

# Map-Matching on Low-Sampling-Rate Trajectories through Frequent Pattern Mining

Lei Yu,<sup>1</sup> Zhiqiang Zhang,<sup>2</sup> Rongtao Ding<sup>3</sup>

<sup>1</sup> School of business management, HangZhou Polytechnic, Hangzhou 314200, China.

<sup>2</sup> School of Software Engineering, Tongji University, Shanghai 2000092, China.

<sup>3</sup> School of E-commerce, Zhejiang Business College, Hangzhou 310053, China.

Correspondence should be addressed to Lei Yu; yulei@mail.hzpt.edu.cn

## Abstract

Map-matching, an important pre-processing task in many location-based services (LBS), projects each point of the global positioning system (GPS) within a trajectory dataset onto a digital map. The state-of-the-art map-matching algorithms typically employ hidden Markov model (HMM) via shortest path computation. But the computation of the shortest path might not work well on low-sampling-rate trajectory data (e.g., one GPS point every 1-5min), leading to low matching precision and high running time. To solve the problem, this paper firstly identifies frequent patterns (FPs) in historical trajectories to capture meaningful mobility behaviors, and then extracts mobile behavior criterion (MBC) of mobile users. Such a criterion generally represents the route choice of mobile users on road networks. Moreover, the temporal information within trajectory data was employed to estimate the speed of mobile users on road segments. The identified FPs, coupled with MBC and moving speed, help to improve the map-matching precision of low-sampling-rate trajectories. In addition, an FP-forest structure was proposed to index the identified FPs. The structure could greatly speed up the lookup of frequent paths for shorter running time. Further, the FP-forest structure was pruned to reduce redundancy with smaller space cost. Finally, experiments were carried out on real-world datasets. The results confirm that our FP-matching method outperforms state-of-the-arts in terms of effectiveness and efficiency.

## 1. Introduction

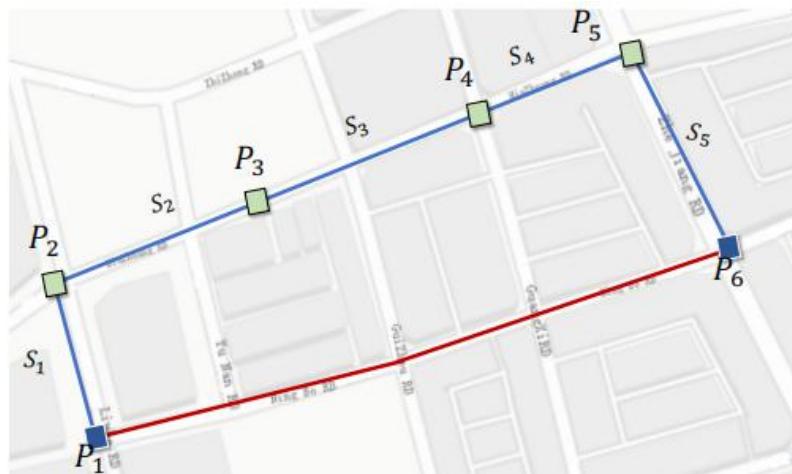
Recent years bear witness to the proliferation of location-based services (LBSs), such as Uber, WeChat and Google Map, on mobile devices. Relying on global positioning system (GPS) sensors, these LBSs record life trajectories of humans, and generate massive trajectory data. The data have been utilized in many applications to understand human mobility patterns, namely, activity recognition, hot route finder, geographical social network, and urban planning.

Nevertheless, GPS trajectories often contain lots of noisy GPS coordinates, owing to the inevitable errors of GPS positioning [1-7]. These coordinates deviate from the true positions of mobile devices, calling for the important pre-processing task of map-matching [8-15]. Map-matching aims to project every recorded GPS point within the trajectory data onto a digital map. In this way, the GPS points of users moving in the road network within the trajectory data can be corrected onto the road network, even if their coordinates deviate from the true values.

40 The frequency of the GPS points of mobile users recorded by mobile devices or third-party  
 41 applications is known as the sampling rate of GPS positions. Low-sampling-rate trajectories  
 42 with sparse GPS positions are ubiquitous. For example, when mobile users often switch off  
 43 GPS sensors to save energy or preserve privacy, their GPS trajectories will contain very  
 44 sparse GPS points with low sampling rate. LBS applications like Foursquare maintain sparse  
 45 check-in data of mobile users. In addition, telecom operators like Telco/Cellular could extract  
 46 sparse GPS coordinates within Web logs, when mobile users are using data services [16, 17].  
 47 All these sparse positions will reduce the sampling rate of the trajectories.

48 It is a challenging task to design an accurate map-matching algorithm for low-sampling-rate  
 49 trajectories with sparse positions [18-20]. The state-of-the-arts essentially adopt the hidden  
 50 Markov model (HMM) and its variants [11, 12] to perform map-matching for low-sampling-  
 51 rate GPS trajectories, under the assumption that the path between two GPS points is the  
 52 shortest path of a certain metric. But the assumption does not necessarily hold in real life,  
 53 especially in the face of low-sampling-rate trajectories.

54 Take Figure 1 for example. There is a trajectory of six GPS points, where P1 and P6 are the  
 55 source and destination, respectively. If map-matching is performed for the trajectory based on  
 56 the shortest path, the correct matching result is the blue route  $P_1 \rightarrow P_2 \rightarrow \dots \rightarrow P_5 \rightarrow P_6$ .  
 57 However, when the trajectory is so sparse as to contain only two points P1 and P6, the  
 58 shortest-path-based map-matching will converge to an incorrect result: the red route  $P_1 \rightarrow P_6$ .



59  
 60 Figure 1: An example of the failure of shortest-path-based map-matching

61 The above example shows the difficulty in choosing between the multiple optional routes  
 62 between the sparse GPS points on a trajectory. Each route could indicate the preference of  
 63 some mobile users, and could be the matching result of the trajectory. Therefore, it is  
 64 impossible to match the preference of all mobile users, using the shortest path solely based on  
 65 one metric.

66 To tackle the issue above, this paper proposes a frequent pattern (FP)-based map-matching  
 67 technique for low-sampling-rate trajectories. Firstly, the FPs were identified in third-party  
 68 data on historical trajectories. Then, the FPs that best match the input low-sampling-rate  
 69 trajectory were identified, and linked as the map-matching result of that trajectory. These FPs  
 70 clearly reflect the mobility behaviors of most mobile users among historical trajectories: each  
 71 of them is with the most familiar route, the fastest route, or the shortest route, rather than with

72 the shortest route alone. By linking the multiple FPs, the map-matching result of the input  
 73 trajectory mixes various preferences, pushing up the precision of map-matching.

74 Nowadays, almost all smartphones and vehicles are equipped with GPS sensors to record  
 75 trajectories. As a result, there are many publicly available third-party trajectories, such as  
 76 those recorded by moving taxis and Open Street Map. Some of these trajectories contain  
 77 high-sampling-rate GPS points. If FPs are mined from these trajectories, it would be highly  
 78 possible to design an accurate map-matching algorithm for low-sampling-rate trajectories.  
 79 For instance, Zhu et al. [17] lowered the positioning errors with the help of the third-party  
 80 data on historical GPS trajectories of taxis, and then realized more precise map-matching for  
 81 the low-sampling-rate trajectories of sparse locations within the big data of Telco/Cellular.

82 Apart from the mined FPs, several other meaningful factors were leveraged in this research.  
 83 One of the most important factors remains the selection of short and straight routes. Firstly,  
 84 two meaningful route features (length and transition) were chosen to define the general  
 85 mobile behavior criterion (MBC). Next, the temporal information within trajectory data was  
 86 employed to estimate the moving speed on road segments, for people tend to stay close to the  
 87 speed limit of each segment. Based on the massive trajectories moving on every segment, the  
 88 overall speed distribution was captured for different road segments. For each input estimated  
 89 speed, it is possible to infer the probability of the speed truly moving on each segment.  
 90 Finally, the probability, together with the two route features (length and transition), was  
 91 adopted to derive the best route of connected road segments, laying the basis for better  
 92 matching precision.

93 Besides the effectiveness of map-matching, the mined FPs make map-matching more  
 94 computationally efficient. Specifically, an FP-forest structure was designed to index the  
 95 identified FPs. The structure offers fast query response to look up a certain pattern used by  
 96 map-matching, allowing the proposed FP-matching algorithm to converge quickly to the  
 97 matching result. To reduce the space cost of FP-forest, the redundant road segments were  
 98 pruned from the mined FPs.

99 The main contributions of this research are as follows:

100 (1) FP-matching, a novel map-matching algorithm, was designed to capture the mobility  
 101 behaviors of most people. Whereas the state-of-the-arts use only one metric, FP-matching  
 102 effectively maps low-sampling-rate trajectories onto a digital map, using FPs mixed with  
 103 multiple metrics, general MBC, and temporal information.

104 (2) An FP-forest indexing structurer was developed to respond quickly to the lookup for a  
 105 certain FP, turning the raw GPS sequence to a road segment sequence. In addition to  
 106 matching effectiveness, our FP-matching algorithm can thereby achieve high computational  
 107 efficiency, as evidenced by short running time and low space overhead.

108 (3) The effectiveness and efficiency of FP-matching were compared with those of the state-  
 109 of-the-arts. The comparison shows that our method outperforms the contrastive methods in  
 110 both matching accuracy and running time.

111 The rest of the paper is organized as follows: Section 2 reviews the related work; Section 3  
 112 defines the problem and overviews the solution; Section 4 describes the FP-forest structure;

113 Section 5 presents the FP-matching algorithm; Section 6 verifies the proposed algorithm;  
 114 Section 7 puts forward the conclusions.

115

## 116 **2. Literature Review**

117 This section mainly reviews the related work on map-matching and FPs.

### 118 **2.1 Map-Matching**

119 Map-matching intends to locate a sequence of GPS points, consisting of timestamp, longitude,  
 120 and latitude onto a digital map. Over the past decades, extensive research has been conducted  
 121 on map-matching. The relevant works can be broadly classified into three categories:  
 122 incremental approach, global approach, and geometrical approach [21].

123 (1) Incremental approaches

124 J. Pei [13] predicted the road segment of a GPS point based on the previous prediction of that  
 125 segment, and decided which road segment should the current GPS point be mapped to, taking  
 126 into account of two factors: the projection distance from the current GPS point to its  
 127 candidate segments; the difference between the GPS heading and direction of the candidates.  
 128 C.E.White [14] evaluated the performance of four incremental algorithms, and found that all  
 129 of them have a low accuracy. The reason is that the incremental approach selects the best  
 130 candidate for each sample at the current time stamp, based on a small range of recent samples.  
 131 However, the selected candidate is not the global optimal candidate, and only works well  
 132 online.

133 (2) Global approach

134 Like the incremental approach, the global approach searches for the candidates of each GPS  
 135 sample, and then computes a weight for each candidate. Considering the entire trajectory, the  
 136 global approach calculates the aggregated weight of the candidate sequence, and looks for the  
 137 candidate with the largest weight. With the aid of the HMM, Lou and Newson et al. [11, 12]  
 138 calculated the emission probability under the normal distribution of projection distances  
 139 between GPS points and segment candidates. Lou et al. [11] also integrated spatial function  
 140 with temporal function to compute the transition probability, while Newson and Krumm [12]  
 141 adopted an exponential function of route distance and straight-line distance to deduce that  
 142 probability. However, the transition probability was calculated with a long running time, due  
 143 to the necessity to compute the shortest distance between candidates.

144 (3) Geometrical approach

145 The geometrical approach frequently searches for the optimal path on a digital map by  
 146 geometric similarity measures, e.g., the Fréchet distance [22], eliminating the need to find a  
 147 candidate set for each GPS point. For example, Alt and Brakatsoulas et al. [8, 9] looked for  
 148 the path with the minimum Fréchet distance to the GPS sequence.

149 FP-matching, our map-matching algorithm, is a global approach. This paper compares our  
 150 algorithm with the HMM algorithms proposed by Lou et al. [11], Newson and Krumm [12],

151 and Zheng et al. [15], which were proved to outshine other map-matching methods on low-  
 152 sampling-rate trajectories. Our algorithm differs from these HMM algorithms in the  
 153 following aspects:

154 First, our algorithm leverages the FPs mined from the historical trajectories moving on at  
 155 least two continuous road segments to replace the transition probability of these segments,  
 156 and maps the GPS points onto a road map through dynamic programming. Second, the FP-  
 157 forest structure was utilized to find the most frequent routes between two mapped segments,  
 158 and all the routes were connected to complement the entire trajectory. FP-matching greatly  
 159 reduces the running time by eliminating the computation of the shortest distance.

160 Zheng et al. [15] presented a history-based route inference system (HRIS) for map-matching  
 161 based on reference trajectories. Different from HRIS, our approach is grounded on the mined  
 162 FPs, which reflect the mobility patterns of those historical trajectories successfully map-  
 163 matched onto road networks. The accuracy and efficiency of our approach are both better  
 164 than those of HRIS, for the mined FPs are with map-matching result, while HRIS is still with  
 165 raw trajectories.

## 166 2.2 FP Mining

167 FPs play an important role in many data mining tasks, and contribute immensely to the  
 168 solution of practical problems [23]. Traditionally, FPs are mined in two ways: sequence  
 169 pattern mining, and association rules mining. The sequence pattern mining, represented by  
 170 example PrefixSpan [24] and generalized sequential pattern (GSP) [25, 26], takes basis on  
 171 divide-and-conquer algorithm, and draws the merits of database projection pattern and  
 172 correlation algorithms. The association rules mining, e.g., Apriori [1], FP-growth [27] and  
 173 Eclat [28], is underpinned by candidate set generation and test. For instance, sequential  
 174 pattern discovery using equivalent class (SPADE) [29] extends the Apriori algorithm into  
 175 sequential FP problem. Lin et al. [30] proposed a hybrid multilevel search algorithm to mine  
 176 long FPs. Prabamanieswari [31] invented the fuzzy-based frequent itemset mining to reduce  
 177 the number of scanning database, and demonstrated that the method is far superior to fuzzy  
 178 Apriori [32].

179 While Apriori [33] faces high time and space costs, FP-growth [27] builds a FP-tree structure  
 180 from the database, and recursively looks for frequent item sets by traversing the FP-tree  
 181 without explicitly generating redundant candidates. In our FP-matching algorithm, FP-forest  
 182 is adopted to index the FPs. Despite some similarities, our FP-forest has some significant  
 183 differences from FP-tree. First, a long sequence of road segments, i.e., the map-matching  
 184 result of the input high-sampling-rate trajectory, was split into multiple subsequences, and an  
 185 FP-forest was established on the multiple FP-trees for these subsequences. The FP-forest  
 186 speeds up the lookup of a certain sub-trajectory from a pair of source and sink inputs (road  
 187 segments). But the fast speed is realized at the cost of many redundant FP-trees. Thus, the  
 188 redundant segment identities (IDs) were pruned to save the space cost. Second, more other  
 189 items were introduced to the FP-forest, including the speed histogram on each indexed  
 190 segment, trajectory count, and mobility behavior weight  $W_m$ , to improve the effectiveness of  
 191 map-matching.

## 192 3. Problem Definition and Solution Overview

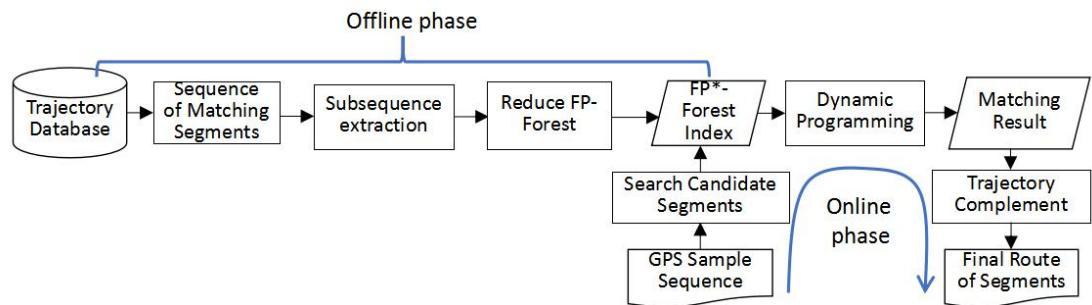
193 This section firstly defines the research problem and then highlights our solution to the  
194 problem.

195 Let  $M = (V, S)$  be a digital map containing a set of intersections (vertices)  $V$  and a set of  
196 segments (edges)  $S$ ;  $T$  be a GPS trajectory containing a sequence of GPS samples  $(p_i, t_i)$ ,  
197  $1 \leq i \leq |T|$ , each of which has a GPS position  $p_i$  and a timestamp  $t_i$ . Position  $p_i$  is a pair of  
198 GPS longitude and latitude. Map-matching aims to project the GPS points within the input  
199 trajectory  $T$  onto a sequence of consecutive segments  $s_i \in S$  in  $M$ . In the example of Figure 1,  
200 a route of connected segments  $s_1 \rightarrow \dots \rightarrow s_5$  can be generated through map-matching for the  
201 raw trajectory containing six GPS points.

202 At a very low sampling rate, the trajectory  $T$  only contains very sparse GPS points. In this  
203 case, it is highly possible to project these points onto disconnected road segments, i.e., the  
204 source and destination segments are not connected. The disconnections should be  
205 complemented with a path in the digital map  $M$ , such as to recover the entire route of  
206 connected segments.

207 To match a sparse trajectory  $T$ , our general idea is to exploit a third-party database containing  
208 historical high-sampling-rate trajectories. Such a database is widely available nowadays. Any  
209 mobile device, namely, taxis and smartphones, with GPS sensors can generate high-  
210 sampling-rate trajectories. Hence, this paper relies on the said database to optimize the map-  
211 matching of the input trajectory  $T$ .

212 Problem1. For a historical trajectory database  $D$  containing  $|D|$  high-sampling-rate trajectories  
213  $T_1 \dots T_{|D|}$ , it is necessary to identify a set of meaningful trajectory patterns from the database,  
214 and derive from these patterns the most likely route for a very sparse input trajectory  $T$  onto  
215 map  $M$ .



216

217 Figure 2: Overview of our solution

218 Figure 2 shows our solution to the above problem: a two-phase map-matching framework. In  
219 the offline phase, an FP-forest structure was maintained based on the historical trajectory  
220 database. On the high-sampling-rate trajectories in the database, a classic map-matching  
221 algorithm was applied to find the sequence of matching segments, onto which GPS points in  
222 those trajectories are projected. Next, subsequences were extracted from the sequences of  
223 matching segments. Then, all the sequences and subsequences were indexed by the FP-forest.

224 In the online phase, the digital map was traversed to find the nearest candidate segments for  
225 the input trajectory of very sparse GPS points. After that, the FPs containing these candidates

were found under the FP-forest structure, and the path that best matches the input trajectory was identified through dynamic programming. The purpose of dynamic programming is to select the best candidate for each GPS sample according to the FPs obtained via the FP-forest. Finally, the disconnected segments in the best matching path, if necessary, were complemented to obtain the entire route.

Overall, the offline phase attempts to maintain the FP-forest based on the historical database D, and the online phase seeks to evaluate the input trajectory T against the FP-forest to generate the most likely route of connected road segments.

In addition, a new FP-forest structure was developed to improve the accuracy of map-matching. To reduce the space cost, the FP-forest was downsized by an efficient method. Our algorithm applies to the new FP-forest structure, because it is similar to the original FP-forest structure.

## 4. FP-Forest

This section introduces the structure and generation of the FP-forest, details the lookup operation, presents a simple and efficient method to reduce the FP-forest, and finally reconstructs the FP-forest.

### 4.1 FP-Forest

In the FP-forest, there are multiple FP-trees rooted at the associated segment IDs. Each FP-tree is similar to the tree in the famous FP-growth structure, whose root is null. Every node inside an FP-tree contains two items: a road segment ID (SID) and a counter (CNT) of trajectories, which records the number of trajectories matching the current segment. Besides, each FP-tree has an associated dictionary, in which each key is an SID pointing to a list of tree nodes, whose SIDs are equal to the key.

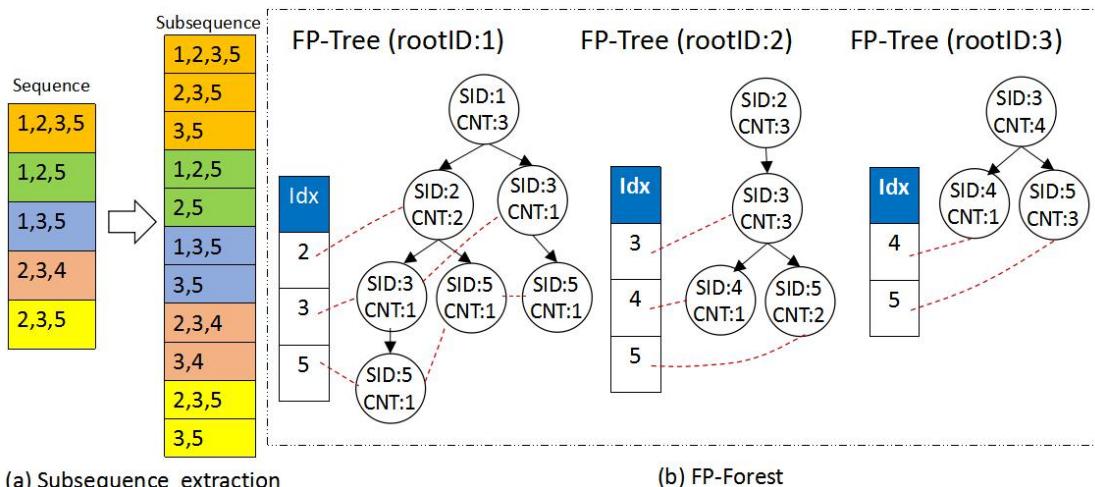


Figure 3: An example of FP-forest

The FP-forest in Figure 3 contains three FP-trees, rooted at SIDs 1, 2 and 3, respectively. For the FP-tree rooted at SID1, the dictionary maintains three keys 2, 3 and 5. Each key points to a list of internal nodes in the FP-tree. For example, key 2 points to the list of only one node

254 with the pair  $\langle 2, 2 \rangle$ , key 3 points to the list of two nodes with the pair  $\langle 3, 1 \rangle$ , and key 5  
 255 points to the list of three nodes with the pair  $\langle 5, 1 \rangle$ .

256 **4.2 Generation of FP-Forest**

257 Under the given FP-forest structure, this subsection demonstrates how to generate an FP-  
 258 forest from historical database  $D$ . The first step is to identify the high-sampling-rate  
 259 trajectories  $T_i$  in the database  $D$ . After matching trajectories  $T_i$  onto map  $M$  through classic  
 260 map-matching [19, 20], the sequences of matching segments are available: once a GPS point  
 261 in  $T_i$  is mapped onto a certain segment  $s_j \in M$ , we have a sequence of distinct segments  $s_j$   
 262 with respect to (w.r.t)  $T_j$ . After that, the FP-forest structure in Figure 3 can be called to index  
 263 all segment sequences with respect to the high-sampling-rate trajectories  $T_i$ .

264 The generation of the FP-forest can be realized in two steps: (1) extracting subsequences  
 265 from each segment sequence after map-matching each high-sampling-rate trajectory  $T_i \in D$  ;  
 266 (2) inserting the subsequences into the FP-forest structure to enable fast lookup of the forest.

267 Step 1. Subsequence extraction:

268 The following  $(k - 1)$  subsequences:  $\{s_2 \rightarrow s_3 \rightarrow \dots \rightarrow s_k\}$ ,  $\{s_3 \rightarrow \dots \rightarrow s_k\}$ , ...,  $\{s_{k-1} \rightarrow s_k\}$  are  
 269 extracted from the given sequence of segments  $S = \{s_1 \rightarrow s_2 \rightarrow \dots \rightarrow s_k\}$ . This step generates  
 270 a total of  $k$  subsequences, including the extracted subsequences plus the original sequence  $S$ .  
 271 The 11 subsequences extracted from 5 sequences of segments are presented in Figure 3.

272 Step 2. Subsequence insertion:

273 A new subsequence  $S = s_1 \rightarrow s_2 \rightarrow \dots \rightarrow s_k$  is inserted into the FP-forest structure in the  
 274 following manner. First, it is necessary to check if the FP-forest contains a tree whose root  
 275 equals SID  $s_1$ . If not, a new FP-tree is built with  $s_1$  as the root ID, with the root counter set to  
 276 1. Otherwise, the counter of this root is increased by 1.

277 After that, the rest of sequence  $S$ , e.g.,  $\{s_2 \rightarrow \dots \rightarrow s_k\}$ , is processed by one of the two  
 278 operations. One of the operations is to update the counter of an existing node in this tree, and  
 279 the other to create a new node as a child of an existing node. Take consider a subsequence  
 280  $\{s_i \rightarrow s_{i+1}\}$  for example. It is assumed that an existing node with SID  $s_i$  has been visited just  
 281 now. Then, it is necessary to visit the segment node  $s_{i+1}$ . If node  $s_i$  contains a child with  
 282 SID  $s_{i+1}$ , then the counter of the child  $s_{i+1}$  is increased by 1. Otherwise, a new child node is  
 283 created with the SID, and its counter is increased by 1.

284 The pseudocode of the above steps is given as Algorithm 1. Lines 1-14 provide the details on  
 285 how to generate an FP-forest. In Algorithm 1, each trajectory of the input database  $D_t$  is with  
 286 a sequence  $S \in D$  of connected road segments  $S = \{s_1 \dots s_{|S|}\}$ , which are obtained through  
 287 classic map-matching [19,20]. Lines 15-18 search for a specific FP-tree in FP-forest, and  
 288 Lines 19-22 look for the child of a given node.

289 The complexity of Algorithm 1 can be evaluated as follows: Suppose each trajectory  
 290 sequence  $S \in D$  contains at most  $k$  segments. Thus, the number of generated subsequences is  
 291 no greater than  $k$ , i.e., the total number of operations amounts to  $O(k)$  (i.e., `findTree`). Thus,  
 292 As shown in Figure 4 the running time of Algorithm1 is  $O(|D| \times k)$ , where  $|D|$  is the number  
 293 of trajectories in  $D$ .

294 The Conclusions section should clearly explain the main findings and implications of the  
 295 work, highlighting its importance and relevance.

296

---

**Algorithm 1:** `CreateForest`(Trajectory Database  $D$ )

---

```

Output: FP-model F
1  $F = \emptyset;$ 
2 foreach  $S \in D$  do
3   TREE  $t \leftarrow \text{FindTree}(F, S.s_1)$ ; NODE  $r = \text{null}$  ;
4   if  $t != \text{null}$  then { $r = t.\text{root}$ ;  $r.\text{counter}++$ };
5   else  $r.\text{SID} = S.s_1$ ;  $r.\text{counter} = 1$ ;  $t.\text{root} = r$ ; add  $t$  to  $F$ ;
6   NODE  $\text{curNode} = r$ ;
7   for  $i = 2$ ;  $i \leq |S|$ ;  $i + +$  do
8     NODE  $n \leftarrow \text{FindChild}(\text{curNode}, s_i)$ ;
9     if  $n != \text{null}$  then { $\text{curNode} = n$ ;  $n.\text{counter}++$ };
10    else
11       $n.\text{counter} = 1$ ;
12       $n.\text{parent} = \text{curNode}$ ;  $\text{curNode}.\text{addChild}(n)$ ;
13       $t.\text{addIndex}(n)$ ;  $\text{curNode} = n$ ;
14 return  $F$ ;
15 FindTree (FOREST  $F$ , SID  $s$ )
16   if  $F$  contains a tree whose root has the segment ID  $s$  then
17     return the tree in  $F$  whose root is the SID  $s$ ;
18   else return NULL;
19 FindChild (NODE  $\text{curNode}$ , SID  $s$ )
20   if  $\text{curNode}$  has a child whose segmentId is  $s$  then
21     return the child of  $\text{curNode}$ ;
22   else return NULL;
23 GetFreq (FOREST  $F$ , SID  $sid$ , SID  $eid$ )
24    $f = 0$ ; TREE  $t \leftarrow \text{findTree}(F, sid)$ ; PATH  $p = []$ ;
25   for NODE  $n$  in  $t.\text{index}(eid)$  do  $f \leftarrow f + n.\text{counter}$ ;
26 return  $f$ ;
```

---

297

298

Figure 4: Algorithm 1

299 **4.3 FP-Forest Lookup**

300 This subsection introduces two lookup operations on FP-forest: frequency lookup and  
 301 trajectory complementation.

302 (1) Frequency lookup

303 The goal is to find the counter (frequency) of trajectories passing from one segment sid to  
 304 another eid. To this end, the returned frequency is initiated as zero. Then, the FP-tree is  
 305 looked for, whose root with the SID equal to sid. After that, a list of tree nodes with the SIDs  
 306 equal to eid is found with the help of the dictionary. For each tree node with the SID eid, if

307 there exists a path from the root (with SID sid) to the tree node (with SID eid), the counter of  
 308 such a tree node is added to the returned frequency.

309 Here is an example for the frequency lookup from sid = 1 to eid = 5. Firstly, the FP-tree with  
 310 the root ID 1 is looked up. Next, it is assured that the associated dictionary contains key 5 in  
 311 the FP-tree. Key t points to the list of three internal nodes with the pair <5, 1,>. For each  
 312 node in the list, it is necessary to judge if there exists a path from the root 1 to the node with  
 313 ID 5. After identifies all such paths, the counters of these nodes are added up, and frequency  
 314 3 is returned.

315 Apparently, the frequency lookup is very fast, because the operation only needs to check the  
 316 connectivity of the path from the root sid to each internal node with SID eid, without any  
 317 further pruning. Note that the fast lookup is enabled by the above-mentioned subsequence  
 318 extraction: each extracted subsequence starting from sid is inserted to the FP-tree rooted at  
 319 sid. That is why the lookup only needs to check the connectivity from root sid to the node eid.  
 320 The time complexity of the frequency lookup is  $O(|Tree|)$ , where  $|Tree|$  is the number of  
 321 nodes in an FP-Tree.

## 322 (2) Trajectory complementation

323 If two segments are disconnected, it is necessary to complement a route between them, such  
 324 as to linking up all the segments in the entire route. The main idea is to select the most  
 325 suitable route, in the light of FPs, general MBC, and speed limits. For this purpose, the first  
 326 step is to find the FP-tree whose root has the SID sid. Then, the dictionary is looked up to  
 327 find the list of tree nodes whose SIDs are eid. Based on the found tree nodes, the associated  
 328 paths from the root to the tree nodes are searched for. Then, the paths are sorted by counters  
 329 and weights, and the most suitable path is returned.

330 As shown in Figure 5, the trajectory complementation is detailed in Algorithm2, whose time  
 331 complexity is  $O(|Tree| * d)$ , where  $d$  is the depth of an FP tree.

---

### Algorithm 2: Complement (FOREST F, SID sid, SID eid)

---

```

Output: BestPath
1   $f = 0$ ; TREE  $t \leftarrow \text{findTree}(F, sid)$ ;
2   $pa = []$ ;  $path = \{\}$ ;
3  for Node  $n$  in  $t.\text{index}(eid)$  do
4    while  $n \neq \text{NULL}$  do
5       $pa.append(n.SID)$ ;
6       $n \leftarrow n.parent()$ ;
7     $pa.reverse()$ ;
8     $compute w(pa)$  by Equation(2);
9     $path[w] = pa$ ;
10   return  $\max(path)$ ;

```

---

332

333

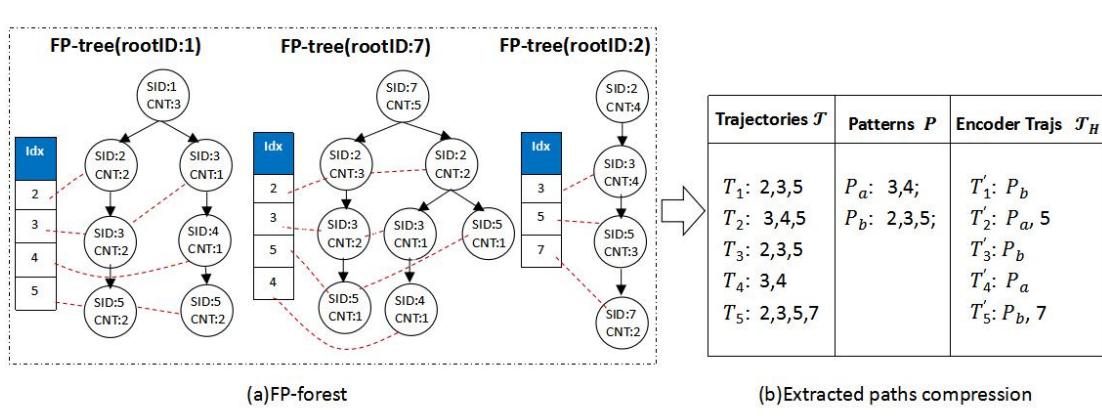
Figure 5. Algorithm 2

## 334 4.4 Redundancy Reduction

335 Recall that the FP-forest is created based on historical trajectory data. All the sub- sequences  
 336 of trajectories are indexed by an FP-forest index, which could cause road segment IDs

337 redundancy. To solve this problem, we propose a set of techniques to optimize the space cost  
 338 of FP-forest. Our basic idea is to first extract no branch paths from FP-forest, which the tree  
 339 nodes only contains one child node. Next, we mine the meaningful patterns on the extracted  
 340 sequences, and then encode a long trajectory by the patterns in order to reduce the redundant  
 341 road segment IDs, leading to smaller space cost.

342 As mentioned before, the FP-forest is generated from the data on historical trajectories. All  
 343 the subsequences of trajectories are indexed by an FP-forest index, which could cause SIDs  
 344 redundancy. To solve the problem, a set of techniques were proposed to optimize the space  
 345 cost of FP-forest. Firstly, no branch paths are extracted from FP-forest, where the tree nodes  
 346 only contain one child node each. Next, meaningful patterns are mined from the extracted  
 347 sequences, and used to encode a long trajectory. The redundant SIDs are thus reduced,  
 348 resulting in a decline in space cost.



349  
 350 Figure 6: An example of redundancy reduction

351 As shown in Figure 6, five no branch paths  $T_1 \dots T_5$  can be extracted from the FP forest in  
 352 subgraph 4(a). The paths can be viewed as map trajectories. In subgraph 4(b), the middle  
 353 column presents the two patterns  $P_a$  and  $P_b$  mined from the trajectories. Each of them,  
 354 namely,  $p_a = \{3, 4\}$ , represents 2 consecutively connected edges, that is, a path, in the road  
 355 network  $3 \rightarrow 4$ . Thus, the two edge IDs can be represented by a single virtual SID  $P_a$ ,  
 356 thereby reducing the space cost. The rightmost column depicts the encoded trajectories as  
 357 patterns. For each trajectory, namely,  $T_2$ , the redundant edge IDs  $\{3, 4\}$  can be replaced by the  
 358 segment  $P_a$ . Next is an introduction to the search for meaningful patterns.

359 The FPs of sub-trajectories indicates the redundancy of partially connected segments. Thus, it  
 360 is natural to leverage them to reduce redundancy. Inspired by Han et al. [34], a simple and  
 361 efficient approach was proposed to mine FPs under the constraint of the road network.  
 362 Considering the adjacent edges in the network, the redundant candidate patterns, i.e., those  
 363 containing disjoint edges, are removed to improve operating efficiency. Let  $M = (V, S)$  be  
 364 the road network;  $\tau = \{T_1, T_2 \dots T_N\}$  be the set of  $N$  trajectories from the mobile objects  
 365 traveling on  $M$ . The following concepts must be defined before finding the FPs.

366 Definition 1. Let  $P = \{T_i / T_i(s) \in \tau\}$  be the base pattern, with  $s$  being the segment. Then, the  
 367 trajectory containing the base pattern can be denoted as  $T_i(s)$ .

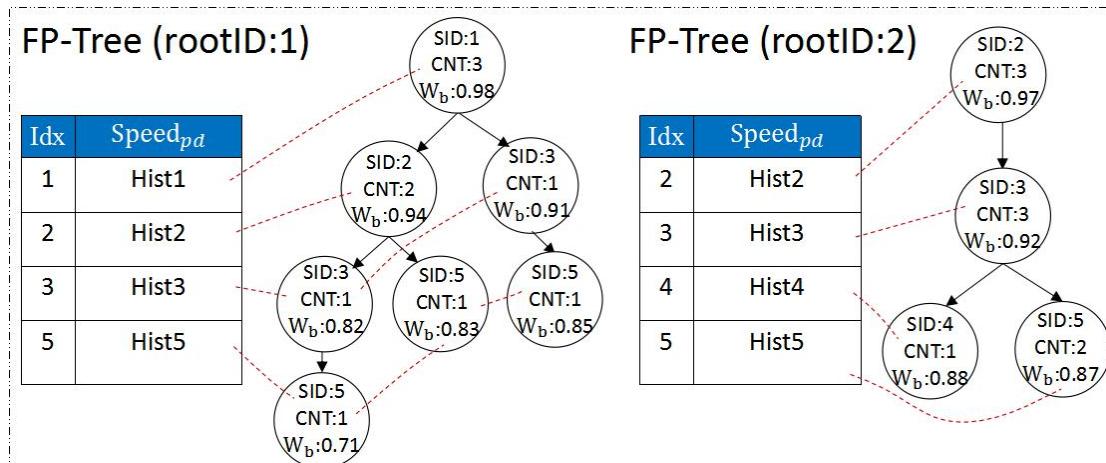
368 Definition 2. The FP, denoted by  $FP = \{P_1, P_2 \dots P_n\}$ , is an ordered list of base patterns.

369 Definition 3. The support  $\text{sup}(P)$  of pattern  $P$  is the number of the trajectories containing  
370 pattern  $P$ .

371 Then, all the frequent sub-trajectories can be found through the following steps. Firstly, the  
372 segments are organized into base patterns from the map trajectories, and the support of these  
373 base patterns (Definition 3) is calculated. Then, the base patterns are sorted in descending  
374 order of support. Thus, the generation phase of base pattern outputs an ordered list of base  
375 patterns, which serve as the building blocks of FPs in the next phase. Secondly, the base  
376 pattern with the maximum support in the said list is selected iteratively, and merged into FPs  
377 whose support is above the given lower bound. (If the base patterns are picked randomly, the  
378 merged FPs might be filtered out, due to their small supports). After that, the selected base  
379 patterns are connected into the FPs based on the road network, and all the frequency sub-  
380 trajectories are returned.

#### 381 4.5 Reconstruction of FP-Forest

382 This subsection improves the effectiveness of FP-forest structure, creating a new structure  
383 called FP-forest\*. The improvement takes account of the temporal information within  
384 trajectories, and the general mobility behavior of users. The FPs intuitively depict the spatial  
385 patterns within the trajectories. Therefore, the FP-forest\* can represent the temporal  
386 information within trajectories and the general mobility behavior. The temporal information,  
387 i.e., the timestamps within the trajectories, is taken to estimate the moving speed of the users.  
388 Meanwhile, length and transition features of each route are leveraged to infer the general  
389 mobility behavior of the users.



391 Figure 7: An example of FP-forest\*

392 As shown in Figure 7, FP-forest\* has only two differences from FP-tree: (1) Each node in  
393 FP-forest\* maintains an additional attribute, i.e., the mobility behavior weight  $W_m$  about the  
394 importance of the path from the current node to the root node; (2) The dictionary record the  
395 speed probability distribution of each referred segment. For the FP-tree\* rooted at SID 1, the  
396 node with SID 2 has a mobility behavior weight of 0.94, and the item with segment SID 2 in  
397 the dictionary exists as an equal-width histogram of moving speeds on the associated segment.  
398 In this histogram, each equal-width bucket is with a rate of those trajectories that move on the

399 associated segment within the speed interval of the bucket. For simplicity, the same number  
400 of buckets (e.g., 8 in Figure 7) is adopted for all histograms in FP-forest\*.

401 The speed histogram is maintained through the following procedure. For each segment in the  
402 third-party database D, there is a set of high-sampling-rate historical trajectories that can be  
403 projected onto that segment. Then, it is necessary to compute the mean moving speed of each  
404 trajectory passing through the segment. On this basis, a speed histogram can be comfortably  
405 maintained based on the moving speeds of all the trajectories passing through that segment.

406 Once the speed histogram is available for a segment, the probability for an input trajectory to  
407 match the segment can be derived from the histogram. In the example of Figure 7, an input  
408 trajectory is projected to the segment SID 2, and the mean speed of the trajectory moving on  
409 the segment is estimated as 6m/s. Referring to the speed histogram of SID 2, it can be  
410 observed that 6m/s falls in the bucket interval [4, 8), and the probability 0.2 of the input  
411 trajectory matches this segment.

412 Next is to estimate the weight of mobility behavior. The mobility behavior weight of each  
413 node in the FP-forest\* can be computed intuitively. Since people generally prefer short and  
414 straight routes, the mobility behavior should be weighted according to the distance and  
415 transition angle of the route. Let  $P = \{s_1 \rightarrow s_2 \rightarrow \dots \rightarrow s_k\}$  be the route from a root node to a  
416 tree node in the FP-forest\*. Then, the mobility behavior weight  $W_m$  of node  $s_k$  can be  
417 calculated by:

$$418 \quad W_m = (W_{len} + W_{turn})/2$$

$$419 \quad W_{len} = e^{-\frac{\sum_{u=1}^k s_u \cdot len}{l_{max}}} \quad (1)$$

$$W_{turn} = e^{-\frac{\sum_{u=1}^{k-1} \theta_u}{\theta_{max}}}$$

420

421 As shown in formula (1),  $W_m$  can be computed by taking the average of the weight  $W_{len}$  of  
422 route length and the weight  $W_{turn}$  of the transition angle. The weight  $W_{len}$  of route length is  
423 the sum of the length  $s_u \cdot len$  of all segments  $s_u$  within route  $P$ . The weight  $W_{turn}$  of the  
424 transition angle is the sum of the transition angles  $\theta_u$  from the preceding segments to the  
425 current segment. Minmax normalization is performed to unify the vastly different scales of  
426 the two sums. Here,  $l_{max}$  (resp.  $\theta_{max}$ ) indicates the maximal route length (resp. transition  
427 angles) from a root node to a tree node within the FP-forest\*. Following the route weights  
428 defined above, a short and straight route can be selected by reducing  $W_{len}$  and  $W_{turn}$ , i.e.,  
429 increasing  $W_m$ .

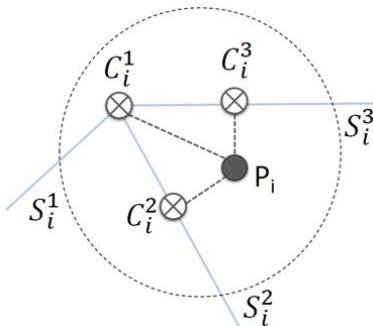
## 430 **5. FP-Matching Algorithm**

431 In this section, we will present the detail of the proposed FP-matching algorithm with two  
432 matching steps: 1) candidate search, and 2) finding a best route.

433 **5.1 Candidate Search**

434 For a given sequence of GPS points  $\rho = \{p_1 \rightarrow p_2 \rightarrow \dots \rightarrow p_k\}$ , a set of candidate segments  
 435 is searched for each point  $p_i$  by a radius  $r$  ( $1 \leq i \leq k$ ). As shown in Figure 8, GPS point  $p_i$  is  
 436 three segments  $S_i^1, S_i^2, S_i^3$  within the search radius  $r$ . Then, the distance from  $p_i$  to each  
 437 candidate  $S_i^j$  ( $1 \leq j \leq 3$ ) can be calculated by  $dist(p_i, S_i^j) = \min_{c \in S_i^j} dist(p_i, c)$ , where  $c_i^j$  is  
 438 the nearest point inside  $S_i^j$  to  $p_i$ .

439 To find candidates efficiently, an R-tree index is built for the whole digital map. Then, the  
 440 top-k nearest segment to  $p_i$  can be found in two steps: The first is to search for a set of  
 441 candidate minimal boundary rectangles (MBR) nearest to  $p_i$ . Next is to scan the segments  
 442 inside (or intersected by) the MBRs. By exploring more MBRs and segments, it is possible to  
 443 select the top-k nearest segments as candidates.

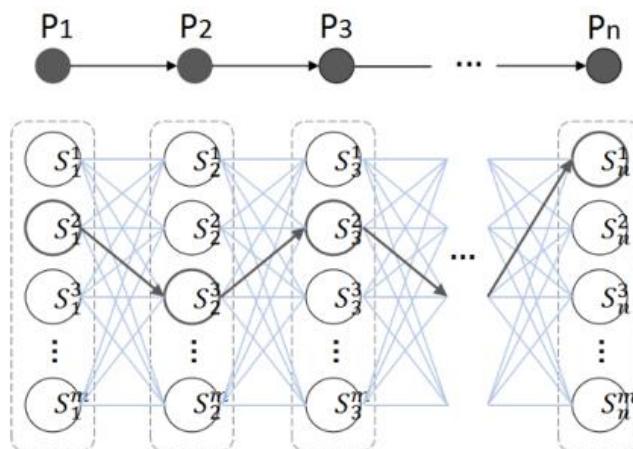


444

Figure 8: Candidate Search

446 **5.2 Route Optimization**

447 After finding the candidate segments with respect to a sequence of GPS points, it is necessary  
 448 to derive a route of candidates that best match the entire sequence of GPS points.



449

Figure 9: Route optimization

451 Note: The solid black circles are raw GPS points  $p_1 \dots p_4$ ; the hollow black circles under the  
 452 solid black circles are candidate segments.

453 Let  $S_i^j$  denote the j-th candidate segment of a GPS point  $p_i$ ;  $p_i$  and  $p_{i+1}$  be two neighboring  
 454 GPS points;  $S_i^j$  and  $S_{i+1}^{j'}$  be the candidate segments associated with the two neighboring GPS  
 455 points. Then, the two candidate segments  $S_i^j$  and  $S_{i+1}^{j'}$  might be directly connected or  
 456 disconnected on the road network. If they are disconnected, the two segments need to be  
 457 complemented by candidate patterns on FP-forest\* (Subsection 4.3). If no FP is available for  
 458 the complementation, the two segments need to be completed with the best route, under the  
 459 following weight. Up to now, all candidate segments are linked into an entire route, such as  
 460 to maximize the cumulative weight of the segments connected in that route. The route weight  
 461 between two candidate segments,  $S_i^j$  and  $S_{i+1}^{j'}$ , depends on four factors:

462 (1)  $W_d(S_i^j, S_{i+1}^{j'})$ , the weight of the distance between  $p_i$  and candidate segment  $S_i^j$  (resp. the  
 463 distance from  $S_i^j$  to  $S_{i+1}^{j'}$ ).

464 (2)  $W_t(S_i^j \rightarrow S_{i+1}^{j'})$ , the weight of the moving speed on the route between  $S_i^j$  and  $S_{i+1}^{j'}$ .

465 (3)  $W_f(S_i^j \rightarrow S_{i+1}^{j'})$ , the weight of the frequency between  $S_i^j$  and  $S_{i+1}^{j'}$ , which can be determined  
 466 through the lookup operation.

467 (4)  $W_m(S_i^j \rightarrow S_{i+1}^{j'})$ , the weight of the mobility behavior on the route between  $S_i^j$  and  $S_{i+1}^{j'}$  which  
 468 can be obtained in the FP-forest\*.

469 Whereas the previous approaches [19, 20] only consider weight  $W_d(S_i^j, S_{i+1}^{j'})$ , this paper  
 470 compute the overall route weight  $W(S_i^j \rightarrow S_{i+1}^{j'})$  between two candidates  $S_i^j$  and  $S_{i+1}^{j'}$ , which  
 471 combines all the four weights above:

$$472 \quad W(S_i^j \rightarrow S_{i+1}^{j'}) = W_f(S_i^j \rightarrow S_{i+1}^{j'}) * W_d(S_i^j, S_{i+1}^{j'}) \\ * W_r(S_i^j \rightarrow S_{i+1}^{j'}) * Wm(S_i^j \rightarrow S_{i+1}^{j'}) \quad (2)$$

473 Following Li et al.'s approach [19], the distance weight  $W_d(S_i^j, S_{i+1}^{j'})$  is computed under the  
 474 assumption that the distance between a GPS point and its real position obeys the normal  
 475 distribution  $N=(\mu, \sigma^2)$ . Thus, the probability of candidate  $S_i^j$  being the correct map match can  
 476 be estimated by:

$$477 \quad P(S_i^j) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{[Dist(i,j)-\mu]^2}{2\sigma^2}}$$

478 where,  $Dist(i,j) = Dist(p_i, S_i^j)$  is the nearest distance between  $p_i$  and its candidate segment  $S_i^j$ .  
 479 The route  $S_i^j \rightarrow S_{i+1}^{j'}$  involves two candidate segments:  $S_i^j$  and  $S_{i+1}^{j'}$ . Thus, a parameter  $\beta \geq 0$   
 480 is arranged to mediate the factors with respect to the two candidates  $S_i^j$  and  $S_{i+1}^{j'}$ , and the  
 481 following weight is defined:

482

$$W_d(S_i^j, S_{i+1}^{j'}) = \frac{(1+\beta) * P(S_i^j) * P(S_{i+1}^{j'})}{\beta * P(S_i^j) + P(S_{i+1}^{j'})} \quad (3)$$

483 where,  $\beta$  is the relative importance of  $S_i^j$  over  $S_{i+1}^{j'}$ .  $\beta = 1.0$  means the two candidates  $S_i^j$   
 484 and  $S_{i+1}^{j'}$  are equally important on the route between them;  $\beta < 1.0$  (resp.  $\beta > 1.0$ ) means  $S_i^j$  is  
 485 less (resp. more) important than  $S_{i+1}^{j'}$  on the route.

486 To compute  $W_t(S_i^j \rightarrow S_{i+1}^{j'})$ , it is necessary to estimate the mean moving speed on the route  
 487 between  $S_i^j$  and  $S_{i+1}^{j'}$ . Here, two candidate segments  $S_i^j$  and  $S_{i+1}^{j'}$  are given for two neighboring  
 488 GPS sampling points  $p_i$  and  $p_{i+1}$ , respectively. With the aid of the FP-forest\*, it is possible to  
 489 obtain the route  $P*(S_i^j, S_{i+1}^{j'})$  from  $S_i^j$  to  $S_{i+1}^{j'}$ . On this basis, the mean speed  $\bar{V}_{(i,j) \rightarrow (i+1,j')}$  of  
 490  $P*(S_i^j, S_{i+1}^{j'})$  can be calculated by:

$$\bar{V}_{(i,j) \rightarrow (i+1,j')} = \frac{w_{(i,j) \rightarrow (i+1,j')}}{\Delta t_{i,i+1}}$$

491 where,  $w_{(i,j) \rightarrow (i+1,j')}$  is the traversed distance along route  $P*(S_i^j, S_{i+1}^{j'})$ ;  $\Delta t_{i,i+1} = p_i \cdot t - p_{i+1} \cdot t$  is the  
 492 time interval between  $p_i$  and  $p_{i+1}$ . Based on the mean speed  $\bar{V}_{(i,j) \rightarrow (i+1,j')}$ , the probability  
 493  $P_{S_i^j}(\bar{V}_{(i,j) \rightarrow (i+1,j')})$  resp.  $P_{S_{i+1}^{j'}}(\bar{V}_{(i,j) \rightarrow (i+1,j')})$  of mobile devices moving on  $S_i^j$  (resp.  $S_{i+1}^{j'}$ ) with the  
 494 speed  $\bar{V}_{(i,j) \rightarrow (i+1,j')}$  can be inferred from the speed histogram with respect to the candidate  
 495 segments  $S_i^j$  (resp.  $S_{i+1}^{j'}$ ). After that, the speed weight can be calculated by:

496

$$W_t(S_i^j \rightarrow S_{i+1}^{j'}) = P_{S_i^j}(\bar{V}_{(i,j) \rightarrow (i+1,j')}) * P_{S_{i+1}^{j'}}(\bar{V}_{(i,j) \rightarrow (i+1,j')}) \quad (4)$$

497 Furthermore, the frequency weight  $W_f(S_i^j \rightarrow S_{i+1}^{j'})$  can be calculated through minmax  
 498 normalization, such that the frequency weight is linearly proportional to the frequency of  
 499 historical routes from  $S_i^j$  to  $S_{i+1}^{j'}$ , i.e.,  $Freq(S_i^j \rightarrow S_{i+1}^{j'})$ :

500

$$W_f(S_i^j \rightarrow S_{i+1}^{j'}) = \frac{Freq(S_i^j \rightarrow S_{i+1}^{j'}) - Freq_{min}}{Freq_{max} - Freq_{min}} \quad (5)$$

501 where,  $Freq_{max}$  (resp.  $Freq_{min}$ ) is the maximum (resp. minimum) frequency from a root to a  
 502 tree node within FP-forest\*.

503 The calculation of mobility behavior weight  $W_m(S_i^j \rightarrow S_{i+1}^{j'})$  is already detailed in the  
 504 preceding subsections.

505 Until now, the FP-matching problem is to find a route with the highest overall route weight  
 506 (e.g., the black line in Figure 9), such that the sum of all edge weights is maximized among  
 507 all optional routes. Formally, the best matching route  $P$  of segments can be defined as:

508

$$P = \arg \max \sum_{i=1}^{n-1} W(S_i^{best_i} \rightarrow S_{i+1}^{best_{i+1}}) \quad (6)$$

509 Our purpose is to find one route with the highest overall weight in a directed acyclic graph  
 510 (DAG). As shown in Figure 10 Algorithm 3 shows how to find the best route through  
 511 dynamic programming. After the algorithm outputs the sequence of matching segments, it is  
 512 necessary to judge whether these segments are disconnected within the sequence. If some of  
 513 them are disconnected, the full route must be recovered through the complementation  
 514 operation in Algorithm 2.

---

**Algorithm 3: FindBestRoute (DAG:  $(S_1^1 \dots S_1^k) \rightarrow \dots \rightarrow (S_n^1 \dots S_n^{k'})$ , FOREST F)**


---

**Output:** Best Route  $P: S_1^{best_1} \rightarrow \dots \rightarrow S_n^{best_n}$

```

1 Let f[] denote the route weight computed so far;
2 Let pre[] denote the previous of current candidate;
3 for i = 1; i ≤ k; i + + do f[Si1] = 0;
4 for i = 2; i ≤ n; i + + do
5   for j = 1; j ≤ k; j + + do
6     max = -1;
7     for l = 1; l ≤ k; l + + do
8       temp = f[Sli-1] + W(Sli-1 → Sji);
9       if max ≤ temp then max = temp; pre[Sji] = Sli-1;
10      f[Sji] = max;
11 P = []; cur = argmaxf[Sxn]Sxn;
12 for i = n; i ≥ 2; i - - do { P.append(cur); cur = pre[cur] };
13 P.append(cur); P.reverse();
14 return P;
```

---

515

516

Figure 10: Algorithm 3

517 Depending on the specific input DAG, the time complexity of Algorithm3 is  $O(n * m^2)$ ,  
 518 where n is the number of GPS points within the input trajectory, and m is the mean number of  
 519 candidate segments in DAG.

520 **6. Evaluation**521 **6.1 Datasets, Contrastive Algorithms, and Metrics**

522 Our experiments were carried out on two datasets. The first dataset is the digital map of  
 523 Shanghai road network, extracted from OpenStreetMap. The dataset contains 146,804  
 524 vertices and 95,950 edges. The second dataset, known as TaxiData, covers 92,602 taxi  
 525 trajectories in one day. Each GPS sample contains a timestamp, occupied/empty state, GPS  
 526 longitude and latitude, speed, and direction. For the taxi trajectories, the GPS sampling rate is  
 527 as high as 1 point per 10 seconds. The candidate segments were searched for mainly based on  
 528 GPS longitude and latitude. With help of the timestamp, GPS points were sampled from the  
 529 trajectories to simulate various sampling rates.

530 The ground truth of the experiments, i.e., the routes with high accuracy, was obtained by  
 531 applying classic map-matching algorithms to the original high-sampling-rate trajectories.

532 Based on the occupied/empty state, each long trajectory of a taxi were divided into multiple  
 533 short sub-trajectories, each of which represents one route of a passenger. The division makes  
 534 sense, because the sub-trajectories with passengers directly indicate the travel route from  
 535 sources to destinations. Then, a FP-forest was established based on the sub-trajectories. A  
 536 total of 155,725 trajectories were randomly selected to build the forest, and 337 were chosen  
 537 to evaluate the map-matching approaches. Note that historical trajectories involve GPS points  
 538 of high sampling rate; with help of the timestamp, this paper samples GPS points within  
 539 testing trajectories to simulate various sampling intervals from 1 to 5mins.

540 For comparison, our FP-matching algorithm was contrasted against three approaches: HMM-  
 541 matching [12], ST-Matching [11], and HRIS [15]. The performance of each map-matching  
 542 approach was evaluated against the following metrics about matching accuracy and running  
 543 time.

544 (1) Accuracy by number  $A_N$

545 The accuracy by number is the ratio of the number of correctly matched segments to the total  
 546 number of segments within the trajectories.

547 (2) Accuracy by length  $A_L$

548 The accuracy by length is the ratio of the length of matched segments to the total length of all  
 549 segments within the trajectories.

550 (3) Route mismatch fraction (RMF)

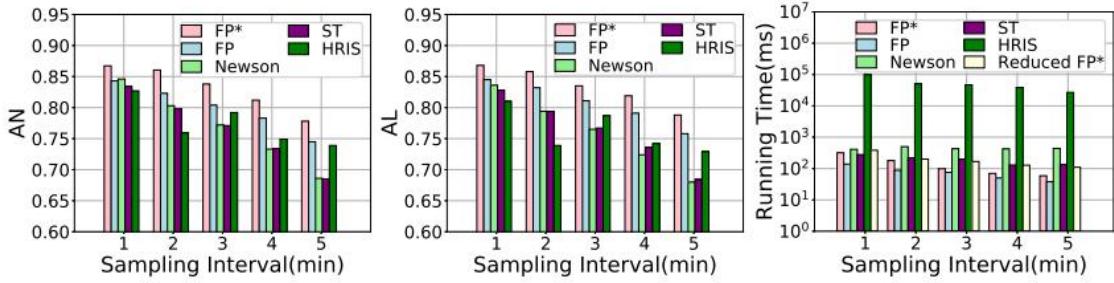
551 Let  $d_+$  and  $d$  be the total length of incorrectly matched segments added to that of the correct  
 552 route  $d_0$  and the total length of incorrectly matched segments subtracted from  $d_0$ , respectively.  
 553 Then, the summation between  $d_+$  and  $d$  is the total mismatched distance. The RMF is defined  
 554 as the fraction of  $(d_+ + d)/d_0$ , which quantifies the map-matching error, i.e.,  $RMF = (d_+ + d)/d_0$   
 555

556 (4) Minimum Fréchet distance (MFD)

557 To evaluate the correctness of the matched path, the MFD was computed between the true  
 558 path and the matched path. The smaller the MFD, the greater the similarity between the two  
 559 paths.

## 560 **6.2 Baseline Experiment**

561 The baseline experiment aims to evaluate the accuracy of five approaches, namely, FP\*-  
 562 matching, FP-matching, ST-matching [11], HMM-matching [12], and HRIS [15] at various  
 563 sampling rates. Note that FP\*-matching (resp. FP-matching) indicates the map-matching  
 564 algorithm using FP\*-forest (resp. FP-forest).



565

566

Figure 11: Effect of sampling intervals: AN, AL and running time (from left to right)

567 As shown in Figure 11 (left and middle), the accuracy of every algorithm declined with the  
568 expansion of the sampling interval, i.e., the reduction of sampling rate. Among the five  
569 approaches, FP\*-matching and FP-matching achieved the best performance, while ST-  
570 matching and HMM-matching exhibited similar trends. After the interval of 3mins, ST-  
571 matching and HMM-matching witnessed drastic degradation of accuracy, while FP\*-  
572 matching and FP-matching still performed the best. This means FP\*-matching works  
573 excellently on very sparse GPS points. In addition, HRIS operated well on sparse trajectories,  
574 when the sampling interval was greater than 3mins. This is because HRIS can exploit  
575 historical trajectories to improve the map-matching precision for the input trajectories, even if  
576 the sampling rate is very low.

577 Next, the five map-matching approaches were compared in terms of running time. As shown  
578 in Figure 11 (right), all approaches realized a shorter running time, with the expansion of the  
579 sampling interval, i.e., the increase of sparse GPS samples. Among them, HRIS had the  
580 longest running time at any sampling interval. The result is consistent with the findings of  
581 Zheng [15]. FP\*-matching generally outperformed HMM-matching [12] and ST-matching  
582 [11], for the latter two consume too much time in computing the shortest path. Besides, FP\*-  
583 matching consumed a slightly longer running time than FP-matching, for a few time is  
584 required for the redundancy reduction in FP-forest\*.

585 Table 1 compare the map-matching effectiveness of the five approaches in terms of RMF and  
586 MFD. With the increase of sampling interval, i.e., the growing number of sparse GPS  
587 samples, the RMF and MFD rose for all five approaches. Among them, FP\*-matching and  
588 FP-matching consistently outperformed the contrastive approaches.

589

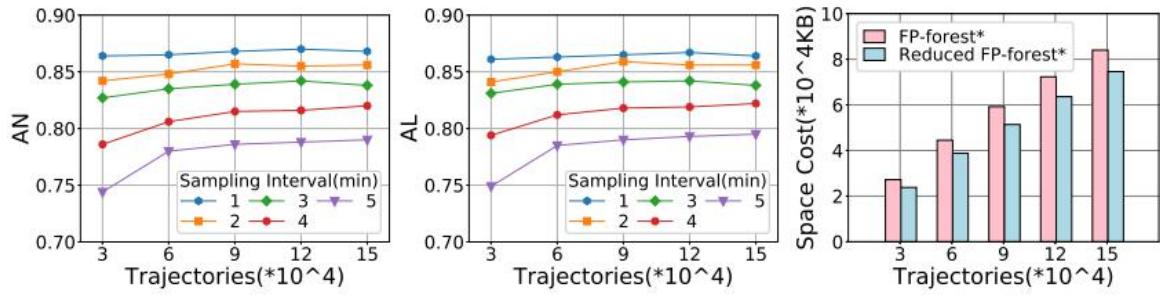
Table 1: Effect of sampling intervals from 1, 2, 3, 4 to 5 (mins)

	RMF					MFD (m)				
	1	2	3	4	5	1	2	3	4	5
HMM [20]	0.21	0.29	0.35	0.43	0.5	354.03	383.37	436.35	495.09	569.61
ST [19]	0.16	0.28	0.36	0.44	0.49	314.91	398.03	477.64	517.34	586.18
HRIS [31]	0.25	0.32	0.31	0.39	0.45	379.42	427.48	411.32	458.46	534.42
FP	0.18	0.26	0.33	0.38	0.42	296.21	368.34	390.97	451.66	513.67
FP*	0.16	0.25	0.29	0.36	0.40	285.16	349.77	375.13	416.40	473.53

590

### 6.3 Sensitivity Analysis

591 This subsection analyzes the sensitivity of FP-matching to several key parameters.



592

593 Figure 12: Effect of historical trajectories and sampling rate on map-matching accuracy and space cost (from  
594 left to right)

595

596

### 597 (1) Effect of historical data

598 The size of the historical trajectory dataset D was varied from  $3 \times 10^4$  to  $15 \times 10^4$  to test the  
599 effect of D on map-matching performance. As shown in Figure 12, when the sampling rate  
600 was fixed, the map-matching accuracy, as measured by AN and AL, increased with the  
601 number of historical trajectories; when the number of trajectories was fixed, better accuracy  
602 was achieved with the reduction of sampling interval, i.e., the growing density of GPS  
603 samples.

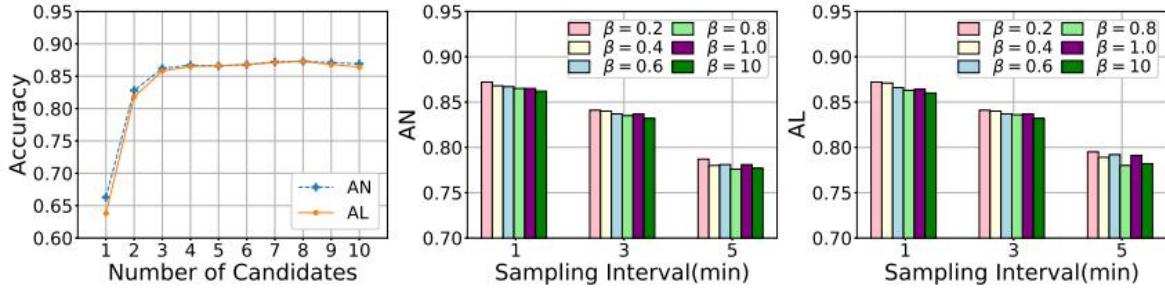
604 The authors are also interested in the space of FP-forest and the one by redundancy reduction.  
605 As shown in Figure 12 (right), when the size D grew, the space cost of FP-forest\* was clearly  
606 smaller than that of the original FP-forest.

### 607 (2) Reliability of mobility behavior weight

608 Table 2: Reliability of mobility behavior weight

Distance (m)	200	400	600	800	1000	1200
Proportion	0.999	0.975	0.952	0.924	0.899	0.852

609 In subsection 4.5, the mobility behavior weight  $W_m$  is designed for a tree node in FP-forest\*,  
610 based on people's preference for short, straight paths within a certain distance. To verify the  
611 reliability of this weight, 10,000 source-destination pairs were randomly sampled from the  
612 GPS trajectory database at different distances. For each sampled pair, the authors calculated  
613 the trajectories that traverse the shortest and straightest path as a proportion of all trajectories  
614 passing through these two points. The shortest path was measured by Euclidean distance. As  
615 shown in Table 2, the said proportion was greater than 0.92, when the two points were less  
616 than 800m apart. The proportion was still as high as 0.852, when the distance between the  
617 two points grew to 200m. Hence, the proposed mobility behavior weight is rather reliable.



618

619 Figure 13: Effect of candidate segments, AN and AL at various sampling rates and  $\beta$  (from left to right)

## 620 (3) Effect of candidate segments

621 In subsection 5.1, the candidate segments are searched for to perform map-matching. Here,  
 622 the number of candidate segments is varied to test the matching accuracy of FP\*-matching.  
 623 As shown in Figure 13, more candidates led to a higher accuracy for both AN and AL. As the  
 624 number of candidates increased from 1 to 8, the map-matching accuracy surged up to the  
 625 peak at 3 candidates, and remained stable thereafter.

626

627 (4) Effect of parameter  $\beta$ 

628 In formula (3), the road distance weight  $W_d$  employs a parameter  $\beta$  to tune the importance of  
 629 two neighboring candidate segments. As shown in Figure 13 (middle and right), when the  $\beta$   
 630 was fixed, the greater the sampling interval (the more the sparse GPS samples), the lower the  
 631 map-matching accuracy; when the sampling interval was fixed, the map-matching accuracy  
 632 increased with the decrease of  $\beta$ . These trends can be explained as follows:

633 For example, taxi drivers decide the route mainly based on the current position. That is, they  
 634 choose the next segments leading to the final destination, in view of the traffic situation  
 635 around the current position. As a result, it makes sense for FP-matching to perform better at a  
 636 smaller  $\beta$ , which leads to a higher weight for the source segment in the route  $S_i^j \rightarrow S_{i'}^{j'}$ .

637 **7. Conclusions**

638 To solve the problem of map-matching for low-sampling-rate trajectories, this paper mines  
 639 the FPs out of a huge amount of historical taxi routes, extracts the MBC, and proposes a  
 640 novel matching algorithm. The proposed algorithm has a high map-matching accuracy, as it  
 641 incorporates the weights of the distance from GPS points to road segments, moving speeds,  
 642 trajectory frequency and mobility behaviors on the segments. To speed up the map-matching,  
 643 the identified FPs were indexed by the FP-forest structure. Then, redundant paths were  
 644 pruned within the indexing structure to save space. Comparative experiments demonstrate  
 645 that our algorithm outperformed existing HMM-based and historical data-based algorithms in  
 646 terms of matching accuracy and running time. The future work plans to extend our algorithm  
 647 online, and develop a map-matching approach by the popular deep neural networks (DNNs)

## 648 Data Availability

649 The data used to support the findings of this study are available from the corresponding  
 650 author upon request.

## 651 Conflicts of Interest

652 The authors declare that they have no conflicts of interest.

## 653 Acknowledgments

654 General Scientific Research Project of Zhejiang Provincial Department of Education in 2020  
 655 (Y202043149).

## 656 References

- 657 [1] A. M. Pinto, A. P. Moreira, and P. G. Costa, “A localization method based on map-matching and particle  
 658 swarm optimization,” *Journal of Intelligent & Robotic Systems*, vol. 77, no. 2, pp. 313–326, 2015.
- 659 [2] H. Gonzalez, J. Han, X. Li, M. Myslinska and J. P. Sondag, “Adaptive fastest path computation on a road  
 660 network: A traffic mining approach,” In *Proceedings of the 33rd International Conference on Very Large  
 661 Data Bases, University of Vienna, Austria*, September 23-27, pp. 794–805, 2007.
- 662 [3] J. Guo, X. Li, Z. Zhang, and J. Zhang, “Traffic flow fluctuation analysis based on beijing taxi GPS data,”  
 663 In *Knowledge Science, Engineering and Management - 11th International Conference, KSEM 2018,*  
 664 Changchun, China, August 17-19, Proceedings, Part II, pp. 452–464, 2018.
- 665 [4] X. Li, J. Han, J. Lee, and H. Gonzalez, “Traffic density-based discovery of hot routes in road networks,” In  
 666 *Advances in Spatial and Temporal Databases, 10th International Symposium, SSTD 2007*, Boston, MA,  
 667 USA, July 16-18, Proceedings, pp. 441-459, 2007.
- 668 [5] L. Liao, D. Fox, and H. A. Kautz, “Hierarchical conditional random fields for gps-based activity  
 669 recognition,” In *Robotics Research: Results of the 12th International Symposium, ISRR 2005*, October 12-  
 670 15, San Francisco, CA, USA, pp. 487-506, 2005.
- 671 [6] K. Zhao, M. Musolesi, P. Hui, W. Rao, and S. Tarkoma, “Explaining the power-law distribution of human  
 672 mobility through transportation modality decomposition,” *Scientific reports*, vol. 5, pp. 9136, 2015.
- 673 [7] Y. Zheng, L. Wang, R. Zhang, X. Xie, and W. Ma, “Geolife: Managing and understanding your past life  
 674 over maps,” In *9th International Conference on Mobile Data Management (MDM 2008)*, Beijing, China,  
 675 April 27-30, pp. 211–212, 2008.
- 676 [8] H. Alt, A. Efrat, G. Rote, and C. Wenk, “Matching planar maps,” In *Proceedings of the Fourteenth Annual  
 677 ACM-SIAM Symposium on Discrete Algorithms*, January 12-14, 2003, Baltimore, Maryland, USA., pp.  
 678 589-598, 2003.
- 679 [9] S. Brakatsoulas, D. Pfoser, R. Salas, and C. Wenk, “On map-matching vehicle tracking data,” In  
 680 *Proceedings of the 31st International Conference on Very Large Data Bases, Trondheim, Norway*, August  
 681 30 - September 2, 2005, pp. 853–864, 2005.
- 682 [10] Y. Huang, W. Rao, Z. Zhang, “Frequent pattern-based map-matching on low sampling rate trajectories,” In  
 683 *19th IEEE International Conference on Mobile Data Management, MDM 2018*, Aalborg, Denmark, June  
 684 25-28, 2018, pp. 266–273, 2018.
- 685 [11] Y. Lou, C. Zhang, Y. Zheng, “Map-matching for low-sampling-rate GPS trajectories,” In *17th ACM  
 686 SIGSPATIAL International Symposium on Advances in Geographic Information Systems, ACM-GIS 2009*,  
 687 November 4-6, 2009, Seattle, Washington, USA, Proceedings, pp. 352–361, 2009.

- 688 [12] P. Newson and J. Krumm, "Hidden markov map matching through noise and sparseness," In *17th ACM*  
 689 *SIGSPATIAL International Symposium on Advances in Geographic Information Systems, ACM-GIS 2009*,  
 690 November 4-6, 2009, Seattle, Washington, USA, Proceedings, pp. 336–343, 2009.
- 691 [13] M. A. Quddus, W. Y. Ochieng, L. Zhao, and R. B. Noland, "A general map matching algorithm for  
 692 transport telematics applications," *Gps Solutions*, vol. 7, no. 3, pp. 157–167, 2003.
- 693 [14] C. E. White, D. Bernstein, and A. L. Kornhauser, "Some map matching algorithms for personal navigation  
 694 assistants," *Transportation research part c: emerging technologies*, vol. 8, no. 1-6, pp. 91–108, 2000.
- 695 [15] K. Zheng, Y. Zheng, X. Xie, and X. Zhou, "Reducing uncertainty of low-sampling-rate trajectories," In  
 696 *IEEE 28th International Conference on Data Engineering (ICDE 2012)*, Washington, DC, USA (Arlington,  
 697 Virginia), 1-5 April, 2012, pp. 1144–1155, 2012.
- 698 [16] F. Zhu, C. Luo, M. Yuan "City-scale localization with telco big data," In *Proceedings of the 25th ACM*  
 699 *International Conference on Information and Knowledge Management, CIKM 2016*, Indianapolis, IN,  
 700 USA, October 24-28, 2016, pp. 439–448, 2016.
- 701 [17] F. Zhu, M. Yuan, X. Xie, "A data-driven sequential localization framework for big telco data," *IEEE*  
 702 *Transactions on Knowledge and Data Engineering (Early Access)*, pp. 1-1, 2019.
- 703 [18] Z. Li, K. Liu, Y. Zhao, and Y. Ma, "Mapit: An enhanced pending interest table for NDN with mapping  
 704 bloom filter," *IEEE Communications Letters*, vol. 18, no. 11, pp. 1915–1918, 2014.
- 705 [19] Z. Li, L. Song, and H. Shi, "Approaching the capacity of k-user MIMO interference channel with  
 706 interference counteraction scheme," *Ad Hoc Networks*, vol. 58, pp. 286–291, 2017.
- 707 [20] Y. Liu and Z. Li, "A novel algorithm of low sampling rate GPS trajectories on map-matching," *EURASIP*  
 708 *Journal on Wireless Communications and Networking*, vol. 2017, no. 1, pp. 1-5, 2017.
- 709 [21] H. Wei, Y. Wang, G. Forman, Y. Zhu, and H. Guan, "Fast viterbi map matching with tunable weight  
 710 functions," In *SIGSPATIAL 2012 International Conference on Advances in Geographic Information*  
 711 *Systems (formerly known as GIS)*, SIGSPATIAL'12, Redondo Beach, CA, USA, November 7-9, 2012, pp.  
 712 613–616, 2012.
- 713 [22] T. Eiter, H. Mannila, *Computing discrete Fréchet distance*, Technical Report CD-TR 94/64, Christian  
 714 Doppler Laboratory for Expert Systems, TU Vienna, Austria, 1994.
- 715 [23] J. Han, H. Cheng, D. Xin, "Frequent pattern mining: current status and future directions," *Data mining and*  
 716 *knowledge discovery*, vol. 15, no. 1, pp. 55–86, 2007.
- 717 [24] J. Pei, J. Han, B. Mortazavi-Asl, "Prefifixspan: Mining sequential patterns by prefix-projected growth," In  
 718 *Proceedings of the 17th International Conference on Data Engineering*, April 2-6, 2001, Heidelberg,  
 719 Germany, pp. 215–224, 2001.
- 720 [25] R. Agrawal and R. Srikant, "Mining sequential patterns," In *Proceedings of the Eleventh International*  
 721 *Conference on Data Engineering*, March 6-10, 1995, Taipei, Taiwan, pp. 3–14, 1995.
- 722 [26] R. Srikant and R. Agrawal, "Mining sequential patterns: Generalizations and performance improvements,"  
 723 In *Advances in Database Technology - EDBT'96, 5th International Conference on Extending Database*  
 724 *Technology*, Avignon, France, March 25-29, 1996, Proceedings, pp. 3–17, 1996.
- 725 [27] M. J. Zaki, S. Parthasarathy, M. Ogihara, and W. Li, "New algorithms for fast discovery of association  
 726 rules," In *Proceedings of the Third International Conference on Knowledge Discovery and Data Mining*  
 727 *(KDD-97)*, Newport Beach, California, USA, August 14-17, 1997, pp. 283–286, 1997.
- 728 [28] J. Han, J. Pei, and Y. Yin, "Mining frequent pattern without candidate generation," In *Proceedings of the*  
 729 *2000 ACM SIGMOD International Conference on Management of Data*, May 16-18, 2000, Dallas, Texas,  
 730 USA., pp. 1–12, 2000.
- 731 [29] M. J. Zaki, "SPADE: An efficient algorithm for mining frequent sequences," *Machine Learning*, vol. 42,  
 732 no. 1/2, pp. 31–60, 2001.

- 733 [30] S. Lin, Y. Chen, D. Yang, and J. Wu, “Discovering long maximal frequent pattern,” In *Eighth*  
734 *International Conference on Advanced Computational Intelligence, ICACI 2016*, Chiang Mai, Thailand,  
735 February 14-16, 2016, pp. 136–142, 2016.
- 736 [31] R. Prabamanieswari, “A combined approach for mining fuzzy frequent itemset,” *International Journal of*  
737 *Computer Applications*, vol. 975, pp. 8887, 2013.
- 738 [32] C. M. Kuok, A. W. Fu, and M. H. Wong, “Mining fuzzy association rules in databases,” *SIGMOD Record*,  
739 vol. 27, no. 1, pp. 41–46, 1998.
- 740 [33] R. Agrawal and R. Srikant, “Fast algorithms for mining association rules in large databases,” In *VLDB’94,*  
741 *Proceedings of 20th International Conference on Very Large Data Bases*, September 12-15, 1994,  
742 Santiago de Chile, Chile, pp. 487–499, 1994.
- 743 [34] B. Han, L. Liu, and E. Omiecinski, “Road-network aware trajectory clustering: Integrating locality, flow,  
744 and density,” *IEEE Transactions on Mobile Computing*, vol. 14, no. 2, pp. 416–429, 2015.

745