

# Efficient Contact Similarity Query over Uncertain Trajectories

Xichen Zhang, Suprio Ray, Farzaneh Shoeleh, Rongxing Lu

Canadian Institute for Cybersecurity, Faculty of Computer Science, University of New Brunswick

Fredericton, NB, Canada

{xichen.zhang,sray,farzaneh.shoeleh,rlu1}@unb.ca

## ABSTRACT

The problem of studying the effective contact between different moving objects (MOs) has received great interest in recent years. However, the uncertain nature of trajectory data poses significant challenges for the researchers. Most of the existing studies focus on range query or similarity join. But they cannot find MOs who may contact each other for a long enough time in the same location. In this paper, we study how to evaluate the effective contact among MOs. More specifically, we provide a purposeful definition for measuring the contact effectiveness, called **Contact Similarity**. Based on this notion, we introduce a novel query called **Contact Similarity Query (CSQ)**. A necklace-based model is used for representing uncertain trajectories, in which a bead (or an ellipse) indicates the possible locations of an object in-between two successive trajectory observations. To enable efficient query processing, we design a novel data structure called UTM-tree (i.e., Uncertain Trajectory M-tree) for indexing the necklace-based trajectories. Experiments have been conducted on a real-world dataset, and the results demonstrate that our proposed solution significantly outperforms the baseline approaches.

## 1 INTRODUCTION

With the increasing availability and rapid development of global positioning technologies, moving objects (MOs)' trajectories can be utilized by Location-Based Services (LBS) in many applications, such as vehicle navigation, traffic management, and co-occurrence analysis. However, due to sensor devices' physical and resource limitations or privacy considerations, MOs' trajectories are often captured at low sampling rates, and the time interval between two consecutive observations is quite long. In such uncertain courses, no information can be found about the whereabouts of MOs in-between two successive points. Recently, the problem of querying uncertain trajectory has been studied by many works, e.g., [1, 2, 5, 6, 10, 13]. However, most of them focused on retrieving qualified results regarding either probabilistic range query or similarity join. The existing approaches are not applicable to the problem of modeling the effective contact among individuals, which is a research topic of great importance. As follows, a real-world example is presented to describe a motivational scenario.

**Example 1.** Fig. 1 shows the uncertain trajectories of three visitors  $m_1$ ,  $m_2$ , and  $m_3$  in one day. They are equipped with a GPS device in their car, and they periodically reported their real-time locations along their trajectories. But we don't know where they are in-between two consecutive reports. We can see that  $m_1$  and  $m_3$  may have a very high trajectory similarity since the spatial and temporal distance between their observations can be very small. However,  $m_1$  and  $m_3$  may never come in contact with each other. On the other hand, assume that along the trip,  $m_1$  parked the car at

location A, disembarked the car, walked into the Park and stayed there for a while. Then  $m_1$  went back to A to board the car and continued with the trip. Perhaps  $m_2$  can take similar action at location B, i.e.,  $m_2$  may also stay in the Park for a while during his/her trip. Therefore, even though  $m_1$  and  $m_2$  have very low trajectory similarity, they may have an effective contact either in the Park or in the Shopping District. But such a result can never be captured by the traditional trajectory similarity query. Besides, assume that several days later,  $m_1$  found him/herself as an infected individual of a highly contagious disease (e.g., Covid 19). To help determine other susceptible persons who might be infected,  $m_1$  provided more details about his/her trajectory to an authorized third-party, such as a local public health agency (LPHA). For example,  $m_1$  stayed in the Park from 10:00 am to 10:50 am, and then went to the Shopping District from 11:35 am to 12:30 pm. Next, given this information, LPHA discovered that  $m_1$  and  $m_2$  may have possible contact with each other, and may like to know the effectiveness of their contact (e.g., the chance of their contact and how long they may contact with each other). Moreover, hundreds of people may have visited the neighboring areas on the same day. LPHA wants to find the susceptible individuals from the big group of candidates who may have effective contact with  $m_1$ .

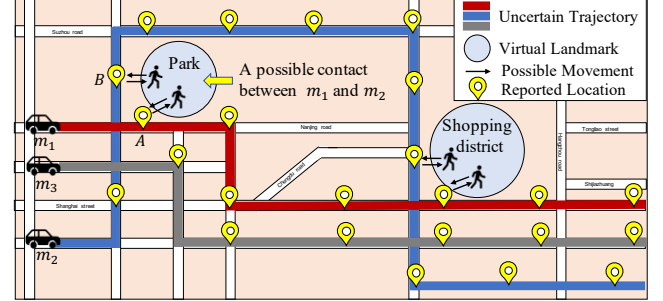


Figure 1: An real-world example of possible contacts between MOs  $m_1$ ,  $m_2$ , and  $m_3$

To address the problem described above, it is important to answer the following questions successively:

Q1: How to calculate the potential contact area and the possible longest contact time duration between  $m_1$  and  $m_2$ ?

Q2: How to measure the effective contact between  $m_1$  and  $m_2$ ?

Q3: How to efficiently retrieve the qualified results from a large number of people who visited the park on the same day as  $m_1$ ?

To the best of our knowledge, little or no systematic and theoretical study has been conducted to address the abovementioned questions. Aiming to fill this research gap, we propose a study on analyzing the effective contact between MOs over uncertain trajectories. Specifically, the main contributions of this work are three-fold as follows.

- We propose a definition of **Contact Similarity** over uncertain trajectories and a novel query called **Contact Similarity Query (CSQ)**. With CSQ, we can assess the effectiveness of contact between MOs in free space.

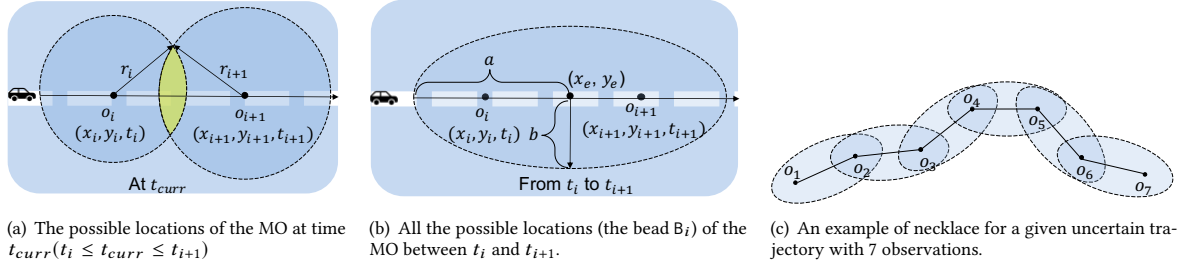


Figure 2: An example of bead  $B_i$  given the information of  $o_i, o_{i+1}$ , and  $v_{max}$ .

Table 1: The summary of notations

Notation	Definition
<b>Trajectory Notations</b>	
MO	Moving object
$T = \{o_1, o_2, \dots, o_n\}$	An uncertain trajectory
$t_i/(x_i, y_i)$	Timestamp/coordinates of $o_i$ in $T$
$v_{max}$	The maximum speed of a MO
$T_{neck}$	The necklace of $T$
$B_i$	A bead in $T_{neck}$
$P$	A virtual landmark
$t^s/t^e$	The arrival/departure time of $P$
<b>Query Notations</b>	
$R_{int}$	The intersection area ratio
$\tau(P, B_i)$	The possible longest time duration
$Sim_{cont}^{P \cap B_i}$	Contact similarity at bead level
$D$	A MO database
$Sim_{cont}^{P \cap T}$	Contact similarity at trajectory level
$\bar{t}^s$	The earliest arrival time
$\bar{t}^e$	The latest departure time

- We design a novel data structure called UTM-tree (Uncertain Trajectory M-tree) for indexing uncertain trajectories, which significantly improves the efficiency for processing CSQ.

- We conduct an extensive experimental evaluation of the proposed approaches, demonstrating the effectiveness and efficiency of the methodology proposed in this work.

The rest of this paper is organized as follows. In Section 2, we introduce the preliminaries and problem definition. In Section 3, we propose the novel indexing structure for managing uncertain trajectories. In Section 4, we present the details of our experiments and the evaluation results, followed by the related works in Section 5. Finally, we recap the conclusions in Section 6.

## 2 PRELIMINARIES

In this section, we briefly introduce the preliminaries of this work. The frequently used notations are listed in Table 1.

### 2.1 Trajectories, Beads, and Necklaces

We first introduce the definitions of trajectory, beads, and necklaces. A trajectory is defined as follows.

**Definition 1.** A *trajectory*  $T$  of a MO is defined as finite, time-ordered observations  $T = \{o_1, o_2, \dots, o_n\}$ , where  $o_i = (x_i, y_i, t_i)$  for  $i \in [1, n]$ , with  $(x_i, y_i)$  being a sample point in Euclidean space, and  $t_i$  being a timestamp.

The possible locations of a MO in-between two observations can be defined by a bead (or ellipse) [10] as follows.

**Definition 2.** Let  $v_{max}$  denote the maximum speed that an object can take between  $o_i$  and  $o_{i+1}$ . A *bead*  $B_i = \{(x_i, y_i, t_i), (x_{i+1}, y_{i+1}, t_{i+1})\}$

is defined as all points  $(x, y, t)$  that satisfy the following constraints:

$$\begin{cases} t_i \leq t \leq t_{i+1} \\ (x - x_i)^2 + (y - y_i)^2 \leq (t - t_i)^2 \times v_{max}^2 \\ (x - x_{i+1})^2 + (y - y_{i+1})^2 \leq (t_{i+1} - t)^2 \times v_{max}^2 \end{cases} \quad (1)$$

From Fig. 2 (a), we can see that from time  $t_i$  to  $t_{curr}$  ( $t_{curr} > t_i$ ), the MO's possible travel area is a circle with  $(x_i, y_i)$  being the center and  $r_i = (t_{curr} - t_i) \times v_{max}$  being the radius. Similarly, from time  $t_{curr}$  to  $t_{i+1}$  ( $t_{curr} < t_{i+1}$ ), the possible locations are included in a circle centered at  $(x_{i+1}, y_{i+1})$  with radius  $r_{i+1} = (t_{i+1} - t_{curr}) \times v_{max}$ . So, the overlapped area in Fig. 2 (a) of the two circles includes MO's possible locations at current time  $t_{curr}$ . Based on [8], all the possible locations from time  $t_i$  to  $t_{i+1}$  form an ellipse (or bead) with foci at  $(x_i, y_i)$  and  $(x_{i+1}, y_{i+1})$  (see in Fig. 2 (b)). According to [10], the equation of the bead  $B$  is:

$$(x - x_e)^2/a^2 + (y - y_e)^2/b^2 = 1, \quad (2)$$

where  $(x_e, y_e)$  is the center of the ellipse;  $a$  and  $b$  represents the semi-major axis and semi-minor axis of the ellipse, as follows.

$$\begin{aligned} x_e &= \frac{x_i + x_{i+1}}{2}, y_e = \frac{y_i + y_{i+1}}{2}, a = \frac{v_{max} \times (t_{i+1} - t_i)}{2} \\ b &= \frac{\sqrt{v_{max}^2 \times (t_{i+1} - t_i)^2 - (x_{i+1} - x_i)^2 - (y_{i+1} - y_i)^2}}{2}. \end{aligned} \quad (3)$$

**Definition 3.** A trajectory  $T = \{o_1, o_2, \dots, o_n\}$  can be represented as a sequence of beads, which is called a *necklace* and denoted as  $T_n = \{B_1, B_2, \dots, B_{n-1}\}$  (see an example in Fig. 2 (c)), such that  $o_i$  and  $o_{i+1}$  form the bead  $B_i$  for  $i \in [1, n-1]$ .

### 2.2 Contact Similarity

In this section, we introduce the essential definitions for contact similarity.

**Definition 4.** A *virtual landmark*  $P$  is a place where two MOs can contact with each other.  $P$  can be either meaningful locations (e.g., a shopping district or a park) or non-meaningful locations (e.g., a crossroad). For simplicity sake, any  $P$  in this work is considered as a circle that is centered at  $O_c = (x_c, y_c)$  with a given radius  $r_c$ . Hence, the equation of  $P$  is  $(x - x_c)^2 + (y - y_c)^2 = r_c^2$ .

**Definition 5.** For a MO  $m_1$ , he/she stayed at  $P$  for a certain period. For another MO  $m_2$ , his/her necklace is  $T_n = \{B_1, B_2, \dots, B_{n-1}\}$ . If  $P$  spatially overlaps with one bead  $B_i \in T_n$  for  $i \in [1, n-1]$  (denoted as  $P \cap B_i \neq \emptyset$ ), then the *intersection area ratio*  $R_{int}$  is

$$R_{int} = \text{Area}_{int} / \text{Area}_P, \quad (4)$$

where  $\text{Area}_{int}$  is the intersection area between the  $P$  and  $B_i$ , and  $\text{Area}_P$  is the area of  $P$ , i.e.,  $\text{Area}_P = \pi \cdot r_c^2$ .

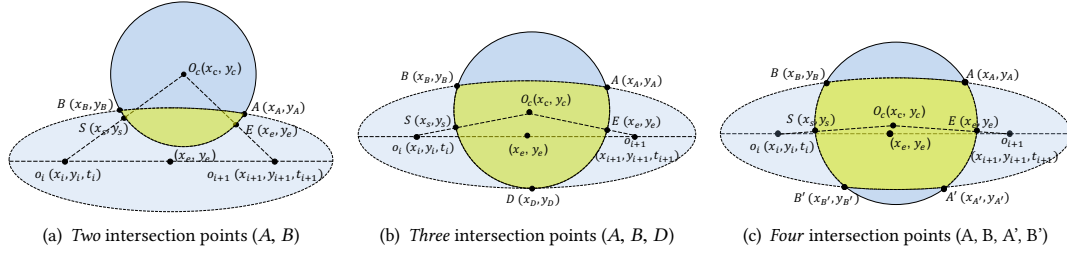


Figure 3: Examples of overlap between a virtual landmark and a bead with different intersection points

Fig. 3 shows different scenarios that P (the circle) overlaps with a bead B<sub>i</sub> (the ellipse), with two, three, and four intersection points, respectively. Then, we define the possible longest contact time duration between MOs as follows.

**Definition 6.** For MO m<sub>1</sub>, he/she stayed at P from time t<sup>s</sup> to t<sup>e</sup>. For another MO m<sub>2</sub>, his/her necklace T<sub>n</sub> = {B<sub>1</sub>, B<sub>2</sub>, ..., B<sub>n-1</sub>}. Moreover, if one of the beads B<sub>i</sub> ∈ T<sub>n</sub> (i.e., i ∈ [1, n-1]) temporally overlaps with P (denoted as P ∩<sub>T</sub> B<sub>i</sub> ≠ ∅), then the **possible longest time duration** τ(P, B<sub>i</sub>) between m<sub>1</sub> and m<sub>2</sub> can be calculated as

$$\tau(P, B_i) = \min\{t^e, \bar{t}^e\} - \max\{t^s, \bar{t}^s\}, \quad (5)$$

where  $\bar{t}^s$  means the possible earliest time for m<sub>2</sub> to arrive at P, and  $\bar{t}^e$  means the possible latest time for m<sub>2</sub> to leave P.

**Example 2.** If m<sub>1</sub> stayed at P from t<sup>s</sup> = 10:00 am and t<sup>e</sup> = 10:50 am. For m<sub>2</sub>,  $\bar{t}^s$  = 10:05 am and  $\bar{t}^e$  = 10:35 am, then τ(P, B<sub>i</sub>) = min{10:50, 10:35} - max{10:00, 10:05} = 30 min.

Next, inspired by [7], we formally propose the definition of contact similarity as below.

**Definition 7.** For m<sub>1</sub>, he/she stayed at P from time t<sup>s</sup> to t<sup>e</sup>. For m<sub>2</sub>, his/her necklace is T<sub>n</sub> = {B<sub>1</sub>, B<sub>2</sub>, ..., B<sub>n-1</sub>}. If ∃ B<sub>i</sub> ∈ T<sub>n</sub> such that P ∩<sub>S</sub> B<sub>i</sub> ≠ ∅ and P ∩<sub>T</sub> B<sub>i</sub> ≠ ∅, then the **contact similarity at bead level** Sim<sub>cont</sub><sup>P∩B<sub>i</sub></sup> between m<sub>1</sub> and m<sub>2</sub> at P in terms of B<sub>i</sub> can be calculated as:

$$\text{Sim}_{\text{cont}}^{\text{P} \cap \text{B}_i} = 1 - (1 - R_{\text{int}})^{\tau(\text{P}, \text{B}_i)} \quad (6)$$

We assume that the intersection area ratio R<sub>int</sub> indicates the possible opportunities that m<sub>1</sub> may be in contact with m<sub>2</sub> at P per unit time. Therefore, (1 - R<sub>int</sub>)<sup>τ(P, B<sub>i</sub>)</sup> denotes the probability that m<sub>1</sub> and m<sub>2</sub> have no effective contact during the time period τ(P, B<sub>i</sub>). Consequently, Sim<sub>cont</sub><sup>P∩B<sub>i</sub></sup>(m<sub>1</sub>, m<sub>2</sub>) = 1 - (1 - R<sub>int</sub>)<sup>τ(P, B<sub>i</sub>)</sup>. A larger value of Sim<sub>cont</sub><sup>P∩B<sub>i</sub></sup>(m<sub>1</sub>, m<sub>2</sub>) indicates a higher probability that an effective contact may happen between m<sub>1</sub> and m<sub>2</sub> at P.

**Definition 8.** Let  $\mathbb{B}$  denote a set of beads in T<sub>n</sub>, such that for ∀ B<sub>i</sub> ∈  $\mathbb{B}$ , P ∩<sub>S</sub> B<sub>i</sub> ≠ ∅ and P ∩<sub>T</sub> B<sub>i</sub> ≠ ∅. Then, the **contact similarity at trajectory level** Sim<sub>cont</sub><sup>P∩T</sup> can be calculated as:

$$\text{Sim}_{\text{cont}}^{\text{P} \cap \text{T}} = \frac{\sum_{B_i \in \mathbb{B}} \text{Sim}_{\text{cont}}^{\text{P} \cap \text{B}_i}}{|\mathbb{B}|}. \quad (7)$$

Here, | $\mathbb{B}$ | denotes the number of beads in  $\mathbb{B}$ .

### 2.3 Problem Definition

**Definition 9.** For m<sub>1</sub>, the **query trajectory** T<sub>q</sub> is defined as a sequence of P(s) that m<sub>1</sub> visited, i.e., T<sub>q</sub> = {P<sub>1</sub>, P<sub>2</sub>, ..., P<sub>d</sub>}. P<sub>k</sub> = (x<sub>k</sub><sup>c</sup>, y<sub>k</sub><sup>c</sup>, r<sub>k</sub><sup>c</sup>, t<sub>k</sub><sup>s</sup>, t<sub>k</sub><sup>e</sup>) for k ∈ [1, d]. Here, m<sub>1</sub> stayed at P<sub>k</sub> from t<sub>k</sub><sup>s</sup> to t<sub>k</sub><sup>e</sup>. (x<sub>k</sub><sup>c</sup>, y<sub>k</sub><sup>c</sup>) and r<sub>k</sub><sup>c</sup> are the center and radius of P<sub>k</sub>, respectively.

**Definition 10.** Given a query trajectory T<sub>q</sub> = {P<sub>1</sub>, P<sub>2</sub>, ..., P<sub>d</sub>}, a MO database D that contains s necklaces (e.g., D = {T<sub>n1</sub>, T<sub>n2</sub>, ..., T<sub>ns</sub>}), and a predefined threshold α, for each P<sub>k</sub> ∈ T<sub>q</sub> where k ∈ [1, d], the **Contact Similarity Query (CSQ)** finds all the necklaces T<sub>nj</sub> ∈ D for j ∈ [1, s], such that Sim<sub>cont</sub><sup>P<sub>k</sub> ∩ T<sub>nj</sub></sup> ≥ α.

**Example 3.** In Fig. 1, T<sub>q</sub> is m<sub>1</sub>'s query trajectory. For example, m<sub>1</sub> stayed in the park P<sub>1</sub> from 10:00 to 10:50 am, then m<sub>1</sub> went to the shopping district P<sub>2</sub> and stayed there from 11:35 am to 12:30 pm. D is a MO database that contains 1000 people's trajectory necklaces who visited the neighboring areas on the same day, and let α = 0.8. For every virtual landmark P<sub>k</sub> that m<sub>1</sub> visited, the CSQ finds all the trajectory necklaces T<sub>nj</sub> ∈ D such that Sim<sub>cont</sub><sup>P<sub>k</sub> ∩ T<sub>nj</sub></sup> ≥ 0.8.

### 2.4 Calculation of R<sub>int</sub> and τ(P, B<sub>i</sub>)

In this section, we discuss how to calculate R<sub>int</sub> and τ(P, B<sub>i</sub>), followed by the workflow for CSQ.

Given the P with equation (x - x<sub>c</sub>)<sup>2</sup> + (y - y<sub>c</sub>)<sup>2</sup> = r<sub>c</sub><sup>2</sup> and the B<sub>i</sub> with equation (x - x<sub>e</sub>)<sup>2</sup>/a<sup>2</sup> + (y - y<sub>e</sub>)<sup>2</sup>/b<sup>2</sup> = 1, it is straightforward to calculate R<sub>int</sub> as follows.

- There are two or three intersection points between the virtual landmark P and the bead B<sub>i</sub> (see Fig. 3 (a) and (b)).

$$R_{\text{int}} = \int_{x_B}^{x_A} \left( \frac{b}{a} \sqrt{a^2 - (x - x_e)^2} + y_e + \sqrt{r_c^2 - (x - x_c)^2} - y_c \right) dx \quad (8)$$

- There are four intersection points between the virtual landmark P and the bead B<sub>i</sub> (see Fig. 3 (c)).

$$R_{\text{int}} = \pi r_c^2 - \int_{x_B}^{x_A} \left( \sqrt{r_c^2 - (x - x_c)^2} + y_c - \frac{b}{a} \sqrt{a^2 - (x - x_e)^2} - y_e \right) dx - \int_{x'_B}^{x'_A} \left( y_e - \frac{b}{a} \sqrt{a^2 - (x - x_e)^2} + \sqrt{r_c^2 - (x - x_c)^2} - y_c \right) dx \quad (9)$$

In Fig. 3, we can easily know that S is the earliest point that one MO can arrive at P, and E is the latest point that MO can leave P. Therefore, given the maximum speed v<sub>max</sub>, we can calculate  $\bar{t}^s$  and  $\bar{t}^e$  as follows.

$$\bar{t}^s = t_i + \frac{|o_i S|}{v_{\text{max}}}, \quad \bar{t}^e = t_{i+1} - \frac{|o_{i+1} E|}{v_{\text{max}}} \quad (10)$$

Based on  $\bar{t}^s$  and  $\bar{t}^e$ , we can calculate τ(P, B<sub>i</sub>) by Eq. (5). After calculating R<sub>int</sub> and τ(P, B<sub>i</sub>), given an uncertain trajectory necklace database D and a trajectory query T<sub>q</sub>, we can compute Sim<sub>cont</sub><sup>P∩B<sub>i</sub></sup> and Sim<sub>cont</sub><sup>P∩T</sup> by Eq. (6) and (7), respectively. If Sim<sub>cont</sub><sup>P∩T</sup> ≥ α, the result should be returned. We repeat the above mentioned process until all the necklaces have been visited.

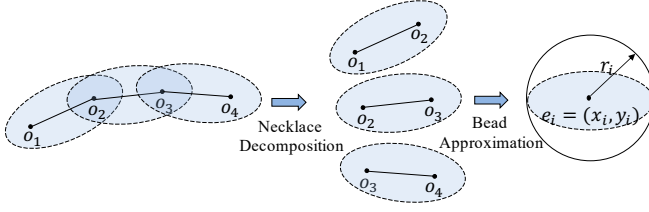


Figure 4: Decomposition of necklaces into circles

### 3 THE INDEXING STRUCTURE FOR CSQ

M-tree [4] is a popular indexing structure for efficient range query and k-nearest neighbor query. But we cannot directly adopt it in this work, because M-tree cannot handle entries like ellipse or beads. Therefore, based on M-tree, we propose a novel indexing structure called Uncertain Trajectory M-tree (UTM-tree), which is discussed next.

#### 3.1 Decomposing the Necklace

First, we decompose the necklaces into beads. After decomposition, we use a circle to approximate each bead, and the circle is centered at the center of the bead  $(x, y)$  with radius  $r = a/2$  (see in Fig. 4). Each circle is considered an individual entry and is stored independently in the leaf node of the UTM-tree.

---

#### Algorithm 1: The Insertion Algorithm: $\text{Insert}(E_i, N_j)$

---

**Input** : A new entry  $E_i = [e_i, r_i, t_i, t_{i+1}, \text{TID}, \text{Ptr}(E_p), \text{Ptr}(E_n)]$ , a tree node  $N_j = [E_j, R_j, t_{\min}, t_{\max}, \text{ptr}]$   
**Output** : A UTM-tree with  $E_i$  added

```

1 if  $N_j$  is not a leaf node then
2   for all the child nodes  $N_j^{\text{sub}}$  of  $N_j$  do
3     Let  $R_j^{\text{sub}}$  denote the covering radius of  $N_j^{\text{sub}}$ , calculate  $d(E_i, N_j^{\text{sub}})$ ;
4     if  $\exists$  at least one  $N_j^{\text{sub}}$  s.t.  $d(E_i, N_j^{\text{sub}}) + r_i \leq R_j^{\text{sub}}$  then
5       Insert( $E_i, N_j^{\text{sub}}$ ) where  $d(E_i, N_j^{\text{sub}})$  is minimum;
6     else
7       Insert( $E_i, N_j^{\text{sub}}$ ) where  $d(E_i, N_j^{\text{sub}}) + r_i - R_j^{\text{sub}}$  is minimum;
8 else
9   if  $N_j$  is not full then
10    Add  $E_i$  into  $N_j$ , update  $t_{\min}$  and  $t_{\max}$ ;
11    if  $d(E_i, N_j) + r_i > R_j$  then  $R_j = d(E_i, N_j) + r_i$ ;
12  else Split( $E_i, N_j$ );
13 return An updated UTM-tree after inserting  $E_i$  into  $N_j$ 

```

---

#### 3.2 Building the UTM-tree

Next, we describe how to construct a UTM-tree, including the format of an entry/node, entry insertion, splitting policy, selection of routing objects, and query processing. An illustration of the UTM-tree is shown in Fig. 5.

**3.2.1 Format of the entry.** Same as the M-tree, in a UTM-tree, all the entries are stored in the leaf nodes. Each node in the tree can store at most  $M$  entries, which is also called the capacity of the tree. Specifically, we consider each circle as the entry of the UTM-tree, and each entry is in format of  $E_i$ . Here,  $E_i = [e_i, r_i, t_i, t_{i+1}, \text{TID}, \text{Ptr}(E_p), \text{Ptr}(E_n)]$ , where  $e_i = (x_i, y_i)$  and  $r_i$  are the center and radius of the approximating circle,  $t_i$  and  $t_{i+1}$  are the observation timestamps from the original trajectory that form the bead, TID is the trajectory identifier which the bead belongs to, and  $\text{Ptr}(E_p)$  and  $\text{Ptr}(E_n)$  are doubly linked list

pointers for the previous/next entry that corresponding to beads in the original trajectory.

In the UTM-tree, each node selects an entry from the leaf node as its routing object (similar to M-tree). The format of a node is represented as  $N_j = [E_j, R_j, t_{\min}, t_{\max}, \text{ptr}(N_j^{\text{sub}})]$ , where  $E_j$  is its routing object,  $R_j$  is the covering radius that covers all the entries stored in  $N_j$ ,  $t_{\min}$  and  $t_{\max}$  are the minimum and maximum time for all the entries stored in  $N_j$ , and  $\text{ptr}$  is a pointer pointing to its child node  $N_j^{\text{sub}}$  if  $N_j$  is a non-leaf node.

**3.2.2 Insert an Entry.** First, we define the distance between an entry  $E_i$  and a node  $N_j$  as  $d(E_i, N_j) = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$ , where  $(x_i, y_i)$  is the center of  $E_i$ , and  $(x_j, y_j)$  is the center of the routing object of  $N_j$ . The insertion algorithm recursively descends the tree to locate the most suitable leaf node for  $E_i$ . We first find the leaf node such that after adding  $E_i$ , no enlargement of the covering radius is needed, i.e.,  $d(E_i, N_j) + r_i \leq R_j$ . If more than one child nodes are found, then the node that is closest to  $E_i$  will be chosen. However, if no such node exists, the choice is to minimize the increase of the covering radius, i.e.,  $d(E_i, N_j) + r_i - R_j$  is minimum. After insertion, if  $d(E_i, N_j) + r_i > R_j$ , then we update  $R_j$  as  $R_j = d(E_i, N_j) + r_i$ . Besides,  $t_{\min}$  and  $t_{\max}$  for all the visited nodes should be updated as well. After inserting all the entries into the UTM-tree, we use a doubly linked list to connect those entries from the same trajectory (see in Fig. 5). The details of the insertion algorithm is shown in Algorithm 1.

---

#### Algorithm 2: The Split Algorithm: $\text{Split}(E_i, N_j)$

---

**Input** : An entry  $E_i = [e_i, r_i, t_i, t_{i+1}, \text{TID}, \text{Ptr}(E_p), \text{Ptr}(E_n)]$ , a leaf node  $N_j = [E_j, R_j, t_{\min}, t_{\max}, \text{ptr}]$   
**Output** : A UTM-tree that splits  $N_j$  into  $N_j^1$  and  $N_j^2$

```

1 Let  $N_j$  be the set of all  $N_j$ 's entries plus  $E_i$ ;
2 Select two routing objects  $E_j^1$  and  $E_j^2$  from  $N_j$  based on  $m\_RAD$  algorithm [4], and then partition  $N_j$  into two sets,  $N_j^1$  and  $N_j^2$ ;
3 Allocate a new node  $N_j'$ , store  $N_j^1$  in  $N_j$  and  $N_j^2$  in  $N_j'$ ;
4 Replace  $N_j$ 's routing object  $E_j$  with  $E_j^1$ , update  $N_j$ 's  $R_j$ ;
5 Set  $E_j^2$  as  $N_j'$ 's routing object, and update  $N_j'$ 's  $R_j'$ ;
6 if  $N_j$  is not the root of the tree ( $N_j^{\text{par}}$  is  $N_j$ 's parent) then
7   if  $N_j^{\text{par}}$  is full then Split( $N_j', N_j^{\text{par}}$ );
8   else add  $N_j'$  into  $N_j^{\text{par}}$ ;
9 else Allocate a new root node  $N_r$ , set  $N_r$  as parent node for  $N_j$  and  $N_j'$ ;
10 Update  $t_{\min}$  and  $t_{\max}$  for all the visited nodes;
11 return An updated UTM-tree after splitting  $N_j$  by  $E_i$ 

```

---

**3.2.3 Split a Node.**  $E_i$  cannot be inserted into a leaf node  $N_j$  if there are  $M$  entries in  $N_j$ . We need to split  $N_j$  into two nodes. Particularly, let  $N_j$  denote the set of all  $N_j$ 's entries plus  $E_i$ . We find two routing objects  $E_j^1$  and  $E_j^2$  from  $N_j$  and partition  $N_j$  into two nodes  $N_j^1$  and  $N_j^2$  based on the  $m\_RAD$  algorithm [4]. After partitioning, the sum of the covering radius  $R_j^1 + R_j^2$  is minimum. The split algorithm is shown in Algorithm 2.

**3.2.4 Data Filtering for CSQ Query.** In this part, we discuss how to filter qualified candidates for CSQ using UTM-tree.

Given a root node  $N_r$  and a query trajectory  $T_q = \{P_1, P_2, \dots, P_d\}$  where  $P_k = (x_k^c, y_k^c, r_k^c, t_k^s, t_k^e)$  for  $k \in [1, d]$ , we can first filter all the entries  $E_i = [e_i, r_i, t_i, t_{i+1}, \text{TID}, \text{Ptr}(E_p), \text{Ptr}(E_n)]$  such that (1)  $d(e_i, (x_k^c, y_k^c)) \leq r_i + r_k^c$  and (2)  $t_i \leq t_k^s$  and  $t_{i+1} \geq t_k^e$  for all pairs of  $(N_r, P_k)$ . Next, based on the filtered entries, we can get a list of trajectory identifiers TID (s). For each trajectory  $T$  from this list, we can calculate  $\text{Sim}_{\text{cont}}^{P_k \cap T}$ , and then check if  $\text{Sim}_{\text{cont}}^{P_k \cap T} \geq \alpha$  or not. The UTM-tree indexing structure can significantly reduce



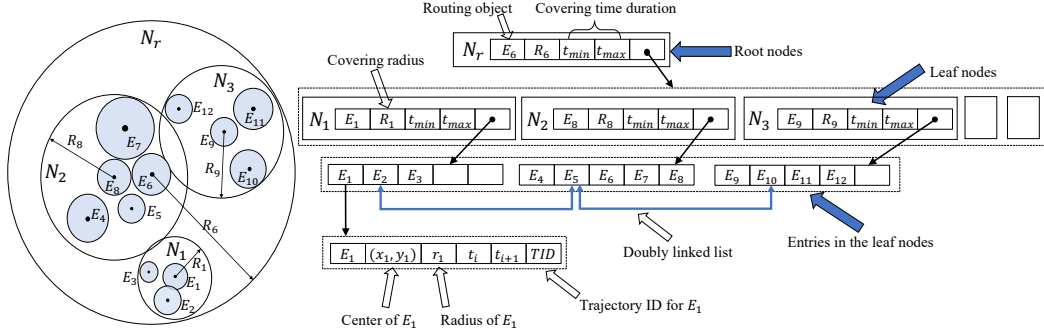


Figure 5: An example of the UTM-tree with capacity of 5

the number of trajectories to be visited, thus it can greatly increase the performance of the CSQ query. The details of the data filtering algorithm is shown in Algorithm 3.

**Algorithm 3:** The filtering algorithm:  $\text{Filter}(N_r, P_k)$

**Input** : A root node  $N_r$ , and  $P_k$  where  $P_k = (x_k^c, y_k^c, r_k^c, t_k^s, t_k^e)$   
**Output**: A list of the trajectories IDs that satisfy the query

```

1 if  $N_r$  is not a leaf node then
2   for all the children nodes  $N_r^{sub}$  in  $N_r$  where
3      $N_r^{sub} = [E_r^{sub}, R_r^{sub}, t_{min}^{sub}, t_{max}^{sub}, \text{ptr}(N_r^{sub(sub)})]$  do
4       if  $(t_{min}^{sub}, t_{max}^{sub}) \cap (t_k^s, t_k^e) \neq \emptyset$  then
5         if  $d(E_r^{sub}, (x_k^c, y_k^c)) \leq r_k^c + R_r^{sub}$  then
6           Filter( $N_r^{sub}, P_k$ );
7         else Prune  $N_r^{sub}$ ;
8   else Prune  $N_r^{sub}$ ;
9 else
10  for every  $E_i = [e_i, r_i, t_i, t_{i+1}, TID, \text{Ptr}(E_p), \text{Ptr}(E_n)]$  in  $N_r$  do
11    if  $(t_i, t_{i+1}) \cap (t_k^s, t_k^e) \neq \emptyset$  then
12      if  $d(e_i, (x_k^c, y_k^c)) \leq r_k^c + r_i$  then Add TID into an ID list;
13      else Prune  $E_i$ ;
14 else Prune  $E_i$ ;
15 return The ID list

```

## 4 EXPERIMENTAL EVALUATION

We implemented the systems with Python on an Intel(R) Core(TM) i7-6700 CPU @3.60GHz running Windows 64-bit OS with 32 GB RAM. The details of experimental settings and the performance evaluation are discussed as follows.

### 4.1 Settings

In this part, we discuss the datasets, experimental settings, and the basic approaches in this work.

- **Beijing Taxi Dataset:** This dataset [12] contains the GPS trajectories of 10357 taxis from Feb. 2 to Feb. 8, 2008 in Beijing. The total number of points in this dataset is about 15 million, and the total distance traversed by the trajectories is 9 million kilometers. The average sampling interval is 177 seconds, with an average length of 623 meters.

Table 2: The summary of simulation settings

Parameters	Settings
# of observations in one trajectory ( $\text{Num}_o$ )	30
# of trajectories ( $\text{Num}_T$ )	25K, 50K, 100K, 250K, 500K
Radius of each virtual landmark ( $r_k^c$ )	5m, 10m, 15m, 20m, 25m
# of virtual landmark in a query ( $\text{Num}_p$ )	20, 25, 30, 35, 40
Query time interval ( $\tau_q$ )	1h, 2h, 4h, 6h, 8h

- **Experimental Settings:** There are two metrics for evaluating the performance of the proposed UTM-tree, they are (1) running time and (2) the number of visited trajectories. Table 2 shows all the parameter settings in the experiment. Specifically, the number of observations in each trajectory is set to 30. For a virtual landmark  $P_k = (x_k^c, y_k^c, r_k^c, t_k^s, t_k^e) \in T_q$ , query time interval equals  $\tau_q = t_k^e - t_k^s$ . Moreover, considering a day with 24 hours,  $\tau_q$  is randomly selected from 8:00 am to 6:00 pm in the same day. For each  $P_k$ , its latitude coordinate  $x_k^c$  and longitude coordinate  $y_k^c$  are randomly selected in the city of Beijing such that  $x_k^c \in (116.1650, 116.6201)$  and  $y_k^c \in (39.7133, 40.0070)$ . Besides, the maximum speed  $v_{max}$  in-between any two continuous observations is set to a random number in  $[10 \text{ km/h}, 50 \text{ km/h}]$ .

- **Other approaches evaluated:** To evaluate the performance of the UTM-tree, we implemented two other approaches for comparison purposes, they are (1) baseline method and (2) temporal-first matching (TF-matching). In the baseline method, we do not index the trajectories at all. More specifically, we evaluate every pair of  $(B_i, P_k)$  for the CSQ query. In the TF-matching method, we follow the steps in [9] to build a temporal filtering tree. Specifically, we split a day into 12 time slots by considering every two hours as a time slot. And then, we store each trajectory into the corresponding nodes. In the TF-matching method, all the trajectories are processed by the temporal filtering tree.

### 4.2 Performance Evaluation

This section compares the proposed UTM-tree with the baseline and the TF-matching in terms of query running time and the number of visited trajectories. The number of visited trajectories means how many trajectories from the original database still need to be considered after using the UTM-tree for spatial filtering and temporal filtering.

First of all, for the three approaches, the query running time (in Fig. 6 (a)) and the number of visited trajectories (in Fig. 7 (a)) are increasing as  $\text{Num}_T$  increase. Second, we can see that the UTM-tree is much faster than the baseline method and the TF-matching method. Third, we found that for the UTM-tree, as  $\text{Num}_T$  becomes larger, the increase for the query running time and the number of visited trajectories is not significant.

Next, the comparison results for three approaches under the different  $\text{Num}_p$  are shown in both Fig. 6 (b) and Fig. 7 (b). We can see that  $\text{Num}_p$  has the least influence on the performance of the UTM-tree. In contrast, the performance of the baseline method and the TF-matching approach is degraded dramatically.

The influence of  $\tau_q$  on the three approaches under different experimental settings are shown in both Fig. 6 (c) and 7 (c). The running time and the number of visited trajectories in TF-matching

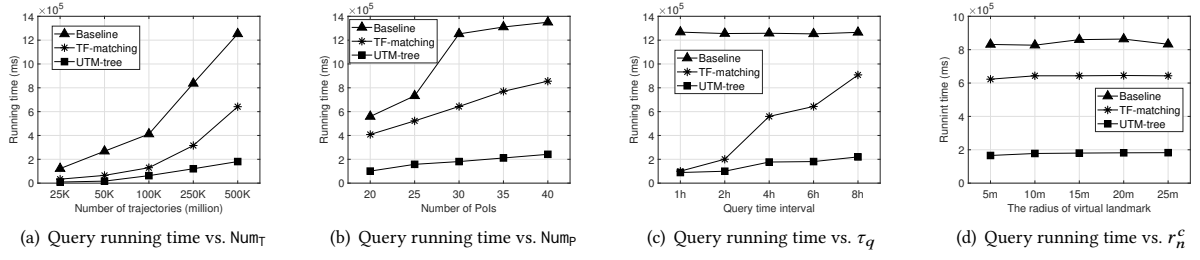


Figure 6: The comparison of baseline method, TF-matching, and UTM-tree for query running time

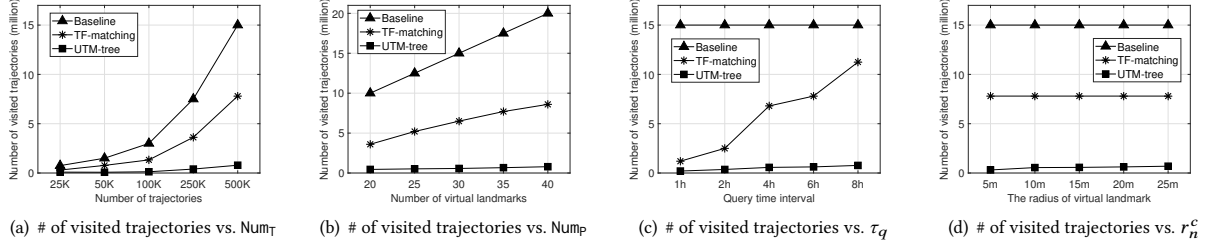


Figure 7: The comparison of baseline method, TF-matching, and UTM-tree for number of visited trajectories

approach increase significantly when  $\tau_q$  becomes wider. However, the UTM-tree is more stable against the changes in  $\tau_q$ .

Last, Fig. 6 (d) and 7 (d) compare the performance of the three approaches under different  $r_k^c$ . Both the three approaches do not change significantly with different  $r_k^c$ . Overall speaking, the UTM-tree is more efficient and stable than the baseline method and the TF-matching approach in terms of running time and the number of visited trajectories under different experiment settings. Specifically, the UTM-tree runs faster than other approaches and is less sensitive to query time intervals, virtual landmark radius, and the number of virtual landmarks in each query.

## 5 RELATED WORK

(1) *Indexing uncertain trajectories.* There are a number of indexing structures, which have been proposed for uncertain spatio-temporal data, such as UTH [13], Grid-based indexing [3], and UST-tree [5]. However, UTH-tree was used for indexing uncertain trajectories on the road network. UST-tree was used for approximating diamond-based moving objects that follow the Markov-chain model. Grid-based approach was not efficient for indexing beads and necklaces since it is time-consuming to compute the overlapped region between grids and ellipses. Therefore, we devise a novel indexing structure called UTM-tree, which is based on the classic M-tree and is efficient for indexing the uncertain trajectories in the form of beads.

(2) *Queries for uncertain trajectories.* There are many studies that were proposed for querying uncertain trajectories, such as spatio-temporal similarity join [9], semantic similarity join [3], nearest-neighbor queries [11], and top- $k$  similarity query [6]. However, most of them retrieved qualified trajectories based on some spatial/temporal criteria. None of the previous works study the problem of contact similarity in terms of spatial intersection and longest contact time duration among trajectories. This work is the first research to formally define the problem of contact similarity and propose the corresponding CSQ for the problem.

## 6 CONCLUSION

In this work, we have formally defined the concept of contact similarity and proposed a novel query, called CSQ. Next, we designed a novel indexing structure called UTM-tree, for managing and querying uncertain trajectories. Besides, we conducted extensive experiments on the Beijing Taxi dataset. The experimental results demonstrated the efficiency and stability of the UTM-tree on CSQ. There are many interesting future directions, e.g., (1) contact modeling between MOs on road networks, (2) indoor contact modeling between MOs, (3) CSQ in terms of contact frequency, (4) second-generation contact between MOs, and (5) performance evaluation with more real-world datasets.

## REFERENCES

- [1] Prithu Banerjee, Sayan Ranu, and Sriram Raghavan. 2014. Inferring uncertain trajectories from partial observations. In *2014 IEEE ICDM*. IEEE, 30–39.
- [2] Ionut Cardei, Cong Liu, Jie Wu, and Quan Yuan. 2008. DTN routing with probabilistic trajectory prediction. In *WASA*. Springer, 40–51.
- [3] Lisi Chen, Shuo Shang, Christian S Jensen, Bin Yao, and Panos Kalnis. 2020. Parallel Semantic Trajectory Similarity Join. In *2020 IEEE 36th ICDE*. 997–1008.
- [4] Paolo Ciacchia, Marco Patella, Fausto Rabitti, and Pavel Zezula. 1997. Indexing metric spaces with m-tree. In *SEBD*, Vol. 97. 67–86.
- [5] Tobias Emrich, Hans-Peter Kriegel, Nikos Mamoulis, Matthias Renz, and Andreas Züfle. 2012. Indexing uncertain spatio-temporal data. In *Proceedings of the 21st ACM CIKM*. 395–404.
- [6] C. Ma, H. Lu, L. Shou, and G. Chen. 2013. KSQ: Top-k Similarity Query on Uncertain Trajectories. *IEEE TKDE* 25, 9 (2013), 2049–2062.
- [7] Mark EJ Newman. 2002. Spread of epidemic disease on networks. *Physical review E* 66, 1 (2002), 016128.
- [8] Dieter Pfoser and Christian S Jensen. 1999. Capturing the uncertainty of moving-object representations. In *SSD*. Springer, 111–131.
- [9] Shuo Shang, Lisi Chen, Zhewei Wei, Kai Zheng, and Panos Kalnis. 2017. Trajectory similarity join in spatial networks. *VLDB Endowment* (2017).
- [10] Goce Trajcevski, Alok Choudhary, Ouri Wolfson, Li Ye, and Gang Li. 2010. Uncertain range queries for necklaces. In *2010 11th MDM*. IEEE, 199–208.
- [11] Goce Trajcevski, Roberto Tamassia, Hui Ding, Peter Scheuermann, and Isabel F Cruz. 2009. Continuous probabilistic nearest-neighbor queries for uncertain trajectories. In *Proceedings of the 12th EDBT*. 874–885.
- [12] Jing Yuan, Yu Zheng, Xing Xie, and Guangzhong Sun. 2011. Driving with knowledge from the physical world. In *Proceedings of the 17th ACM SIGKDD*. 316–324.
- [13] Kai Zheng, Goce Trajcevski, Xiaofang Zhou, and Peter Scheuermann. 2011. Probabilistic range queries for uncertain trajectories on road networks. In *Proceedings of the 14th EDBT*. 283–294.