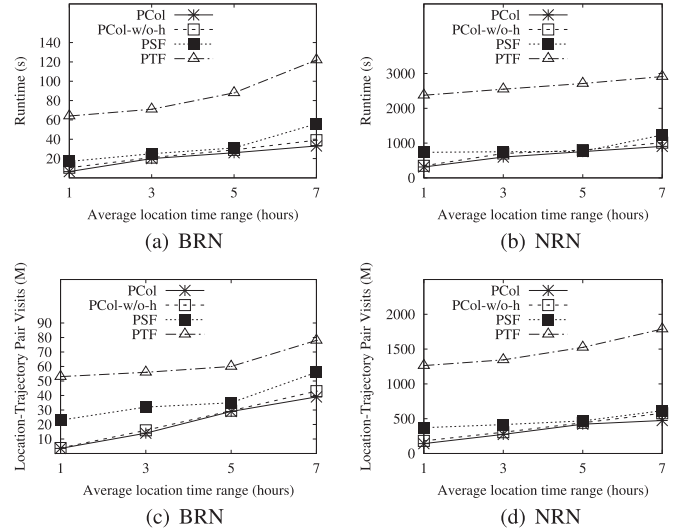
The runtime is not fully aligned with the number of (trajectory, location) pair visits because the algorithms expend computational effort on maintaining the bounds and priority labels (for PCol) used to prune the search space. The resulting cost may offset the benefits of the reduction in the number of (trajectory, location) pair visits. In particular, the filter phase of PTF computes and maintain bounds for almost all trajectory pairs.

6.4 Effect of Location Cardinality $|O|$

Next, we study the effect of location cardinality $|O|$ on the performance of the algorithms. Similar to the effect of the trajectory cardinality $|T|$, a larger $|O|$ implies a longer runtime and more (trajectory, location) pairs to be processed for all algorithms. From Fig. 6, we see that PCol has a clear advantage over PTF, PSF, and PCol-w/o-h. PCol is able to process 2 M locations ($|T| = 1$ M and $|O| = 2$ M) in 815 seconds (see Fig. 6b).

6.5 Effect of Average Location Time Range $\sigma.R$

Fig. 7 shows the effect of varying location time range $\sigma.R$ on efficiency. A larger $\sigma.R$ may lead to a higher temporal correlation (cf. Equation (3)). So we may have more qualified (trajectory, location) pairs to refine, meaning that the runtime and the number of (trajectory, location) pair visits are expected to increase for all algorithms. Moreover, for PTF, a larger $\sigma.R$ leads to more locations to be stored in non-leaf nodes, which offsets the benefit of parallel processing. For PCol, a larger $\sigma.R$ may weaken its pruning power (Equation (11)). So the runtime and the number of (trajectory, location) pair visits of

Fig. 7. Effect of average time range $\sigma.R$.

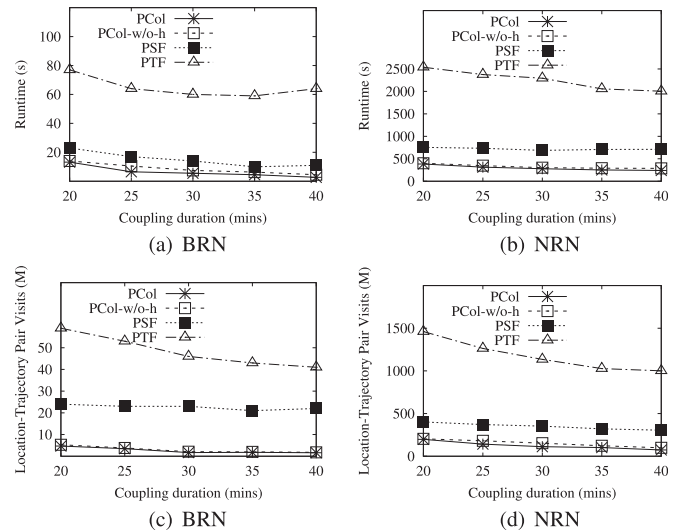
PTF and PCol increase faster than those of PSF. But PCol still holds a clear advantage over PTF, PSF, and PCol-w/o-h.

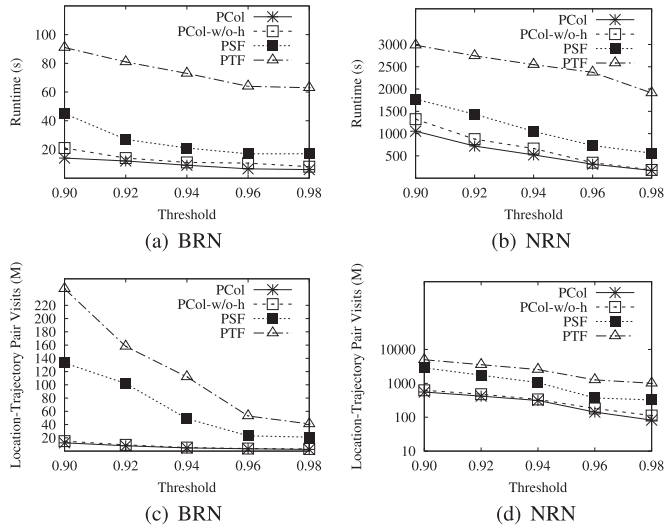
6.6 Effect of Coupling Duration d_c

The next study concerns the effect of coupling duration d_c on the efficiency of the algorithms. As can be seen in Fig. 8, a larger d_c leads to a larger k_r and may lead to a smaller temporal correlation (cf. Equation (3)). So we may have fewer qualified (trajectory, location) pairs to refine, meaning that the runtime and number of (trajectory, location) pair visits are expected to decrease for all algorithms. In addition, a larger d_c may enhance the pruning power of PCol (cf. Equations (11) and (12)).

6.7 Effect of Threshold θ

We show results when varying threshold θ in Fig. 9. A larger θ leads to higher pruning effectiveness (cf. Sections 3.3, 4.2, and 5.3). Thus, the larger θ becomes, the smaller the search space becomes. Therefore, the runtime and number of (trajectory, location) pair visits are expected to decrease correspondingly for all algorithms. In addition, in PTF, a

Fig. 8. Effect of coupling duration d_c .

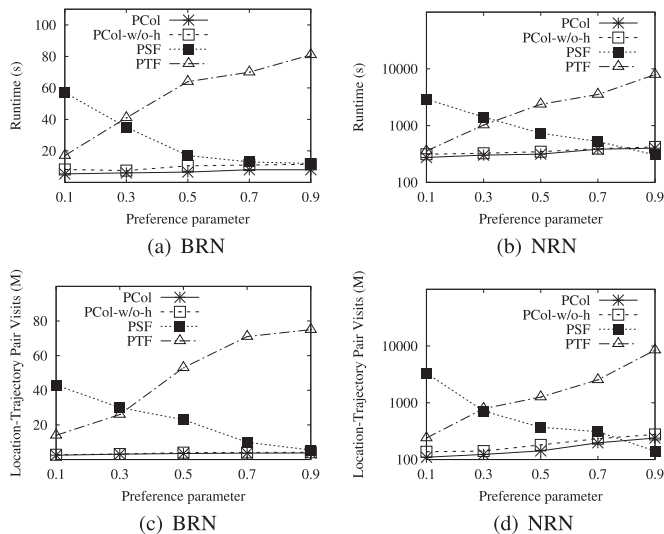
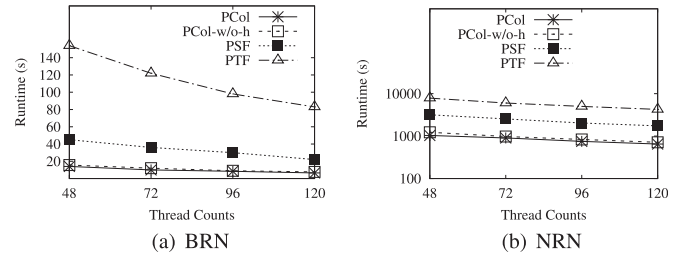
Fig. 9. Effect of threshold θ .

larger θ is useful in reducing the similarity computation (see Equation (6)), which further enhances the efficiency. In Fig. 9b, we see that when $\theta = 0.98$, PCol is able to process 1 M trajectories ($|T| = 1$ M and $|O| = 0.5$ M) in 174 seconds.

We also test that when $\theta = 0.5$ in BRN and NRN, PCol is able to process 100 K and 1 M trajectories under 240 threads in 33 seconds and 2440 seconds.

6.8 Effect of Preference Parameter λ

Fig. 10 shows the effect of varying preference parameter λ . Parameter λ enables adjusting the relative preference of spatial and temporal similarity (see Equation (4)). When $\lambda = 1$, the TL-Join is in the spatial domain only, and when $\lambda = 0$, only temporal similarity is considered. Fig. 10 shows that the spatial domain needs more search effort than the temporal domain. When λ increases, the pruning power of PTF is weakened because its pruning is driven by the temporal domain (cf. Section 3.3). On the other hand, the pruning power of PSF is enhanced as it uses spatial upper bound to prune the search space (cf. Section 4.2). When λ is close to 1, the efficiency of PSF is very close to that of PCol.

Fig. 10. Effect of preference parameter λ .Fig. 11. Effect of thread count m .

6.9 Effect of Thread Count m

We study the effect of thread count m on the efficiency of the algorithms using large trajectory data sets ($|T| = 200$ K and $|O| = 100$ K for BRN and $|P| = 10$ M and $|O| = 0.5$ M for NRN). The results are shown in Fig. 11. We see that PCol outperforms PTF by almost an order of magnitude in term of runtime and outperforms PSF by almost 300 percent in term of runtime. In BRN, PCol is able to process 200 K \times 100 K (trajectory, location) pairs with 120 threads in 8.3 seconds, while in NRN, the PCol is able to process 10 M \times 0.5 M (trajectory, location) pairs with 120 threads in 651 seconds.

We increase the thread count from 48 to 120 (2.5 times). This improves the runtimes of PSF and PCol by a factor of around 2.1, while the runtime of PTF is improved by a factor of around 1.8. The main reason for the smaller improvement of PTF is that more threads (more leaf nodes) leads to a higher merging cost (see Section 3).

7 RELATED WORK

7.1 Trajectory-to-Location Matching

Existing trajectory-to-location matching studies typically consider (i) matching solely in the spatial domain [3], [4], [5], [6], [8], [17], [18] or (ii) use point-to-point matching [2], [3], [4], [8] in the spatial or temporal domain. For the first case, the matching results do not support time-aware applications, while for the second case, the matched (trajectory, location) pairs are unable to capture the continuous correlations between trajectories and locations in the spatial and temporal domains. The so-called Semantic Enrichment approach [7] utilizes the stay time at a location to infer a traveler's activity. It uses point-to-point matching in the spatial domain and range matching in the temporal domain. This matching scheme is not feasible for location recommendation because it relies on a constraint on the stay time (e.g., 30 minutes) of travelers at a location. For example, if a traveler stay at some points of interest (e.g., restaurants, shopping malls, and sightseeing places) for more than 30 minutes, we can infer the trajectory accompanied activities (e.g., dinner, shopping, and tourist).

Trajectory-to-location matching may bring significant benefits to diverse applications. RPNN (reverse path nearest neighbor query [3]) targets the application of location ranking and recommendation. For example, when setting a new facility, RPNN uses the number of matched trajectories to define the influence factors of location candidates, and then finds the most influential location for the new facility to maximize its commercial value. ATSQ [4], UOTS [8], and PTM [2] are location-based trajectory search queries and they are useful in travel planning and carpooling

recommendation (e.g., using historic trajectories for travel planning, or recommending travelers with similar travel trajectories for carpooling). The Semantic Enrichment [7] uses trajectories to analyze traveler's activities.

Most existing centralized trajectory-to-location join algorithms (e.g., VID Joins [5], [6]) operate in euclidean space and cannot process large trajectory data sets. From the experiments reported in the literature [5], [6], the VID joins can process at most 12 K trajectories. In contrast, the TL-Join is performed in a spatial network and can process 10 M trajectories with a reasonable runtime, some three orders of magnitude more trajectories than for the VID joins.

7.2 Trajectory Similarity Join

Trajectory similarity joins [1], [13], [19], [20], [21], [22], [23] target applications such as trajectory near-duplicate detection, data cleaning, ridesharing recommendation, and traffic congestion prediction. Developing such joins typically involves a definition step and a query processing step. First, a similarity function, e.g., Sim , is defined to evaluate the spatial and temporal similarities between two trajectories, e.g., τ and τ' . Second, an efficient algorithm is developed to retrieve the spatiotemporally similar trajectory pairs. The trajectory similarity function should be symmetrical, i.e., $\text{Sim}(\tau, \tau') = \text{Sim}(\tau', \tau)$. Most existing trajectory similarity joins (e.g., [13], [19], [20], [21], [22], [23]) use a time interval threshold to constrain the temporal proximity of two trajectories. In contrast, the TS-Join [1] defines trajectory similarity in a continuous manner. The best connected trajectory (BCT) [24] and its variants [2], [4], [8] cannot be used in the trajectory similarity joins due to being asymmetric. Several similarity functions for time-series data also exist, including Dynamic Time Warping [25], Longest Common Subsequence [26], and Edit Distance on Real sequence [27].

The TS-Join [1], [28] is based on a divide-and-conquer strategy. For each trajectory τ , the algorithm retrieves trajectories that are similar to τ . The trajectory-search processes are independent of each other and are performed in parallel. The TS-Join algorithm cannot process the TL-Join due to their different query arguments (trajectories versus trajectories and locations), and their different matching functions (point-to-point matching versus range-based matching). The TL-Join needs its own specific solutions.

8 CONCLUSION AND FUTURE WORK

We studied the efficient processing of a novel Trajectory-to-Location join (TL-Join) operation in spatial networks, which may benefit diverse applications such as location recommendation, and trajectory activity analysis. We developed three parallel algorithms: parallel temporal-first search (PTF), parallel spatial-first search (PSF), and parallel collaborative search (PCol). We also defined upper and lower bounds and a heuristic scheduling strategy to enable effective search space pruning. The performance of the developed algorithms were studied empirically in extensive experiments on large spatial data sets.

Two future research directions exist. First, it is of interest to take the visiting sequence of trajectory sample points into account when matching trajectories and locations. To do this, new upper and lower bounds on the spatiotemporal

correlation and a new heuristic scheduling strategy are needed. Second, it is of interest to extend existing techniques to support a top- k TL-Join without a matching threshold θ . This calls for updated pruning techniques, including adding pruning to the same thread and updating the corresponding upper and lower bounds (without a given threshold).

REFERENCES

- [1] S. Shang, L. Chen, Z. Wei, C. S. Jensen, K. Zheng, and P. Kalnis, "Trajectory similarity join in spatial networks," *Proc. VLDB Endowment*, vol. 10, no. 11, pp. 1178–1189, 2017.
- [2] S. Shang, R. Ding, K. Zheng, C. S. Jensen, P. Kalnis, and X. Zhou, "Personalized trajectory matching in spatial networks," *VLDB J.*, vol. 23, no. 3, pp. 449–468, 2014.
- [3] S. Shang, B. Yuan, K. Deng, K. Xie, and X. Zhou, "Finding the most accessible locations: Reverse path nearest neighbor query in road networks," in *Proc. ACM 19th ACM SIGSPATIAL Int. Conf. Advances Geographic Inf. Syst.*, 2011, pp. 181–190.
- [4] K. Zheng, S. Shang, N. J. Yuan, and Y. Yang, "Towards efficient search for activity trajectories," in *Proc. IEEE 29th Int. Conf. Data Eng.*, 2013, pp. 230–241.
- [5] K. Xie, K. Deng, and X. Zhou, "From trajectories to activities: A spatio-temporal join approach," in *Proc. Int. Workshop Location Based Social Netw.*, 2009, pp. 25–32.
- [6] S. Shang, K. Xie, K. Zheng, J. Liu, and J. Wen, "VID join: Mapping trajectories to points of interest to support location-based services," *J. Comput. Sci. Technol.*, vol. 30, no. 4, pp. 725–744, 2015.
- [7] B. Furlletti, P. Cintia, C. Renso, and L. Spinsanti, "Inferring human activities from GPS tracks," in *Proc. 2nd ACM SIGKDD Int. Workshop Urban Comput.*, 2013, pp. 5:1–5:8.
- [8] S. Shang, R. Ding, B. Yuan, K. Xie, K. Zheng, and P. Kalnis, "User oriented trajectory search for trip recommendation," in *Proc. 15th Int. Conf. Extending Database Technol.*, 2012, pp. 156–167.
- [9] A. Guttman, "R-trees: A dynamic index structure for spatial searching," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 1984, pp. 47–57.
- [10] C. Wenk, R. Salas, and D. Pfoser, "Addressing the need for map-matching speed: Localizing global curve-matching algorithms," in *Proc. 18th Int. Conf. Sci. Statistical Database Manage.*, 2006, pp. 379–388.
- [11] S. Brakatsoulas, D. Pfoser, R. Salas, and C. Wenk, "On map-matching vehicle tracking data," in *Proc. 31st Int. Conf. Very Large Data Bases*, 2005, pp. 853–864.
- [12] V. P. Chakka, A. Everspaugh, and J. M. Patel, "Indexing large trajectory data sets with SETI," in *Proc. Conf. Innovative Data Syst. Res.*, 2003, pp. 1–12.
- [13] S. Ray, A. D. Brown, N. Koudas, R. Blanco, and A. K. Goel, "Parallel in-memory trajectory-based spatiotemporal topological join," in *Proc. IEEE Int. Conf. Big Data*, 2015, pp. 361–370.
- [14] V. T. de Almeida and R. H. Güting, "Indexing the trajectories of moving objects in networks," *Geoinformatica*, vol. 9, no. 1, pp. 33–60, 2005.
- [15] E. W. Dijkstra, "A note on two problems in connection with graphs," *Numerische Math.*, vol. 1, pp. 269–271, 1959.
- [16] J. Yuan, Y. Zheng, X. Xie, and G. Sun, "T-drive: Enhancing driving directions with taxi drivers' intelligence," *IEEE Trans. Knowl. Data Eng.*, vol. 25, no. 1, pp. 220–232, Jan. 2013.
- [17] S. Shang, K. Zheng, C. S. Jensen, B. Yang, P. Kalnis, G. Li, and J. Wen, "Discovery of path nearby clusters in spatial networks," *IEEE Trans. Knowl. Data Eng.*, vol. 27, no. 6, pp. 1505–1518, Jun. 2015.
- [18] S. Shang, L. Chen, C. S. Jensen, J. Wen, and P. Kalnis, "Searching trajectories by regions of interest," *IEEE Trans. Knowl. Data Eng.*, vol. 29, no. 7, pp. 1549–1562, Jul. 2017.
- [19] H. Ding, G. Trajcevski, and P. Scheuermann, "Efficient similarity join of large sets of moving object trajectories," in *Proc. 15th Int. Symp. Temporal Representation Reasoning*, 2008, pp. 79–87.
- [20] Y. Chen and J. M. Patel, "Design and evaluation of trajectory join algorithms," in *Proc. ACM SIGSPATIAL Int. Conf. Advances Geographic Inf. Syst.*, 2009, pp. 266–275.
- [21] P. Bakalov, M. Hadjieleftheriou, E. J. Keogh, and V. J. Tsotras, "Efficient trajectory joins using symbolic representations," in *Proc. 6th Int. Conf. Mobile Data Manage.*, 2005, pp. 86–93.
- [22] P. Bakalov and V. J. Tsotras, "Continuous spatiotemporal trajectory joins," in *GeoSensor Networks*, Berlin, Germany: Springer, 2006, pp. 109–128.

- [23] N. Ta, G. Li, and J. Feng, "Signature-based trajectory similarity join," *IEEE Trans. Knowl. Data Eng.*, vol. 29, no. 4, pp. 870–883, Apr. 2017.
- [24] Z. Chen, H. T. Shen, X. Zhou, Y. Zheng, and X. Xie, "Searching trajectories by locations: an efficiency study," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2010, pp. 255–266.
- [25] B. Yi, H. V. Jagadish, and C. Faloutsos, "Efficient retrieval of similar time sequences under time warping," in *Proc. 14th Int. Conf. Data Eng.*, 1998, pp. 201–208.
- [26] R. Agrawal, K. Lin, H. S. Sawhney, and K. Shim, "Fast similarity search in the presence of noise, scaling, and translation in time-series databases," in *Proc. 21th Int. Conf. Very Large Data Bases*, 1995, pp. 490–501.
- [27] L. Chen, M. T. Özsu, and V. Oria, "Robust and fast similarity search for moving object trajectories," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2005, pp. 491–502.
- [28] S. Shang, L. Chen, Z. Wei, C. S. Jensen, K. Zheng, and P. Kalnis, "Parallel trajectory similarity joins in spatial networks," *VLDB J.*, vol. 27, no. 3, pp. 395–420, 2018.

Shuo Shang received the PhD degree in computer science from the University of Queensland, Australia. He is a research scientist with KAUST. His research interests include spatiotemporal databases, spatial trajectory computing, urban computing, and location based social media. He is the track co-chair of the 2018 IEEE International Congress on Big Data and the demo co-chair of the of the 2017 APWeb/WAIM Joint Conference. He has served as a PC member, session chair, guest editor, and invited reviewer for many prestigious conferences and journals, including SIGMOD, VLDB, ICDE, CIKM, the *IEEE Transactions on Knowledge and Data Engineering*, the *VLDB Journal*, *ACM TIST*, the *IEEE Transactions on Intelligent Systems and Technology*, the *ACM Transactions on Spatial Algorithms and Systems*, *GeoInformatica*, and *WWW Journal*.

Lisi Chen received the PhD degree in computer science from Nanyang Technological University. He is a lecturer in computer science at the University of Wollongong. His research interests include geo-textual data management, spatial keyword query evaluation, and location based social networks.

Kai Zheng received the PhD degree in computer science from the University of Queensland, in 2012. He is a full professor with the University of Electronic Science and Technology of China. He has been working in the area of spatial-temporal databases, uncertain databases, social-media analysis, and in-memory databases.

Christian S. Jensen is Obel professor in computer science at Aalborg University, Denmark. He was recently at Aarhus University for three years and at Google Inc. for one year. His research concerns data management and data-intensive systems, and its focus is on temporal and spatiotemporal data management. He is an ACM and an IEEE fellow, and he is a member of the Academia Europaea, the Royal Danish Academy of Sciences and Letters, and the Danish Academy of Technical Sciences. He is editor-in-chief of the *ACM Transactions on Database Systems* and was an editor-in-chief of the *VLDB Journal* from 2008 to 2014.

Zhewei Wei received the PhD degree in computer science and engineering from the Hong Kong University of Science and Technology. He is an associate professor with the Renmin University of China. His research interests include streaming algorithms and data structures.

Panos Kalnis received the diploma degree in computer engineering from the Computer Engineering and Informatics Department, University of Patras, and the PhD degree from the HKUST. He is a professor with KAUST. His research interests include database outsourcing and cloud computing, mobile computing, and spatiotemporal and high-dimensional databases. He is an associate editor of *TKDE*, and the *VLDB Journal*.

► **For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.**