

Efficient Clue-Based Route Search on Road Networks

Bolong Zheng, Han Su, Wen Hua, Kai Zheng, *Member, IEEE*,
Xiaofang Zhou, *Senior Member, IEEE*, and Guohui Li, *Member, IEEE*

Abstract—With the advances in geo-positioning technologies and location-based services, it is nowadays quite common for road networks to have textual contents on the vertices. Previous work on identifying an optimal route that covers a sequence of query keywords has been studied in recent years. However, in many practical scenarios, an optimal route might not always be desirable. For example, a personalized route query is issued by providing some clues that describe the spatial context between POIs along the route, where the result can be far from the optimal one. Therefore, in this paper, we investigate the problem of clue-based route search (CRS), which allows a user to provide clues on keywords and spatial relationships. First, we propose a greedy algorithm and a dynamic programming algorithm as baselines. To improve efficiency, we develop a branch-and-bound algorithm that prunes unnecessary vertices in query processing. In order to quickly locate candidate, we propose an AB-tree that stores both the distance and keyword information in tree structure. To further reduce the index size, we construct a PB-tree by utilizing the virtue of 2-hop label index to pinpoint the candidate. Extensive experiments are conducted and verify the superiority of our algorithms and index structures.

Index Terms—Spatial keyword queries, clue, point-of-interest, travel route search, query processing

1 INTRODUCTION

WITH the rapid development of location-based services and geo-positioning technologies, there is a clear trend that an increasing amount of geo-textual objects are available in many applications. For example, the location information as well as concise textual descriptions of some businesses (e.g., restaurants, hotels) can be easily found in online local search services (e.g., yellow pages). To provide better user experience, various keyword related spatial query models and techniques have emerged such that the geo-textual objects can be efficiently retrieved. It is common to search a Point-of-Interest (POI) by providing exact address or distinguishable keyword (i.e., only few POIs contain the keyword) in a region which can uniquely pinpoint the location. For example, we type the address “73 Mary St, Brisbane” or the name “Kadoya” on Google Maps to find a Japanese restaurant in the CBD area. Some existing work [1], [2], [3], [4], [5], [6] extends such query to more sophisticated settings, such as retrieving a group of geo-textual

objects (usually more than 2) or a trajectory covering multiple keywords. However, it is not uncommon that a user aims to find a POI with less distinguishable keyword such as “restaurant”, but she can only provide more or less spatio-textual context information around the POI. Liu et al. [7] formalize such context information as *clues* and use them to identify the most promising POIs. Different with their work, we aim to find a feasible route on road networks by using clues. Particularly, in this paper, we investigate a novel query type, namely *clue-based route search* (CRS), which allows a user to provide clues on textual and spatial context along the route such that a best matching route w.r.t. the clues is returned. More specifically, a CRS query is defined over a road network G , and the input of the query consists of a source vertex v_q and a sequence of clues, where each clue contains a query keyword and a user expected network distance. A vertex contains a clue keyword is considered as a match vertex. The query returns a path \mathcal{P} in G starting at v_q , such that (i.) \mathcal{P} passes through a sequence of match vertices (POIs) w.r.t. the clues and (ii.) the network distances between two contiguous matched vertices are close to the corresponding user specified distance such that the user’s search intention is satisfied.

1.1 Application Scenarios

The existing solutions (e.g., [8], [9], [10]) for trip planning or route search problem are dealing with the scenarios when a user wants to visit a sequence of POIs, each of which contains a user specified keyword. Different optimization constraints are proposed, and the goal is to find an optimal route with minimum cost. In general, the cost can be of various different types, such as travel distance, time or budget. However, to the best of our knowledge, none of the existing solutions (e.g., [8], [9], [10]) on trip planning or route search can be applicable for solving CRS queries since the optimization needs to be conducted based on the clues. As an

- B. Zheng and W. Hua are with the University of Queensland, St Lucia, QLD 4072, Australia. E-mail: {b.zheng, w.hua}@uq.edu.au.
- H. Su is with the Big Data Research Center, University of Electronic Science and Technology of China, Chengdu Shi 610051, China. E-mail: suan.sue@gmail.com.
- K. Zheng is with the School of Computer Science and Technology, Soochow University, Suzhou Shi 215000, China. E-mail: zhengkai@suda.edu.cn.
- X. Zhou is with the University of Queensland, St Lucia, QLD 4072, Australia, and the Macau University of Science and Technology, Taipa, Macau, China. E-mail: zxf@uq.edu.au.
- G. Li is with the Huazhong University of Science and Technology, Wuhan Shi 430074, China. E-mail: guohuili@hust.edu.cn.

Manuscript received 4 Aug. 2016; revised 8 May 2017; accepted 9 May 2017.
Date of publication 11 May 2017; date of current version 3 Aug. 2017.

(Corresponding author: Kai Zheng.)

Recommended for acceptance by K.S. Candan.

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

Digital Object Identifier no. 10.1109/TKDE.2017.2703848

extension of traditional route search queries, CRS query can also be useful in many real scenarios.

Modeling Imprecise User Intention. The clue is typically based on observations that 1) the keywords of PoIs in the clue may be interchangeable or inexact terms (e.g., a user may think of a PoI being *canteen* whereas it may be referred to as a *restaurant*); 2) the spatial relationships between PoIs are approximate, which is a natural phenomenon for human to estimate distance. For example, if the distance between two PoIs in the clue is about 100 meters, the actual distance may be noticeably greater than or less than 100 meters. Consider a scenario in our daily life: a user wants to find a restaurant in a city visited many years ago. She cannot remember the exact name and address but she still recalls that on the way driving to the restaurant from her home, she passed a cafe at about 1 km away, and drove about another 2 km to reach the restaurant. The information given above usually cannot precisely locate a PoI, but intuitively it provides clues to identify the most likely PoIs along the route.

Increased Flexibility in Trip Planning. As mentioned before, most existing work aims to find an optimal route with minimum travel distance. However, in many real scenarios, such an optimal route might not always be desirable. For instance, a user may have some personalized requirements on the distances between PoIs when planning a trip. Consider such a scenario, a user wants to find a buffet restaurant and a nearby cinema only in walking distance, say 3km, thus he can watch a movie after dinner. Therefore, after having delicious food, he can walk to the cinema in order to maintain a healthy lifestyle. These personalized requirements make the route search become distance-sensitive and more flexible such that the distance between PoIs along the route should be as close as possible to the user specified distance.

Clue-Based Route Navigation. Given a description including textual and distance information on an expected route, it is still not direct-viewing enough for users to obtain the exact route. This is usually the case when a user wants to know the way for a specific place and asks the others for help, she may still not be able to exactly figure out the route after obtaining the answers from them, where the answer usually comes in the form, for example, “go straight on the way for about 100 meters, you will see a cafe, and turn right, you will arrive the Japanese restaurant after about 150 meters walk”. Therefore, a novel type of route search which automatically interprets the clues contained in such answers becomes necessary. By augmenting it on current navigation services, a better user experience can be provided.

1.2 Challenge

In order to process the CRS query efficiently, we need to overcome several challenges. The first challenge is concerned with the large amount of possible routes for validation. Basically, the CRS requires candidate vertices that contain query keywords in the route to comply a specific order defined in query. As a feasible path is supposed to cover all the query keywords, the number of feasible paths increases exponentially with the amount of clues. Therefore, a greedy approach to solve our query is proposed, which continuously finds the next candidate vertex with minimum matching distance. Unfortunately, the optimal result can be substantially different from what the greedy algorithm suggests. Then, we propose a dynamic programming algorithm to answer CRS query exactly, but it requires quadratic time and is not scalable especially for more frequent keywords.

To avoid unnecessary route search, we develop a branch-and-bound algorithm which adopts filter-and-refine paradigm, thus much fewer feasible paths are considered.

The second challenge is how to quickly locate candidate vertices in road networks. Given a query vertex u , the matching distance between u and its next candidate v is supposed to be smaller or equal to a threshold. The network expansion approach can be applied here, but it is inefficient due to excessive network traversals. Therefore, we propose a novel index structure, called AB-tree, which stores both keyword and distance information in each node. On top of it, the candidate w.r.t. a query clue can be quickly retrieved. The third challenge is how to reduce the index construction time and space. As AB-tree involves an all-pair matrix computation and has a space cost of $O(|V|^2)$, we propose a PB-tree to further improve the performance. Inspired by the 2-hop label [11], [12], which answers distance queries with a small label index, we modify the structure of original label index to construct a binary tree on each pivot. In addition, we propose a semi-dynamic mechanism for PB-tree to support the index updating.

1.3 Contribution

The principal contributions of this paper can be summarized as follows.

- We propose a greedy clue search algorithm (GCS) to answer the CRS query approximately with no index involved. In GCS, we adopt the network expansion approach to greedily select the current best candidate at each step to construct feasible paths.
- We also develop a clue-based dynamic programming algorithm (CDP) that attempts to enumerate all feasible paths and finally returns the optimal result. In CDP, distance oracle is used to compute the network distance between candidates.
- We further propose a branch-and-bound algorithm (BAB) by applying filter-and-refine paradigm such that only a small portion of vertices are visited, hence improves the search efficiency. In order to quickly locate the candidate vertices, we develop AB-tree and PB-tree structures to speed up the tree traversal, as well as a semi-dynamic index updating mechanism to keep the index maintainable when growing bigger.
- Our experimental evaluation demonstrates the efficiency of our algorithms and index structures for processing the CRS queries on real-world datasets. We show the superiority of our algorithms in answering CRS when compared with the baseline algorithms.

The remainder of this paper is organized as follows. We first formulate the problem of clue-based route search CRS in Section 2. Then we propose a greedy algorithm GCS in Section 3 to answer CRS approximately. Section 4 presents a clue-based dynamic programming algorithm CDP to return exact answer to CRS query. Efficient branch-and-bound algorithm BAB is introduced in Section 5, as well as two index structures AB-tree and PB-tree. Section 6 presents a semi-dynamic mechanism for proposed index structure. Section 7 reports the experimental observations, and Section 8 reviews the related work. Finally, Section 9 concludes the paper.

2 PROBLEM STATEMENT

We model a road network as a weighted undirected graph $G = (V, E)$, where V is the set of vertices and E is the set of edges. Each vertex $v \in V$ contains a set of keywords, denoted as $\Phi(v)$. Each edge $(u, v) \in E$ has a positive weight,

TABLE 1
Summary of Notations

Notation	Definition
$G = (V, E)$	Road network with vertex V and edge E
$\mathcal{P}(u, \dots, v)$	A path from u to v
$\mathcal{FP}(u, \dots, v)$	A feasible path from u to v
$d_G(u, v)$	Network distance between u and v in G
$\mu(w, d, \epsilon)$	A clue with w , d and ϵ
$\sigma(u \rightarrow v)$	A match from u to v
$d_m(\mu, \sigma)$	Matching distance between μ and σ
$d_m(\mathcal{C}, \mathcal{FP})$	Matching distance between \mathcal{C} and \mathcal{FP}
$L(v)$	2-hop label of v
$BT(v)$	Binary tree of v with keyword and distance
$PR(o)$	Pivot-based reserve label of vertex o
$PB(o)$	Binary tree of pivot o

i.e., length or travelling time on the edge, denoted as $\epsilon(u, v)$. Given a path between vertices u and v , denoted as $\mathcal{P}(u, \dots, v)$, the length is the sum of weights of all edges along the path. For any two vertices u and v , the network distance between u and v on G , denoted as $d_G(u, v)$, is the length of the shortest path $SP(u, v)$ between u and v . The notations used in this work is summarized in Table 1.

2.1 Problem Definition

Definition 1 (Clue). A clue is defined as $\mu(w, d, \epsilon)$, where w is a query keyword, d is a user defined distance, and $\epsilon \in [0, 1]$ is a confidence factor.

Definition 2 (Match). Given a source vertex u and a clue $\mu(w, d, \epsilon)$, we say that the vertex pair $\sigma(u \rightarrow v)$ is a match w.r.t. clue μ , if the vertex v contains clue keyword w and the network distance between u and v is in $[d(1 - \epsilon), d(1 + \epsilon)]$, i.e., $w \in \Phi(v)$ and $d_G(u, v) \in [d(1 - \epsilon), d(1 + \epsilon)]$.

Definition 3 (Feasible Path). We define a clue-based route query $\mathcal{Q} = (v_q, \mathcal{C})$ where \mathcal{C} is a sequence of clues denoted as $\mathcal{C} = \{\mu_1(w_1, d_1, \epsilon_1), \dots, \mu_k(w_k, d_k, \epsilon_k)\}$. Given a query \mathcal{Q} , if we find

- (i) a sequence of vertices v_0, v_1, \dots, v_k where $v_0 = v_q$ and $\sigma(v_{i-1} \rightarrow v_i)$ is a match w.r.t. μ_i for $i \in [1, k]$;
- (ii) a path \mathcal{P} starts from v_0 and passes v_1, \dots, v_k one by one;
- (iii) each subpath $\mathcal{P}(v_{i-1}, \dots, v_i)$ of \mathcal{P} is the shortest path $SP(v_{i-1}, v_i)$ for $i \in [1, k]$.

We call such \mathcal{P} a feasible path, denoted as $\mathcal{FP}(v_q, v_1, \dots, v_k)$.

Definition 4 (Matching Distance). The matching distance between a clue $\mu(w, d, \epsilon)$ and its match $\sigma(u \rightarrow v)$ in G , denoted as $d_m(\mu, \sigma)$, is computed by d and the network distance $d_G(u, v)$, such that

$$d_m(\mu, \sigma) = \frac{|d_G(u, v) - d|}{\epsilon \cdot d}. \quad (1)$$

It is worth noting that any monotonic increasing function that normalizes the matching distance into $[0, 1]$ can be applied here.

The matching distance between \mathcal{C} and its feasible path \mathcal{FP} is defined as the maximum matching distance between all clues $\mu \in \mathcal{C}$ and their corresponding matches $\sigma \in \mathcal{FP}$, that is

$$d_m(\mathcal{C}, \mathcal{FP}) = \max_{\mu_i \in \mathcal{C}, \sigma_i \in \mathcal{FP}} d_m(\mu_i, \sigma_i). \quad (2)$$

The motivation of using Equation (2) is that the maximum matching distance of all the clues naturally controls

the overall matching quality of the feasible path, which is a widely adopted method for the problem optimizing an objective score contributed by several components [5], [6].

Clue-Based Route Search. Given a clue-based route search $\mathcal{Q} = (v_q, \mathcal{C})$, it aims to find a feasible path $\mathcal{FP}(v_q, v_1, \dots, v_k)$, such that $d_m(\mathcal{C}, \mathcal{FP})$ is minimized.

The clues in CRS implies that we are supposed to find a feasible path whose distances between two match vertices are as close as possible to user specified distance such that the user's search intention is satisfied. It is worth noting that the CRS query can be easily extended to have a destination by assuming that the query keyword contained in destination is unique within G . In addition, for simplicity, we only discuss the optimal feasible path in this paper, but the algorithms introduced can be easily extended to find top- k feasible paths.

Example 1. Given $\mathcal{Q} = (v_1, \{(w_2, 5, 0.5), (w_1, 4, 0.5), (w_3, 6, 0.5)\})$, thus both $\mathcal{FP}_1(v_1, v_3, v_6, v_7)$ and $\mathcal{FP}_2(v_1, v_3, v_4, v_7)$ are feasible paths with $d_m(\mathcal{C}, \mathcal{FP}_1) = 0.4$ and $d_m(\mathcal{C}, \mathcal{FP}_2) = 0.5$ respectively. Therefore, \mathcal{FP}_1 is reported as the result of CRS.

2.2 Preliminary: Distance Oracle

We adopt the idea of distance oracle \mathcal{DO} to calculate the network distance between two input vertices. Given a source-target pair of vertices, \mathcal{DO} returns the shortest network distance between them. As we know, the algorithms and data structures on \mathcal{DO} have been extensively studied by existing works, which can be roughly summarized into two categories, expansion-based methods and lookup-based methods. The most famous expansion-based method for \mathcal{DO} is Dijkstra's algorithm [13], which, given a s - t pair in road network G , traverses the vertices in G from s to t . However, the problem of using Dijkstra's algorithm is that it must visit every vertex that is closer to s , and the number of such unneeded vertices can be enormous when s and t are far apart, which incurs redundant network traversal. Besides, the lookup-based methods usually have to store some precomputed results. For example, all-pair method is space inefficient that we have to precompute and store a distance matrix, which requires $O(n^2)$ space for a road network G with n vertices. To the best of our knowledge, one of the most notable recent developments is the emergence of practical 2-hop labeling methods [11], [12], [14], [15] for \mathcal{DO} on large networks. It constructs labels for vertices such that a distance query for any vertex pair u and v can be answered by only looking up the common labels of u and v . For each vertex v , we precompute a label, denoted as $L(v)$, which is a set of label entries and each label entry is a pair $(o, \eta_{v,o})$, where $o \in V$ and $\eta_{v,o} = d_G(v, o)$ is the distance between v and o . We say that o is a pivot in label entry if $(o, \eta_{v,o}) \in L(v)$. Given two vertices u and v , we can find a common pivot o that $(o, \eta_{u,o}) \in L(u)$ and $(o, \eta_{v,o}) \in L(v)$

$$d_G(u, v) = \min\{\eta_{u,o} + \eta_{v,o}\}. \quad (3)$$

We say that the pair (u, v) is covered by o and the distance query $d_G(u, v)$ is answered by o with smallest $\eta_{u,o} + \eta_{v,o}$. Therefore, we can compute $d_G(u, v)$ in $O(|L(u)| + |L(v)|)$ time by using a merge-join like algorithm.

3 GREEDY CLUE SEARCH ALGORITHM

We develop a greedy algorithm as a baseline for answering the CRS query, which is called Greedy Clue Search algorithm. Given a query $\mathcal{Q} = (v_q, \mathcal{C})$, we first add v_q into a candidate path. Then we use the Procedure findNextMin() to determine

the next match vertex v_1 that the matching distance between μ_1 and $\sigma_1(v_q \rightarrow v_1)$, i.e., $d_m(\mu_1, \sigma_1)$, is minimized. Afterwards, we insert v_1 into the candidate path, and continue to find its contiguous candidate by `findNextMin()`. This process is repeated until all the match vertices are determined, thus the candidate path forms a feasible path, denoted as \mathcal{FP}_{v_q} . If we assume Procedure `findNextMin()` costs time f , then the time complexity of GCS is $O(k \cdot f)$.

In Procedure `findNextMin()`, we utilize the network expansion algorithm [16] to find the nearby vertices that contain the query keywords and the network distances are in the confidence intervals. The algorithm details are shown in Algorithm 1. Given the source u , and the clue $\mu(w, d, \epsilon)$, we aim to find a match vertex v such that the difference between $d_G(u, v)$ and d is minimized. In the network traversal starting from u , we check every visited vertex to see if it is a match vertex that contains w and locates in the interval $[d(1 - \epsilon), d(1 + \epsilon)]$. If v is the first visited match vertex and $d_G(u, v) > d$, then we stop and return v since the difference incurred by the remaining unvisited vertices cannot be less than $d_G(u, v) - d$. Otherwise, we continue to find the next match vertex v' . If v' is found, then

- (i) If $d_G(u, v') \leq d$, we update v by v' since v' renders a smaller difference than v ;
- (ii) Otherwise, we compare $d - d_G(u, v)$ with $d_G(u, v') - d$ and return the smaller one as the result.

Example 2. In running example Fig. 1, we are given CRS query $\mathcal{Q} = (v_1, \{(w_2, 5, 0.5), (w_1, 4, 0.5), (w_3, 6, 0.5)\})$. First, we fetch v_1 into the candidate path, and call `findNextMin($v_1, w_2, 5, 0.5$)` and return v_3 with $d_m = 0.4$. Therefore, we repeat the process and finally obtain $\mathcal{FP}_{v_1} = (v_1, v_3, v_4, v_7)$ with $d_m(\mathcal{C}, \mathcal{FP}_{v_1}) = 0.4$. Therefore, we have $\mathcal{FP}_{gcs} = \mathcal{FP}_{v_1}$ and $d_m(\mathcal{C}, \mathcal{FP}_{gcs}) = 0.4$.

Algorithm 1. Procedure `findNextMin()`

Input: Source vertex u and clue $\mu(w, d, \epsilon)$
Output: $\min\{d_m(\mu, \sigma)\}$ and match vertex v

- 1 From u , do network traversal;
- 2 **if** a match vertex v is found **then**
- 3 $d_G \leftarrow$ the network distance between u and v ;
- 4 **while** true **do**
- 5 Find next v' contains w , thus obtain d'_G ;
- 6 **if** $d_G < d$ and $d'_G > d$ **then**
- 7 **break**;
- 8 **else**
- 9 $v \leftarrow v'$ and $d_G \leftarrow d'_G$;
- 10 **return** $\min\{d_m(\mu, \sigma)\}$ and v ;

4 CLUE-BASED DYNAMIC PROGRAMMING ALGORITHM

As we know, even though GCS has a short response time, the accuracy of the answer cannot be guaranteed. To achieve better accuracy, we propose an exact algorithm, called Clue-based Dynamic Programming, to answer the CRS query. Generally, it is challenging to develop an efficient exact algorithm for CRS queries, since we cannot avoid exhaustive search for PoIs in road networks. For instance, the number of vertices that contain keyword $w_i \in \mathcal{C}$ is denoted as $|V_{w_i}|$, thus the time complexity of the brute-force approach, which attempts all possible combinations, is $O(\prod_{w_i \in \mathcal{C}} |V_{w_i}|)$.

In CDP, we construct a keyword posting list for each keyword w , which is a list of vertices that contain w . When a

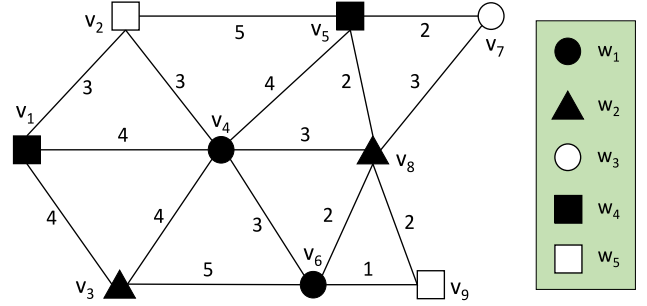


Fig. 1. Running example of G .

CRS query is issued, we sort the posting lists according to the keyword order of $w_i \in \mathcal{C}$. Note that the order of the vertices within each posting list does not matter and can be arbitrary, hence are sorted by vertex id for simplicity. It is easy to see that these posting lists actually construct a k -bipartite graph G' , which in fact shows all feasible paths for a given \mathcal{C} . The weight of each edge in G' is computed as the matching distance. Specifically, for each $u \in V_{w_i}$, we define $D(w_i, u)$ to denote the minimum matching distance one can achieve with a walk that passes the keywords from w_1 to w_i consistent with the order in \mathcal{C} and stops at u . In other words, the weight of vertex $u \in G'$ is computed by $D(w_i, u)$, which is the minimum matching distance of all partial feasible paths end at u . Then we compute $D(w_i, u)$ by the following recursive formula:

- (i) $i = 1$: for match vertices $u \in V_{w_1}$, we have

$$D(w_1, u) = d_m(\mu_1(w_1, d_1), \sigma(v_q \rightarrow u)).$$

- (ii) $i > 1$: for match vertices $v \in V_{w_{i-1}}$ and $u \in V_{w_i}$, we have

$$D(w_i, u) = \min_{v \in V_{w_{i-1}}} \{\max\{D(w_{i-1}, v), d_m(\mu_i, \sigma(v \rightarrow u))\}\}. \quad (4)$$

Algorithm 2. Clue-Based Dynamic Programming CDP

Input: $\mathcal{Q} = (v_q, \mathcal{C} = \{(w_1, d_1), \dots, (w_k, d_k)\})$
Output: \mathcal{FP}_{cdp} with $d_m(\mathcal{C}, \mathcal{FP}_{cdp})$

- 1 **for** each $u \in V_{w_1}$ **do**
- 2 Initial $D(w_1, u)$;
- 3 **for** $1 < i \leq k$ **do**
- 4 **for** each $u \in V_{w_i}$ **do**
- 5 Initial intermediate vector $iv(u)$;
- 6 **for** each $v \in V_{w_{i-1}}$ **do**
- 7 **if** $d_m(\mu_i, \sigma(v \rightarrow u)) < D(w_{i-1}, v)$ **then**
- 8 $iv(u)$ insert $D(w_{i-1}, v)$;
- 9 **else**
- 10 $iv(u)$ insert $d_m(\mu_i, \sigma(v \rightarrow u))$;
- 11 $D(w_i, u) \leftarrow \min\{iv(u)\}$;
- 12 Find $\min\{D(w_k, u)\}$;
- 13 **return** \mathcal{FP}_{cdp} and $d_m(\mathcal{C}, \mathcal{FP}_{cdp}) \leftarrow \min\{D(w_k, u)\}$;

For each iteration, we have at most $|V_{w_{i-1}}| \cdot |V_{w_i}|$ combinations, thus the time required is $O(\sum_{i=2}^k |V_{w_{i-1}}| \cdot |V_{w_i}|)$. The details of CDP are shown in Algorithm 2. In order to compute $D(w_i, u)$, we have to access the posting list of w_{i-1} . For each vertex v in this list, we first check if $\sigma(v \rightarrow u)$ is a match w.r.t. μ_i and then compute $d_m(\mu_i, \sigma(v \rightarrow u))$. Then we compare it with $D(w_{i-1}, v)$, and keep the greater one as intermediate value. Finally, we find the minimum one as $D(w_i, u)$ from these $|V_{w_{i-1}}|$ intermediate values. After we recursively process

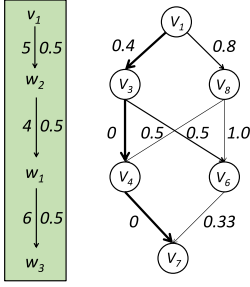


Fig. 2. Matching distances of CDP.

all the keywords, we finally find the minimum $D(w_k, u)$ and backtrace the corresponding vertices that construct \mathcal{FP}_{cdp} .

In each iteration, we have a clue $\mu_i(w_i, d_i)$, therefore we have to compute $d_G(u, v)$ between each $u \in V_{w_i}$ and its precedents $v \in V_{w_{i-1}}$ as prerequisites for determining $d_m(\mu_i, \sigma(v \rightarrow u))$. Here we adopt the distance oracle introduced in Section 2.2 to compute $d_G(u, v)$.

Example 3. As shown in Fig. 2, given CRS query $\mathcal{Q} = (v_1, \{(w_2, 5, 0.5), (w_1, 4, 0.5), (w_3, 6, 0.5)\})$. To compute $D(w_3, v_7)$, we first compare $D(w_1, v_4) = 0.4$ with $d_m(\mu_3, \sigma(v_4 \rightarrow v_7)) = 0$, and obtain intermediate value 0.4. Likewise, we have $D(w_1, v_6) = 0.5$ and $d_m(\mu_3, \sigma(v_6 \rightarrow v_7)) = 0.33$, thus the intermediate value is 0.5. Therefore, CDP returns $\mathcal{FP}_{cdp} = (v_1, v_3, v_4, v_7)$ with $d_m(\mathcal{C}, \mathcal{FP}_{cdp}) = 0.4$.

5 BRANCH AND BOUND ALGORITHM

Although CDP provides an exact solution, the search efficiency cannot be maintained. For instance, consider the worst case, we assume that all vertices contain query keywords, then the time is $O(k \cdot |V|^2)$. To propose a more efficient algorithm, we assume there is an artificial directed graph G' , which is similar to the k -bipartite graph in CDP that formed by all candidate vertices containing keywords in \mathcal{C} , where the edge of G' is a match of one clue and in the meantime its direction complies the keyword order of the clue. Note that, G' is organised into k levels, and each level i corresponds to each keyword w_i . Based on G' , we develop a Branch-and-Bound algorithm to search G' in a depth-first manner by applying the filter-and-refine paradigm, which only visits a small portion of vertices in G' . Fortunately, we can use the result of GCS to speed up the search process since it can serve as an initial upper bound.

5.1 Algorithm Outline

We start the searching from level 1 to k to obtain a feasible path \mathcal{FP} , if the matching distance $d_m(\mathcal{C}, \mathcal{FP})$ is greater than the current upper bound, we continue to search for the next candidate feasible path, otherwise we update the upper bound. It is worth noting that it is not necessary to go through every candidate feasible path. If the matching distance at intermediate level already exceeds the upper bound, it can be removed. This process terminates when the matching distance next to be processed at level 1 can be filtered, since it is impossible to find a feasible path with smaller match distance.

Candidate Feasible Path Updating. Initially, we keep a stack to store the partial candidate path, which contains a sequence of vertices and corresponding matching distances. First, we fetch v_q into the stack, then we continue to find next candidate at level 1. Basically, the key component of this algorithm is to quickly locate the next best match vertex, and

the details of Procedure `findNext()` will be introduced later. Given a partial candidate path $\mathcal{FP}(v_q, v_1, \dots, v_i)$ obtained at level i , we apply `findNext()` to find the next candidate v_{i+1} at level $i+1$. Once v_{i+1} is found, we compute $d_m^{i+1}(v_{i+1})$ which denotes the matching distance at level $i+1$ resulted by v_{i+1} , and compare it with current UB . Note that, v_{i+1} is accepted as a candidate and inserted into the stack if and only if its matching distance $d_m^{i+1}(v_{i+1})$ is smaller than UB . Otherwise, v_i is removed from the stack as well as $d_m^i(v_i)$. In other words, v_i is not valid that the path $\mathcal{FP}(v_q, v_1, \dots, v_{i-1})$ cannot survive by passing v_i , then we have to find an alternative v'_i . As we know v_i is the current best candidate at level i , therefore we have to relax the matching distance by finding v'_i where $d_m^i(v_i) \leq d_m^i(v'_i)$ and $d_m^i(v'_i)$ is minimum among all the rest vertices untouched at level i . Afterwards, if v'_i is valid, we continue to apply `findNext()` on it.

Upper Bound Updating. Specifically, after we obtain a feasible path $\mathcal{FP}(v_q, v_1, \dots, v_{k-1})$ at level $k-1$, if v_k is returned by `findNext()`, then we check if $d_m^k(v_k)$ exceeds UB . If v_k is not valid, we prune v_k and simply repeat the above process. Otherwise, we insert v_k into the stack, and a complete feasible path is determined. Hence, $\mathcal{FP}(v_q, v_1, \dots, v_k)$ is regarded as a temporary result, and UB is updated by the minimum matching distance among all $d_m^i(v_i)$ s. It is easy to see that, we cannot find a better feasible path by alternating v_k with v'_k at level k , since no further level is available to make up the relaxation caused by v'_k . Therefore, in addition to remove v_k , we continue to remove v_{k-1} from the stack and repeat the above process.

In general, the pruning happens from the lower levels to the higher levels, i.e., from level k to level 1. In the end, at level 1, if the matching distance induced by the next candidate vertex is greater than UB , it is impossible to find another feasible path, thus the stack becomes empty after the last vertex v_q is removed, and this process terminates.

Example 4. In the running example, given query $\mathcal{Q} = (v_7, \{(w_1, 6, 0.5), (w_2, 4, 0.5), (w_4, 5, 0.5)\})$. First we fetch v_7 into the stack, and `findNext()` returns v_4 with $d_m^1(v_4) = 0$. Then we insert v_4 into stack and continue to find next candidate vertex, and v_3 is obtained with $d_m^2(v_3) = 0$. The process continues and then we have v_1 with $d_m^3(v_1) = 0.4$. As the size of stack is same as the number of query keywords, a feasible path $\mathcal{FP} = (v_7, v_4, v_3, v_1)$ with $d_m(\mathcal{C}, \mathcal{FP}) = 0.4$ is obtained, and UB is updated by 0.4. Next, we remove v_1 and v_3 from the stack, and continue to find next candidate of v_4 . As $d_m^2(v_3) = 0$, we relax the matching distance and call `findNext()` which returns v_8 with $d_m^2(v_8) = 0.5$. Then we have to remove v_4 from the stack since $d_m^2(v_8)$ already exceeds current upper bound UB . Now we move on to apply `findNext()` on v_7 and returns v_6 with $d_m^1(v_6) = 0.33$. However, the next candidate v_3 has $d_m^2(v_3) = 0.5$ greater than UB , thus we remove v_6 and v_7 from stack. Therefore, the algorithm terminates since no other feasible path exists. We have $\mathcal{FP}_{bab} = (v_7, v_4, v_3, v_1)$ with $d_m(\mathcal{C}, \mathcal{FP}_{bab}) = 0.4$.

5.2 All-Pair Distance Approach

In BAB, the Procedure `findNext()` is applied on v_{i-1} to find the next candidate vertex v_i . We can simply use Procedure `findNextMin()` in GCS to locate the next candidate, but it is inefficient due to redundant network traversal especially when $d_i \in \mu_i$ is large. Moreover, when we prune v_i and attempt to find alternative v'_i , it is easy to see `findNextMin()` cannot be directly applied. Therefore, we propose an All-pair Binary tree (AB-tree) index to improve the search efficiency.

Algorithm 3. Branch and Bound BAB

Input: $\mathcal{Q} = (v_q, \mathcal{C})$
Output: \mathcal{FP}_{bab} with $d_m(\mathcal{C}, \mathcal{FP}_{bab})$

- 1 Initialize stackV, stackD, and search threshold θ ;
- 2 Push v_q into stackV;
- 3 **while** stackV is not empty **do**
- 4 $i \leftarrow \text{stackV.size}()$;
- 5 **if** $\text{findNext}(v_{i-1}, d_i, w_i, \theta) = \text{true}$ **then**
- 6 Obtain v_i and $d_m^i(v_i)$;
- 7 $\theta \leftarrow 0.0$;
- 8 Push v_i into stackV and $d_m^i(v_i)$ into stackD;
- 9 **if** i equals to k **then**
- 10 **if** $\max\{\text{stackD}\} \leq UB$ **then**
- 11 Update UB by $\max\{\text{stackD}\}$;
- 12 Update \mathcal{FP}_{bab} by stackV;
- 13 Update θ by top of stackD;
- 14 Update stackV and StackD;
- 15 **else**
- 16 Update θ by top of stackD;
- 17 Update stackV and StackD;
- 18 **return** \mathcal{FP}_{bab} and $d_m(\mathcal{C}, \mathcal{FP}_{bab}) \leftarrow UB$;

5.2.1 All-Pair Binary Tree

Given a vertex u , we aim to find a vertex v containing keyword w such that the matching distance between $\sigma(u \rightarrow v)$ and query clue μ is slightly greater than and closest to a threshold θ among all vertices containing w in G . Note that, the threshold θ is settled by previous filtered candidate at the same level with v , and it is 0 at initial stage. In other words, we are supposed to find the vertex v that the difference between $d_G(u, v)$ and $d \in \mu$ is close to $\theta \cdot d$. To this end, we construct AB-tree as follows.

For each $v \in V$, we construct a binary tree $BT(v)$ that contains the information of network distances and keywords. After the all-pair distance matrix is obtained, for each v , we have a list of vertices sorted in ascending order of network distance to v . By utilizing the tree structure, the vertices in the list are divided into fragments that the network distances w.r.t. v of the vertices in the same fragment are close to each other, which speeds up the looking up for vertices by network distance. In addition, the keyword information within each fragment is also stored in $BT(v)$ such that the vertices containing query keyword in a fragment can be efficiently retrieved.

We utilize a hash function H that maps keywords and vertices to a binary code with h bits. For each keyword w , one of its h bits in $H(w)$ is set to 1. Hence, the binary code of a vertex v is the superimposition of $H(w)$ for all w it contains, i.e., $H(v) = \bigvee_{w \in \Phi(v)} H(w)$. Likewise, for a set of vertices S , $H(S)$ is the superimposition of $H(v)$ that $v \in S$. It is worth to note that a non-zero value of $H(w) \wedge H(S)$ indicates that there may exist a vertex $v \in S$ containing w , and $H(w) \wedge H(S) = 0$ means w is definitely not contained by any $v \in S$. $BT(v)$ is actually a B^+ -tree with fanout $f = 2$. Each leaf node contains the information of a vertex u with both the network distance $d_G(u, v)$ and binary code $H(u)$ stored. For a non-leaf node, it also keeps a routing element, which equals the maximum network distance of its left subtree. Therefore, $BT(v)$ is constructed recursively in bottom-up manner as shown in Fig. 3a.

Storing $BT(v)$ in an Array. As we know, storing the tree structure as an array enables a better performance than storing pointers. Therefore, we propose a scheme to sequentially store all nodes of $BT(v)$ in an array from nodes on height 0 to

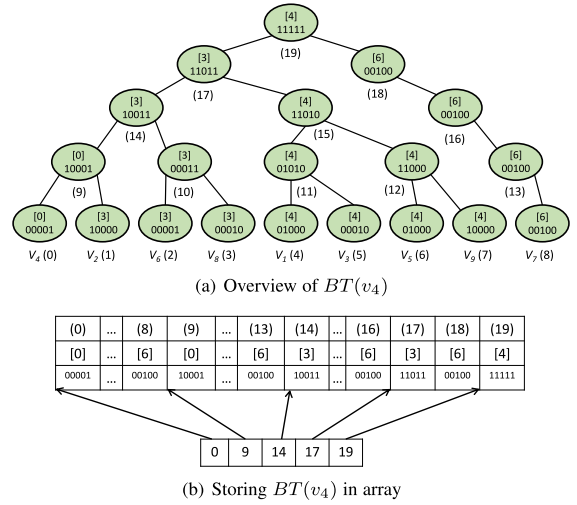


Fig. 3. Overview of all-pair binary tree, with illustration of the two components: (a) Example of $BT(v_4)$. (b) Storing in array.

the root, as shown in Fig. 3b. In addition to this array, we also keep an auxiliary array that indicates the number of nodes in each level of $BT(v)$, by which we can quickly determine the indices of the subnodes of a non-leaf node, or the index of its parent node, in the array. For example, if we want to find the left and right subnodes of node 16 in $BT(v_4)$, we know there are two nodes on its left side by $16 - 14 = 2$ where 14 is the start index of nodes at height 2, so the index of its left subnode is $9 + 2 * 2 = 13$ and the right is $9 + 2 * 2 + 1 = 14$. However, we notice 14 is actually at height 2, then we figure out node 16 does not have a right subnode.

Lemma 1. Given $G = (V, E)$, the space cost of AB-tree is $O(|V|^2 \cdot h)$.

Proof. For each $v \in V$, we have $|V|$ elements in distance matrix, thus each $BT(v)$ has an index size $O(|V| \cdot h)$. It is easy to see the size of AB-tree is $O(|V|^2 \cdot h)$. \square

5.2.2 Predecessor and Successor Queries on AB-Tree

After the construction of AB-tree, we discuss how to use it so that the next vertex in candidate path can be quickly located. Initially, if there is no previous vertices accessed at the next level of v_{i-1} , the network distance $d_G(v_{i-1}, v_i)$ between v_{i-1} and next candidate v_i is supposed to be smaller or equal to $lD = d_i$, or greater or equal to $rD = d_i$, where $d_i \in \mu_i$.

Additionally, consider the aforementioned scenario, we have $\mathcal{FP}(v_q, v_1, \dots, v_i)$, but v_{i+1} returned at level $i+1$ exceeds UB . Then we have to remove v_i from the stack and turn to find v'_i as alternative, where $d_m^i(v_i) \leq d_m^i(v'_i)$. It is easy to see the difference between $d_G(v_{i-1}, v'_i)$ and d_i must be greater or equal to $d_m^i(v_i) \cdot d_i$. In other words, the network distance $d_G(v_{i-1}, v'_i)$ is smaller or equal to lD or greater or equal to rD , where

$$\begin{aligned}
 lD &= d_i - d_m^i(v_i) \cdot d_i, \\
 rD &= d_i + d_m^i(v_i) \cdot d_i,
 \end{aligned} \tag{5}$$

we have $d_G(v_{i-1}, v'_i) \leq lD$ or $d_G(v_{i-1}, v'_i) \geq rD$.

Therefore, the predecessor and successor queries can be issued on $BT(v_{i-1})$ to retrieve next candidate with two boundary network distances lD and rD , respectively.

Predecessor Query. Given $BT(u)$, a query keyword w and network distance lD , we aim to find vertex v that contains w

and $d_G(u, v)$ is smaller or equal to and closest to lD . First, we compute binary code $H(w)$ for query keyword w . Then we start the process of searching $BT(u)$ recursively from top to bottom. For a non-leaf node o , if $H(w) \wedge H(o)$ is non-zero, we continue to search its subtrees. If lD is smaller than the routing element of o , only its left subtree needs to be considered. Otherwise, we first check if we could find v in its right subtree (if exists), if not, we turn to search its left subtree. For a leaf node v , we directly check if v contains w and $d_G(u, v)$ is smaller or equal to lD , therefore, false positives can be avoided. Finally, v is obtained. For example, a predecessor query on $BT(v_4)$ with keyword w_2 and $lD = 4$. First, we have $H(w_2) = 00010$. The search starts from root, and as lD equals to the routing element 4, thus we first search its right subtree. After checking $H(w_2)$ with binary code of Node 18, we find it does not contain w_2 and we turn to search Node 17. As the routing element of Node 17 is smaller than lD , we move to search Node 15. Then we check $H(w_2)$ with the binary code in Node 12, and find Node 12 does not contain w_2 . After checking with Node 11, we have v_3 as result of the predecessor query.

Successor Query. Likewise, we have $BT(u)$, a query keyword w and network distance rD , the goal is to find vertex v that contains w and $d_G(u, v)$ is greater or equal to and closest to rD . For a non-leaf node o , if $H(w) \wedge H(o)$ is non-zero, the subtrees of o need to be considered. If rD is smaller or equal to the routing element of o , we search the left subtree to see if v could be found, if not, we turn to search the right subtree (if exists). Otherwise, we simply search the right subtree (if exists) to locate v . For a leaf node v , if v contains w and $d_G(u, v)$ is greater or equal to rD , v is reported as result. For example, a successor query on $BT(v_4)$ with keyword w_2 and $rD = 4$. We first check the root with $H(w_2) = 00010$, and rD equals to the routing element, which means we first search Node 17 to see if it contains w_2 , then search Node 18. As the routing element of Node 17 is smaller than rD , we only need to check Node 15. Then, since the routing element of Node 15 is same as rD , we turn to search Node 11. Finally, we obtain v_3 as result of the successor query.

As mentioned before, we process a predecessor and a successor queries on $BT(v_{i-1})$ with lD and rD respectively to locate candidate at level i . If both predecessor and successor queries find candidate vertices, we compare their matching distance and report the smaller one as result. If only one of them finds candidate vertex, we directly report it. Otherwise, no candidate is found. Note that, in the process to replace v_i with v'_i , we must skip v_i in the tree traversal to avoid infinite loop caused by the special case $d_m^i(v_i) = d_m^i(v'_i)$.

Lemma 2. *The expected number of nodes visited in a predecessor or successor query on AB-tree is $O(\frac{\log|V| \cdot |W| \cdot |\Phi(V)|}{|V| \cdot h \cdot |V_{w_k}|})$.*

Proof. Assume W is the keyword set of G and the keywords are evenly distributed. The hash function maps each keyword $w \in W$ to a binary code with h bits based on its modulus, which might lead to false positives. We denote the average number of conflicting keywords as $\frac{|W|}{h}$ and the average keyword frequency as $f_{avg} = \frac{|\Phi(V)|}{|V|}$. Therefore, the probability of false positive can be computed as $p = 1 - \frac{|V| \cdot h \cdot |V_{w_k}|}{|W| \cdot |\Phi(V)|}$. When encountering a false-positive at the leaf node, the predecessor query traces back from the right subtree and continues to search the left subtree to find another leaf node. In this case, $O(\log|V|)$ tree nodes

will be visited for each trace back at the worst case. The algorithm terminates when it reaches a true positive. Hence, the expected number of trace backs is equivalent to the expected number of false positives before a true positive, which is $\frac{1}{1-p} = \frac{|W| \cdot |\Phi(V)|}{|V| \cdot h \cdot |V_{w_k}|}$. In total, the expected number of tree nodes visited in a predecessor query is $O(\frac{\log|V| \cdot |W| \cdot |\Phi(V)|}{|V| \cdot h \cdot |V_{w_k}|})$. Assume the time cost for \wedge operation on two h length hash codes is $O(h)$, thus the time of searching AB-tree is $O(\frac{\log|V| \cdot |W| \cdot |\Phi(V)|}{|V| \cdot h \cdot |V_{w_k}|}) \cdot O(h)$. \square

Example 5. In the running example, given query $\mathcal{Q} = (v_7, \{(w_1, 6, 0.5), (w_2, 4, 0.5), (w_4, 5, 0.5)\})$, assume we already have stack (v_7, v_4, v_3) . At level 2, we intend to remove v_3 and find an alternative. Given $d_m^2(v_3) = 0$, we apply a predecessor and successor queries on $BT(v_4)$. For the predecessor query, we take $(w_2, 4, 0.5)$ and 0.0 as input. As v_3 is previous result, we skip it and return v_8 . For the successor query, no vertex is found. Therefore, we report v_8 as our next candidate with $d_m^2(v_8) = 0.5$.

Algorithm 4. Procedure findNext() with AB-Tree

Input: Query vertex v_{i-1} , clue w_i and d_i , threshold θ

Output: Next candidate v_i with $d_m^i(v_i)$

```

1 Obtain  $BT(v_{i-1})$ ;
2  $lD \leftarrow d_i - d_i \cdot \theta$ ;  $rD \leftarrow d_i + d_i \cdot \theta$ ;
3  $v_p$  and  $d_p \leftarrow BT(v_{i-1}).predecessor(lD, w_i)$ ;
4  $v_s$  and  $d_s \leftarrow BT(v_{i-1}).successor(rD, w_i)$ ;
5 if  $d_i - d_p \leq d_s - d_i$  then
6   return  $v_p$  with  $d_m(v_p)$ ;
7 else
8   return  $v_s$  with  $d_m(v_s)$ ;
Procedure Predecessor( $lD, w, Node$ )
1 if  $Node$  is a leaf node then
2   Obtain  $v_p$  and  $d_p$  of current node;
3   if  $v_p$  contains  $w$  and  $d_p \leq lD$  then
4     return  $v_p$  and  $d_p$ ;
5   else
6     return false;
7 else
8   Generate  $H(w)$ ;
9   if  $H(w) \wedge H(Node) = 0$  then
10    return false;
11 if  $lD < Node.routing$  then
12    $lNode \leftarrow$  index of its left subnode;
13   return  $Predecessor(lD, w, lNode)$ ;
14 else
15    $rNode \leftarrow$  index of its right subnode;
16    $lNode \leftarrow$  index of its left subnode;
17 if  $rNode$  exists then
18   if  $Predecessor(lD, w, rNode)$ ;
19   then
20     return  $v_p$  and  $d_p$ ;
21   else
22     return  $Predecessor(lD, w, lNode)$ ;
23 return  $Predecessor(lD, w, lNode)$ ;

```

5.3 Keyword-Based Label Approach

Even though AB-tree is able to answer findNext() query fast, the index space cost is still high and could only be stored in disk, which results in undesired I/O consumption. In this

section, we introduce a main memory based index structure, namely Pivot reverse Binary tree (PB-tree), to deal with `findNext()` query.

5.3.1 Pivot Reverse Binary Tree

As introduced in Section 2.2, we know 2-hop label possesses the nature to process distance queries between any two vertices in network with fast response time, whilst keeping the size of the generated label index as small as possible. The problem of reducing label size is orthogonal to our work, thus we fully utilize the state-of-the-art results to build a small index in this work. As we know, in 2-hop label, the distance between any vertex pair (u, v) can be computed correctly through a common pivot o , in other words, each vertex u can reach any other vertex v in network through a pivot o . Therefore, based on this intuition, we modify the structure of original 2-hop label to construct a *pivot reverse index*, i.e., *PR* index [17] which stores all label entries $(o, \eta_{u,o}) \in \bigcup_{v \in V} L(v)$ regarding vertex o as pivot into the *PR* label of vertex o , i.e., $(v, \eta_{v,o}) \in PR(o)$. In $PR(o)$, we assume that all the label entries $(v, \eta_{v,o})$ are sorted in ascending order of distance. For example, we have $(v_3, 0) \in L(v_3)$ and $(v_3, 4) \in L(v_1)$. Through the transformation, we have $PR(v_3) = \{(v_3, 0), (v_1, 4)\}$.

In order to find vertex by keyword and distance information, each $PR(o)$ is organized as same as the binary tree mentioned before, thus forms $PB(o)$. The structure is shown in Fig. 4, it is worth to note that any network distance $d_G(u, v)$ is divided into two parts, the first part $d_G(u, o)$ between u and its pivot o can be found in $L(u)$, and the other part $d_G(o, v)$ between pivot o and target v can be found in $PB(o)$. Therefore, combined with original label index whose label entries are also sorted in ascending order by network distance, PB-tree could be used to answer predecessor and successor queries more efficiently than AB-tree with a much smaller size.

Lemma 3. Given $G = (V, E)$ and label index $L(v)$ for all $v \in V$, the space cost of PB-tree is $O(|L| \cdot h)$.

Proof. For each $v \in V$, we have $|L(v)|$ label entries, thus each $PB(v)$ has an index size $O(|L(v)| \cdot h)$. It is easy to see the size of PB-tree is $O(|L| \cdot h)$ where $|L|$ is the size of label index. \square

5.3.2 Predecessor and Successor Queries on PB-Tree

With the construction of PB-tree, we discuss the predecessor and successor queries on top of it. Given $PB(v_{i-1})$, we aim to find candidate v_i that contains w_i and $d_G(v_{i-1}, v_i)$ is smaller or equal to lD , or greater or equal to rD . As we know, $d_G(v_{i-1}, v_i)$ can be divided into two parts $d_G(v_{i-1}, o)$ and $d_G(o, v_i)$. Therefore, straightforwardly, we can apply predecessor and successor queries on $PB(o)$ for each pivot $o \in L(v_{i-1})$ with two bound network distances lD_o and rD_o , respectively. Therefore, for $d_G(o, v_i)$,

$$\begin{aligned} lD_o &= lD - d_G(v_{i-1}, o), \\ rD_o &= rD - d_G(v_{i-1}, o), \end{aligned} \quad (6)$$

$$\text{we have } d_G(o, v_i) \leq lD_o \text{ or } d_G(o, v_i) \geq rD_o.$$

For each $PB(o)$, we are supposed to obtain a temporary candidate. Through comparison, we can finally find the next candidate vertex v_i .

Fortunately, it is worth to note that we are not necessarily to access all $PB(o)$ s to process predecessor and successor

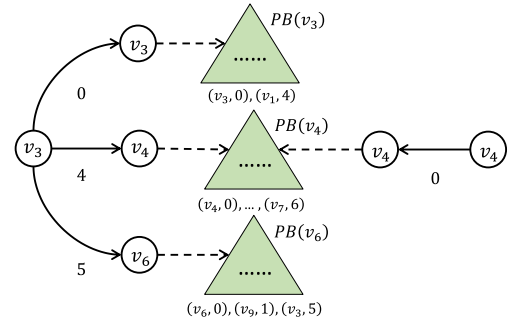


Fig. 4. Overview of pivot reverse binary tree.

queries. Basically, we know $d_m^i(v_i)$ must not exceed upper bound matching distance, therefore current UB can be utilized to prune the search space. That is to say, v_i could only be found if $d_G(v_{i-1}, v_i)$ is greater or equal to lB , or is smaller or equal to rB

$$\begin{aligned} lB &= d_i - d_i * UB, \\ rB &= d_i + d_i * UB, \end{aligned} \quad (7)$$

$$\text{we have } d_G(v_{i-1}, v_i) \geq lB \text{ or } d_G(v_{i-1}, v_i) \leq rB.$$

Particularly, for each $PB(o)$, the bound distances can be computed as

$$\begin{aligned} lB_o &= lB - d_G(v_{i-1}, o), \\ rB_o &= rB - d_G(v_{i-1}, o). \end{aligned} \quad (8)$$

Therefore, the search space can be narrowed down into $[lB_o, lD_o]$ and $[rD_o, rB_o]$. For current pivot o being processed, if we have $rB < d_G(v_{i-1}, o)$, we are impossible to find a candidate in $PB(o)$ since rB_o is negative. In other words, the network distance between v_{i-1} and any vertex in $PB(o)$ is definitely greater than rB thus is not qualified. As we know, the pivots in $L(v_{i-1})$ are sorted in ascending order of network distance, the rest pivots o' after o do not need to be considered since they have even greater network distances to v_{i-1} than o . Therefore, the process terminates.

Predecessor and Successor Queries. Given $PB(o)$, a query keyword w and two network distance bound ranges $[lB_o, lD_o]$ and $[rD_o, rB_o]$, we aim to find a temporary candidate vertex in $PB(o)$. In particular, the difference between AB-tree and PB-tree is that, given a query vertex u , any target v only shows up once in $AB(u)$, but it might appear in multiple $PB(o)$ s. Moreover, if we find a candidate v in $PB(o)$, $d_G(u, o) + d_G(o, v)$ is not necessarily equal to $d_G(u, v)$ since the network distance can only be calculated by the pivot with minimum distance summation. Therefore, we use original label index to check if $\mathcal{P}(u, \dots, o, \dots, v)$ is the shortest path $SP(u, v)$. As mentioned before, if $rB \geq d_G(v_{i-1}, o)$, we first apply a successor query on $PB(o)$. After we obtain a temporary vertex v_{tmp} locates in $[rD_o, rB_o]$, we check if o is on the shortest path $SP(v_{i-1}, v_{tmp})$ by comparing $d_G(v_{i-1}, v_{tmp})$ with $d_G(v_{i-1}, o) + d_G(o, v_{tmp})$. If so, v_{tmp} is reported as a temporary successor result on $PB(o)$. Otherwise, we update rD_o by $d_G(o, v_{tmp})$ and continue to apply a new successor query. This process is repeated until we find a result. After successor query, we compare $d_G(v_{i-1}, o)$ with lD to determine if we need to apply a predecessor query on $PB(o)$. Based on the same intuition, if $lD \geq d_G(v_{i-1}, o)$, the predecessor query is applied in a similar approach as successor query. Finally, we compare the results of predecessor and successor queries, and obtain the temporary candidate found in $PB(o)$. It is worth to

note that we can further narrow down the search space by updating LB and rB . That is, after processing pivot o , if we find a temporary candidate v_{tmp} , LB can be updated by $d_G(v_{i-1}, v_{tmp})$ and rB by $2 * d_i - LB$, which benefits the processing of rest o' .

Lemma 4. *The expected number of nodes visited in a predecessor or successor query on PB-tree is $O(\log \frac{|L|}{|V|} \cdot \frac{|W| \cdot |\Phi(V)|}{|V| \cdot h \cdot |V_{w_k}|})$.*

Proof. According to the proof in Lemma 2, the expected number of false positives before a true positive is $\frac{1}{1-p} = \frac{|W| \cdot |\Phi(V)|}{|V| \cdot h \cdot |V_{w_k}|}$. As we know, the height of each PB-tree is $\log \frac{|L|}{|V|}$. Therefore, expected number of nodes visited in a predecessor or successor query on PB-tree is $O(\log \frac{|L|}{|V|} \cdot \frac{|W| \cdot |\Phi(V)|}{|V| \cdot h \cdot |V_{w_k}|})$. Given that \wedge operation on two h length hash codes costs $O(h)$, thus the time of searching PB-tree is $O(\log \frac{|L|}{|V|} \cdot \frac{|W| \cdot |\Phi(V)|}{|V| \cdot h \cdot |V_{w_k}|}) \cdot O(h)$. \square

Example 6. In the running example, given query $\mathcal{Q} = (v_7, \{(w_1, 6, 0.5), (w_2, 4, 0.5), (w_4, 5, 0.5)\})$, assume we already have stack (v_7, v_4, v_3) . At level 2, we intend to find the next candidate. Initially, θ is set as 0.0, therefore we have $LD = rD = 5$. As current $UB = 0.4$, we have $LB = 4$ and $rB = 6$. As shown in Fig. 4, we first check $PB(v_3)$ with $d_G(v_3, v_3) = 0$. Then we have $LD_{v_3} = rD_{v_3} = 5$, $LB_{v_3} = 4$ and $rB_{v_3} = 6$. A successor query is applied and no vertex is found, and a predecessor query returns v_1 . As $d_G(v_3, v_1) = 4$ does not exceed LB_{v_3} , v_1 is taken as the temporary result for pivot v_3 . Then we continue to search $PB(v_4)$ with $LD_{v_4} = rD_{v_4} = 1$, $LB_{v_4} = 0$ and $rB_{v_4} = 2$ but no vertex is found, neither in $PB(v_6)$. Finally, we report v_1 with $d_m^3(v_1) = 0.4$.

6 DYNAMIC MAINTENANCE

In this section, we discuss how to maintain the PB-tree for road network updating. To avoid recomputing the index structure from scratch, we propose a semi-dynamic mechanism to adjust the PB-tree with a low overhead. As we know, PB-tree is built based on label index, thus the updating is divided into two phases, the updating of label index and the updating of PB-tree. Instead of recomputing a new label index, [18] introduces a dynamic label index scheme for distance queries on time-evolving graphs, and we adopt the algorithm for the first phase label index updating.

6.1 Semi-Dynamic Index Structure

Basically, we have 4 operations to update the network: insert a new vertex with an edge connecting to an existing vertex, delete a vertex with only one edge, insert an edge and delete an edge. As the deletion operation is much harder than insertion, and it seems impossible to find an efficient approach to support deletion in label generation. Moreover, it is rare to see deletion happens in road networks, thus we only take insertion into consideration. As the newly updated vertex is isolated, its label can be viewed as an empty set. Inserting a new vertex can be easily done by inserting an edge connecting to it, thus we only need to focus on edge insertion. As keyword updating is easy to implement, thus we omit it here.

Label Index Updating. Assume we insert an edge (a, b) into G , some shortest paths in old network may change by passing (a, b) . Based on the label generation algorithm, we do not have to remove outdated distances in label but resume BFSs of affected vertices and add new label entries into

index. It is worth to note that only the pivots in $L(a)$ and $L(b)$ are affected by network updating, and it suffices to conduct resumed BFSs originally rooted at pivot v_k if $v_k \in L(a) \cup L(b)$. Different from previous pruning method, a prefix pruning method is proposed to apply in BFS with a new parameter k , where k is the vertex ordering of v_k . The prefix method is to answer the distance query between v_k and u from the pivots in $L(v_k) \cap L(u)$ whose vertex orderings are at most k . Interested readers can refer to [18] for algorithm details.

Algorithm 5. Procedure findNext() with PB-Tree

Input: Query vertex v_{i-1} , clue w_i and d_i , threshold θ
Output: Next candidate v_i with $d_m^i(v_i)$

- 1 $LD \leftarrow d_i - d_i \cdot \theta$; $rD \leftarrow d_i + d_i \cdot \theta$;
- 2 $LB \leftarrow d_i - d_i \cdot UB$; $rB \leftarrow d_i + d_i \cdot UB$;
- 3 **for each** pivot $o \in L(v_{i-1})$ **do**
- 4 Obtain $PB(v_{i-1})$, LD_o , rD_o , LB_o and rB_o ;
- 5 **if** $d_G(v_{i-1}, o) > rB$ **then**
- 6 **break**;
- 7 **else**
- 8 $rD_o \leftarrow rD - d_G(v_{i-1}, o)$;
- 9 **while** $PB(v_{i-1}).succ(rD_o, w_i)$ and $d_G(o, v_{tmp_r}) \leq rB_o$
- 10 **do**
- 11 Obtain v_{tmp_r} ;
- 12 **if** $d_G(v_{i-1}, v_{tmp_r}) \neq d_G(v_{i-1}, o) + d_G(o, v_{tmp_r})$ **then**
- 13 $rD_o \leftarrow d_G(o, v_{tmp_r})$;
- 14 **else**
- 15 Obtain temp suc result on $PB(o)$;
- 16 **break**;
- 17 **if** $d_G(v_{i-1}) < LD$ **then**
- 18 $LD_o \leftarrow LD - d_G(v_{i-1}, o)$;
- 19 **while** $PB(v_{i-1}).pred(LD_o, w_i)$ and $d_G(o, v_{tmp_l}) \geq LB_o$
- 20 **do**
- 21 Obtain v_{tmp_l} ;
- 22 **if** $d_G(v_{i-1}, v_{tmp_l}) \neq d_G(v_{i-1}, o) + d_G(o, v_{tmp_l})$
- 23 **then**
- 24 $LD_o \leftarrow d_G(o, v_{tmp_l})$;
- 25 **else**
- 26 Obtain temp pre result on $PB(o)$;
- 27 **break**;
- 28 **if** $d_i - d_G(v_{i-1}, v_{tmp_l}) \leq d_G(v_{i-1}, v_{tmp_r}) - d_i$ **then**
- 29 $LB \leftarrow d_G(v_{i-1}, v_{tmp_l})$; $rB \leftarrow 2 * d_i - LB$;
- 30 $v_i \leftarrow v_{tmp_l}$;
- 31 **else**
- 32 $rB \leftarrow d_G(v_{i-1}, v_{tmp_r})$; $LB \leftarrow 2 * d_i - rB$;
- 33 $v_i \leftarrow v_{tmp_r}$;
- 34 **return** v_i with $d_m^i(v_i)$;

Pivot-Based Forest. To propose a semi-dynamic index structure, we present a general framework to convert PB-tree into pivot-based forest (PF), which is inspired by the *logarithmic method* [19]. Given $PB(o)$ with m label entries, we divide it into $l = \lfloor \log m \rfloor + 1$ partitions P_0, \dots, P_{l-1} . Each partition P_i either has 2^i label entries or is empty. We first compute a l -bit binary value of m . Interestingly, whether P_i is empty or not is determined by the i th bit, if i th bit is 0 then P_i is empty. For non-empty P_i , we follow the method introduced in Section 5.2.1 to construct a binary tree $PF(o)_i$ on these 2^i label entries. Finally, all these binary trees together form the pivot-based forest structure.

PF Index Updating. After label index updating, we add new label entries or rewrite distances of existing label entries.

TABLE 2
 Statistics of Dataset

	Beijing	New York
$\# V $	168,535	264,346
$\# E $	196,307	733,846
$\# W $	88,910	102,450
$\# \Phi(V) $	1,445,824	3,086,166

TABLE 3
 Parameter Settings

Parameters	Values
Dataset cardinality	4 K, 8 K, 12 K, 16 K
The number of clues	2, 3, 4, 5, 6, 7, 8
h bits hash code	32, 64 , 128, 256, 512
Average keyword frequency	10, 50, 100, 500 , 1,000, 5,000, 10,000
Average distance (km)	2, 4, 6, 8, 10, 12, 14, 16, 18, 20
Average confidence factor	0.2 , 0.4, 0.6, 0.8, 1.0

Assume we add a new label entry $(v, d_G(o, v))$ into $PB(o)$, we first find the smallest i such that $PF(o)_i$ is empty. If i equals to 0, we simply build $PF(o)_0$ with only one label entry $(v, d_G(o, v))$. Otherwise, we union all label entries of $PF(o)_0, \dots, PF(o)_{i-1}$, together with $(v, d_G(o, v))$, into $PF(o)_i$. It is worth to note that $PF(o)_i$ now has 2^i elements and $PF(o)_0, \dots, PF(o)_{i-1}$ become empty. As we know, the label entries in original $PB(o)$ are sorted in ascending order of distance. In $PF(o)$, we do not consider the global distance order but instead consider a local order in each $PF(o)_i$ when we rebuild the index. To rewrite distances of existing label entries, we only need to update the $PF(o)_i$ they belong to.

Query Processing on PF Index. Given query vertex v_{i-1} and a clue $\mu(w_i, d_i)$, we introduce how to answer `findNext()` on PF index. As we know, both the predecessor and successor queries are decomposable. Therefore, we simply apply the predecessor and successor queries on all non-empty $PF(o)_i$. Fortunately, it is not necessary to process queries on all $PF(o)_i$ s. If the query distance is smaller than the minimum network distance stored in $PF(o)_i$, the predecessor query is not required, where the similar case holds for successor query. Finally, we merge these intermediate results to obtain the result.

7 EXPERIMENTS

In this section, we conduct extensive experiments on real road network datasets to study the performance of the proposed index structures and algorithms.

7.1 Experimental Settings

All these algorithms introduced in this paper were implemented in GNU C++ on Linux and run on an Intel(R) CPU i7-4770@3.4 GHz and 32 G RAM.

Datasets. We use two real datasets, the road network datasets of Beijing and New York City from the 9th DIMACS Implementation Challenge.¹ Each dataset contains an undirected weighted graph that represents a part of the road network. The weight of each edge in a graph represents the distance between two endpoints of the edge. We obtain the keywords of vertices from the OpenStreetMap.² As shown in Table 2, for D1 in Beijing, we have 168,535 vertices and

TABLE 4
 Statistics of Keyword Frequency on Beijing Dataset

Frequency \geq	10	50	100	500	1,000	5,000	10,000
Percentage (%)	9.72	4.55	3.35	2.25	1.11	0.82	0.61

TABLE 5
 Performance of Proposed Algorithms and Index Structures on QT (Query Time), IS (Index Size), and IT (Index Time)

Algorithm	QT (ms)		IS (GB)		IT (min)		
	BJ	NY	BJ	NY	BJ	NY	
GCS	14.02	32.17	-	-	-	-	
CDP	Allpair	223.81	223.75	106.1	260.6	-	-
	Label	612.72	693.42	0.51	0.78	-	-
BAB	AB-tree	120.84	153.92	856	2,104	1,045	2,569
	PB-tree	76.12	93.27	2.1	3.21	2	3.1
	PF	89.86	106.14	2.2	3.36	2	3.4

196,307 edges. We also have 88,910 distinct keywords contained by vertices with the total occurrence 1,445,824. For D2 in New York, we have more vertices and edges than D1 in road network with almost twice the size of D1, and the number of keywords contained is larger than D1 as well.

Algorithms. We evaluate the performance of three algorithms, greedy clue search algorithm, clue-based dynamic programming algorithm and branch-and-bound algorithm. In CDP, we use two different distance oracles \mathcal{DO} to compute network distance, i.e., all-pair and 2-hop label. In BAB, we evaluate the performances of three index structures, i.e., AB-tree, PB-tree and PF.

Parameter Settings. We randomly generate 100 queries for each set of experiment and measure their performance by average. To evaluate the algorithms under various settings, we vary the value of some parameters in the query to study the performance, as shown in Table 3. For default settings, we choose 16 K for dataset cardinality (the number of vertices), 4 for the number of clues in query, and 64 for hash code length.

We assume a keyword at most shows up once in a vertex, thus the frequency of a keyword w is the number of vertices that contain w , i.e., $|V_w|$. The statistics of keyword frequency are shown in Table 4, which demonstrates the percentages of keywords with different frequencies. In the query, the keyword frequencies, the average distances and the confidence factors ϵ in clues are randomly generated. More specifically, assume the average keyword frequency for evaluation is $|V_w|$, thus the keyword frequencies we chose in a clue are in the range $[0.9 * |V_w|, 1.1 * |V_w|]$, which is the similar case with average distance and average ϵ .

7.2 Performance Evaluation

Table 5 shows the performance comparison of proposed algorithms and index structures on query time, index size and index construction time. The construction time of all-pair and 2-hop label, which have been studied by existing works, are excluded in our performance comparison. For the query time evaluation, it is easy to see that BAB well outperforms GCS and CDP. Besides, applying all-pair in CDP has a shorter response time but a larger space cost than utilizing 2-hop label, and using PB-tree in BAB has a better performance than using AB-tree and PF. For index size and construction time, label based approaches have a much smaller size and less time than all-pair based approaches. As NY has a larger size than BJ, more time and

1. <http://www.dis.uniroma1.it/challenge9/download.shtml>
 2. <http://www.openstreetmap.org>

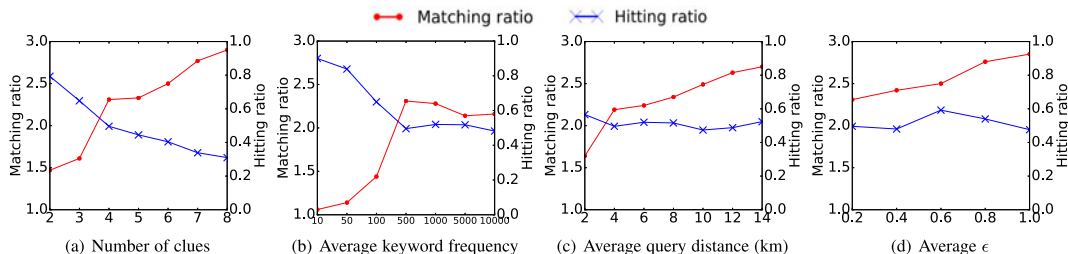


Fig. 5. Accuracy of GCS.

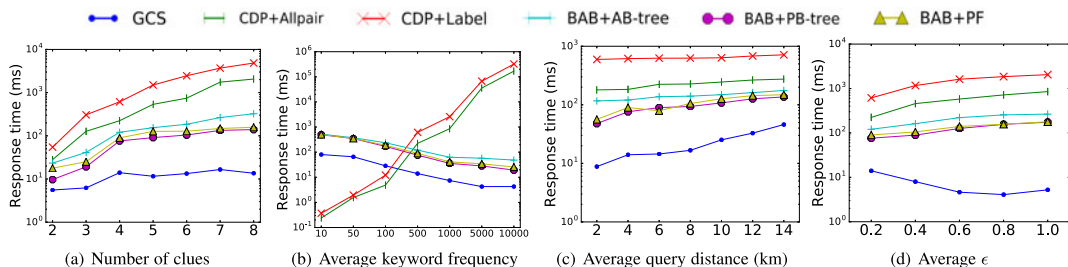


Fig. 6. Query time.

space costs are required. For the rest experiments, we only demonstrate the performance on BJ due to the space limit, where the performance on NY is similar to that on BJ.

7.2.1 Accuracy Measurement of GCS

Fig. 5 shows the accuracy measurement of GCS by varying the parameters in the query, such as the number of clues, average keyword frequencies, expected distances and ϵ . We study the accuracy by two criteria: the matching ratio A_{match} and hitting ratio A_{hit}

$$A_{match} = \frac{d_m(\mathcal{C}, \mathcal{FP}_{gcs})}{\text{optimal matching distance } d_{m,opt}}$$

A_{match} is the ratio of estimated matching distance of GCS on the optimal matching distance. A smaller A_{match} means a better accuracy

$$A_{hit} = \frac{|\mathcal{FP}_{gcs} \cap \mathcal{FP}_{opt}|}{|\mathcal{C}|}$$

A_{hit} focuses on the percentage of match vertices in \mathcal{FP}_{gcs} contained by the optimal feasible path \mathcal{FP}_{opt} . A greater A_{hit} means a better accuracy.

As we can see in Fig. 5, when we enlarge the parameters, A_{match} increases as the result of GCS becomes more inaccurate; and A_{hit} decreases, as less match vertices in optimal feasible path are hit. Both these two criteria becomes less sensitive when the keyword frequency gets larger than 500 in Fig. 5b. Moreover, A_{hit} keeps steady in Figs. 5c and 5d since GCS is not sensitive to average distance and ϵ .

7.2.2 Query Efficiency Measurement

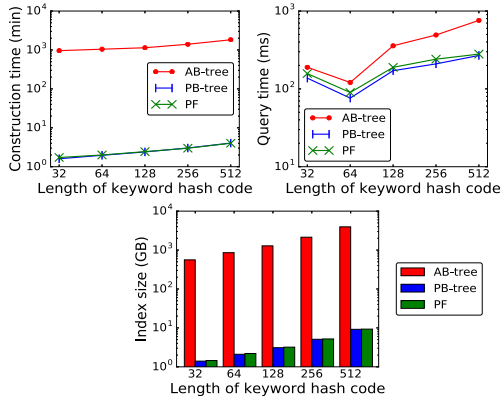
Effect of the Keyword Hash Code Length h . In this set of experiments, we study the effect of keyword hash code length h on performance of AB-tree, PB-tree and PF index structures. As shown in Fig. 7, the pivot-based indices well outperform AB-tree on index construction time, index size and query time. The space of AB-tree is $O(|V|^2 \cdot h)$ and PB-tree is $O(|L| \cdot h)$. When we enlarge h , both the index size and construction time linearly increase. For query time, there are more false

positives in tree traversal when h is small, however, the bit operation costs less time than larger h , which is the case when h is set to 64 comparing with higher values. When we set h to 32, even though we have less bit operation costs, the number of false positives increases such that the query time increases.

Effect of the Dataset Cardinality. Let us take Beijing dataset for example, we randomly extract 5 subgraphs with equal number of vertices from the original Beijing road network where the performance is measured by average for each experiment. For each subgraph, the connectivity and keyword information of each vertex are kept the same as the original road network. In this set of experiments, we vary the size of these subgraphs to study the performance of proposed algorithms and index structures, as shown in Fig. 8. Obviously, the index size and construction time increase when we enlarge the size of datasets. It is worth to note that the size of AB-tree increases exponentially with the number of vertices, and the sizes of PB-tree and PF increase gently especially when the size is enlarged from 120 to 160 K due to the property of 2-hop label. For the query time, the BAB algorithm outperforms CDP by a large margin.

Effect of the Number of Clues. In this set of experiments, Fig. 6a shows the performance of algorithms by increasing the number of clues in CRS query. Not surprisingly, the response time increases when we enlarge the number of clues of all proposed algorithms. For GCS, the response time increases gently since only more rounds of network expansion are induced. For CDP, when we enlarge the number of clues, more iterations are triggered for the computation. For BAB, the number of candidate vertices and feasible paths increase thus takes more computation time.

Effect of the Average Frequency of Keywords. In this set of experiments, we study the performance of algorithms by varying the frequency of query keywords, as shown in Fig. 6b. It suffices to say that for low frequency keywords, say the frequency less than 500, it is more efficient if we adopt CDP with all-pair, and for high frequency keywords, BAB with PB-tree has a much better performance on both response time and index size. This is because, for CDP, there are not too many combinations to consider if the frequency is low, but when we enlarge the frequency, the response time increases


 Fig. 7. Effect of the keyword hash code length h .

exponentially to the frequency. For BAB, there are lots of false positives if the frequency is low, and when we enlarge the frequency, the performance becomes much better since we can quickly locate the candidate by using PB-tree.

Effect of the Average Expected Distance. In this set of experiments, we study the effect of average expected distance on the performance of proposed algorithms, as shown in Fig. 6c. As we know, we apply the network expansion algorithm in GCS, which makes it sensitive to the expected distance. When the distance increases, more vertices are involved that results in more computation cost. For CDP with all-pair or label index, they both have a small dependency on the query distance. Therefore, the computation time of CDP keeps almost steady as the distance increases. For BAB, the effect is still not obvious but if the distance is small, we are supposed to find the next candidate more quickly since there are only a small portion of vertices after filtered by distance.

Effect of the Average ϵ . In this set of experiments, we study the effect of average ϵ on the performance of proposed algorithms, as shown in Fig. 6d. When we enlarge the average ϵ , more match vertices are considered as candidates, thus the time costs of CDP and BAB increase. For GCS, we can do less network traversal to find the current best match vertex, so the query time reduces when we enlarge ϵ .

Evaluation of Index Updating. Here we evaluate the cost of index updating in Table 6. It is easy to observe that the average update time cost is much smaller than reconstruction the index from scratch. The cost comes from two parts, the updating of label index and updating of PF. For each update, we only have to update a very small number of pivot forest structures, that is, the semi-dynamic update is done locally.

8 RELATED WORK

In this section, we introduce two lines of related work, top- k spatial keyword search and travel route search.

8.1 Top- k Spatial Keyword Search

Searching geo-textual objects with query location and keywords has gained increasing attention recently due to the popularity of location-based services. In Euclidean space, *IR²-tree* [20] integrates signature files and *R-tree* to answer boolean keyword queries. *IR-tree* [21] is an *R-tree* augmented with inverted files that supports the ranking of objects based on a score function of spatial distance and text relevancy. Cao et al. [22] proposes a location-aware top- k prestige-based text retrieval (*LkPT*) query, to retrieve the top- k spatial web objects ranked according to both prestige-based text relevance (PR) and location proximity. Chen et al. [23] provides an all-round

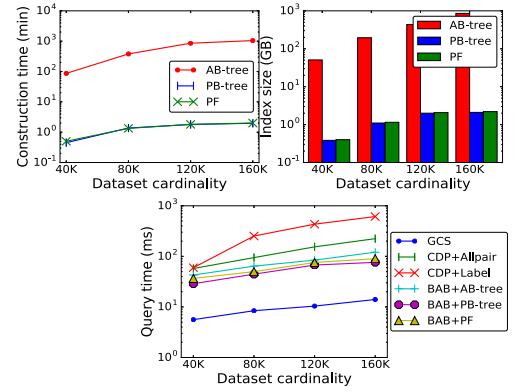


Fig. 8. Effect of the dataset cardinality.

survey of 12 state-of-art geo-textual indices and proposes a benchmark that enables the comparison of the spatial keyword query performance. Zhang et al. [6], [24] proposes the m closet keyword query (mCK query) which aims to find the closest objects that match the query keywords and their distance diameter is minimized. Recently, Guo et al. [5] propose approximation algorithms to solve the mCK query with a ratio of $(\frac{2}{\sqrt{3}} + \epsilon)$. Cao et al. [3] propose a collective spatial keyword query, in which a different semantics is taken such that the group of objects in the result covers the query keywords and has the lowest cost. Li et al. [25] studies the problem of direction-aware spatial keyword search, which aims at finding the k nearest neighbours to the query that contain all input keywords and satisfy the direction constraint. Rocha et al. [26] address the problem of processing top- k spatial keyword queries on road networks where the distance between the query location and the spatial object is the length of shortest path. ROAD [27] organizes the road network as a hierarchy of subgraphs, and connects them by adding shortcuts. For each subgraph, an object abstract is generated for keyword checking. By using network expansion, the subgraphs without intended object are pruned out. G-tree [28] adopts a graph partitioning approach to form a hierarchy. Within each subgraph, a distance matrix is kept, and for any two subgraphs, the distances between all borders of them are stored as well. Based on these distances, it efficiently computes the distance between query vertex and target vertices or tree nodes. Jiang et al. [29] adopt 2-hop label for handling the distance query for k NN problem on large networks, and facilitates KT index to handle the performance issue of frequent keywords. Liu et al. [7] formalize the spatio-textual context information of the querying POI as *clues* and use them to identify the most promising POIs, which is closely related to our CRS problem. Different with their work, we aim to find a feasible route on road networks by using clues. In addition, the spatial distance considered in our work is network distance so that the algorithms in [7] can not be applied.

8.2 Travel Route Search

The travel route search problem has been substantially studied for decades. Traveling Salesman Problem (TSP) [30] is the

TABLE 6
Evaluation of Index Updating

Dataset	Update time	Updated pivots
Beijing	78 ms	3.6
NY	127 ms	5.7

most classic problem in route planning. TSP aims to find the round trip that has the minimum cost from a source point to a set of targets. Li et al. [8] study the problem of Trip Planning Query (TPQ) in spatial databases, where each object is associated with a location and a category. With a starting point S , a destination E and a set of categories C , TPQ retrieves the best trip that starts at S passes through at least one point from each category, and ends at E . TPQ can be considered as a generalization of Travelling Salesman Problem (TSP), thus two approximation algorithms are proposed. Sharifzadeh et al. [10] studies the problem of optimal sequenced route (OSR), which aims to find a route of minimum length starting from a source point and passing through a number of typed locations in a specific sequence imposed on the types of the locations. They propose a LORD and R-LORD algorithms to filter out the locations that cannot be in the optimal route, thus improves the search efficiency. Chen et al. [31] studies the problem of multi-rule partial sequence route (MRPSR), which aims to find an optimal route with minimum distance under some partial category order rules defined in the query. They propose three heuristic algorithms to search for near-optimal solutions for the MRPSR query. Kanza et al. [32] proposes a greedy algorithm to find a route whose length is smaller than a specified threshold while the total text relevance of this route is maximized. Kanza et al. [33] studies the problem of finding a route that visits at least one satisfying entity of each type in an interactive approach. In each step, a candidate is given to user to provide a feedback specifying whether the entity satisfies her. Yao et al. [34] studies the problem of multi-approximate-keyword routing query, which complements the standard shortest path search with multiple keywords and an approximate string similarity function. For each keyword, the matching point is supposed to have an edit distance smaller than a given threshold. Cao et al. [9] defines the problem of keyword-aware optimal route query, which is to find an optimal route such that it covers a set of user-specified keywords, a specific budget constraint is satisfied, and the objective score of the route is optimized. Li et al. [35] proposes two different solutions, namely backward search and forward search, to deal with the general optimal route query without a total order. Zhang et al. [36] proposes the problem of personalized trip recommendation, which aims to find the optimal trip that maximizes users' experiences for a given time budget constraint and also takes the uncertain traveling time into consideration.

9 CONCLUSION AND FUTURE DIRECTIONS

In this paper, we study the problem of CRS on road networks, which aims to find an optimal route such that it covers a set of query keywords in a given specific order, and the matching distance is minimized. To answer the CRS query, we first propose a greedy clue-based algorithm GCS with no index where the network expansion approach is adopted to greedily select the current best candidates to construct feasible paths. Then, we devise an exact algorithm, namely clue-based dynamic programming CDP, to answer the query that enumerates all feasible paths and finally returns the optimal result. To further reduce the computational overhead, we propose a branch-and-bound algorithm BAB by applying filter-and-refine paradigm such that only a small portion of vertices are visited, thus improves the search efficiency. In order to quickly locate the candidate vertices, we develop AB-tree and PB-tree structures to speed up the tree traversal, as well as a semi-dynamic

index updating mechanism. Results of empirical studies show that all the proposed algorithms are capable of answering CRS query efficiently, while the BAB algorithm runs much faster, and the index size of PB-tree is much smaller than AB-tree.

Several directions for future research are promising. First, users may prefer a more generic preference model, which combines PoI rating, PoI average menu price, etc, in the query clue. Second, it is of interest to take temporal information into account and further extend the CRS query. Each PoI is assigned with a opening hours time interval $[T_o, T_c]$, and each clue contains a visiting time t , where the resulting query aims to find a path such that the time interval of each matched PoI covers the visiting time. Third, requiring users to provide exact keyword match is difficult sometimes as they are just providing "clue", which may be inaccurate in nature. Thus, it is of interest to extend our model to support the approximate keyword match. Hence, the matching distance can be modified by incorporating both spatial distance and textual distance together through a linear combination.

ACKNOWLEDGMENTS

This research is partially supported by the Natural Science Foundation of China (Grants No. 61502324, No. 61532018, No. 61232006, and No. 61572215), and the Australian Research Council (Grants No. LP130100164).

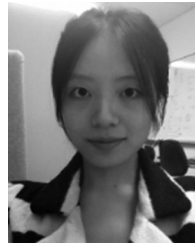
REFERENCES

- [1] B. Zheng, N. J. Yuan, K. Zheng, X. Xie, S. Sadiq, and X. Zhou, "Approximate keyword search in semantic trajectory database," in *Proc. IEEE Int. Conf. Data Eng.*, 2015, pp. 975–986.
- [2] K. Zheng, S. Shang, N. J. Yuan, and Y. Yang, "Towards efficient search for activity trajectories," in *Proc. IEEE Int. Conf. Data Eng.*, 2013, pp. 230–241.
- [3] X. Cao, G. Cong, C. S. Jensen, and B. C. Ooi, "Collective spatial keyword querying," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2011, pp. 373–384.
- [4] C. Long, R. C.-W. Wong, K. Wang, and A. W.-C. Fu, "Collective spatial keyword queries: A distance owner-driven approach," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2013, pp. 689–700.
- [5] T. Guo, X. Cao, and G. Cong, "Efficient algorithms for answering the m-closest keywords query," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2015, pp. 405–418.
- [6] D. Zhang, Y. M. Chee, A. Mondal, A. K. Tung, and M. Kitsuregawa, "Keyword search in spatial databases: Towards searching by document," in *Proc. IEEE Int. Conf. Data Eng.*, 2009, pp. 688–699.
- [7] J. Liu, K. Deng, H. Sun, Y. Ge, X. Zhou, and C. S. Jensen, "Clue-based spatio-textual query," *Proc. VLDB Endowment*, vol. 10, no. 5, pp. 529–540, 2017.
- [8] F. Li, D. Cheng, M. Hadjieleftheriou, G. Kollios, and S.-H. Teng, "On trip planning queries in spatial databases," in *Proc. 9th Int. Conf. Advances Spatial Temporal Databases*, 2005, pp. 273–290.
- [9] X. Cao, L. Chen, G. Cong, and X. Xiao, "Keyword-aware optimal route search," *Proc. VLDB Endowment*, vol. 5, no. 11, pp. 1136–1147, 2012.
- [10] M. Sharifzadeh, M. Kolahdouzan, and C. Shahabi, "The optimal sequenced route query," *VLDB J.*, vol. 17, no. 4, pp. 765–787, 2008.
- [11] I. Abraham, D. Delling, A. V. Goldberg, and R. F. Werneck, "Hierarchical hub labelings for shortest paths," in *Proc. 20th Annu. Eur. Conf. Algorithms*, 2012, pp. 24–35.
- [12] T. Akiba, Y. Iwata, and Y. Yoshida, "Fast exact shortest-path distance queries on large networks by pruned landmark labeling," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2013, pp. 349–360.
- [13] E. W. Dijkstra, "A note on two problems in connexion with graphs," *Numerische Mathematik*, vol. 1, no. 1, pp. 269–271, 1959.
- [14] M. Jiang, A. W.-C. Fu, R. C.-W. Wong, and Y. Xu, "Hop doubling label indexing for point-to-point distance querying on scale-free networks," *Proc. VLDB Endowment*, vol. 7, no. 12, pp. 1203–1214, 2014.

- [15] T. Akiba, Y. Iwata, K.-I. Kawarabayashi, and Y. Kawata, "Fast shortest-path distance queries on road networks by pruned highway labeling," in *Proc. Meeting Algorithm Eng. Experiments*, 2014, pp. 147–154.
- [16] C. S. Jensen, J. Kolářvr, T. B. Pedersen, and I. Timko, "Nearest neighbor queries in road networks," in *Proc. 11th ACM Int. Symp. Advances Geographic Inf. Syst.*, 2003, pp. 1–8.
- [17] B. Zheng, et al., "Keyword-aware continuous kNN query on road networks," in *Proc. IEEE Int. Conf. Data Eng.*, 2016, pp. 871–882.
- [18] T. Akiba, Y. Iwata, and Y. Yoshida, "Dynamic and historical shortest-path distance queries on large evolving networks by pruned landmark labeling," in *Proc. 23rd Int. Conf. World Wide Web*, 2014, pp. 237–248.
- [19] J. L. Bentley and J. B. Saxe, "Decomposable searching problems I. Static-to-dynamic transformation," *J. Algorithms*, vol. 1, no. 4, pp. 301–358, 1980.
- [20] I. De Felipe, V. Hristidis, and N. Rishe, "Keyword search on spatial databases," in *Proc. IEEE Int. Conf. Data Eng.*, 2008, pp. 656–665.
- [21] G. Cong, C. S. Jensen, and D. Wu, "Efficient retrieval of the top-k most relevant spatial web objects," *Proc. VLDB Endowment*, vol. 2, pp. 337–348, 2009.
- [22] X. Cao, G. Cong, and C. S. Jensen, "Retrieving top-k prestige-based relevant spatial web objects," *Proc. VLDB Endowment*, vol. 3, pp. 373–384, 2010.
- [23] L. Chen, G. Cong, C. S. Jensen, and D. Wu, "Spatial keyword query processing: An experimental evaluation," *Proc. VLDB Endowment*, vol. 6, pp. 217–228, 2013.
- [24] D. Zhang, B. C. Ooi, and A. K. Tung, "Locating mapped resources in web 2.0," in *Proc. IEEE Int. Conf. Data Eng.*, 2010, pp. 521–532.
- [25] G. Li, J. Feng, and J. Xu, "DESKS: Direction-aware spatial keyword search," in *Proc. IEEE Int. Conf. Data Eng.*, 2012, pp. 474–485.
- [26] J. B. Rocha-Junior and K. Nørvg, "Top-k spatial keyword queries on road networks," in *Proc. 15th Int. Conf. Extending Database Technol.*, 2012, pp. 168–179.
- [27] K. C. Lee, W.-C. Lee, and B. Zheng, "Fast object search on road networks," in *Proc. 12th Int. Conf. Extending Database Technol.*, 2009, pp. 1018–1029.
- [28] R. Zhong, G. Li, K.-L. Tan, and L. Zhou, "G-tree: An efficient index for KNN search on road networks," in *Proc. 22nd ACM Int. Conf. Inf. Knowl. Manage.*, 2013, pp. 39–48.
- [29] M. Jiang, A. W.-C. Fu, and R. C.-W. Wong, "Exact top-k nearest keyword search in large networks," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2015, pp. 393–404.
- [30] N. Christofides, "Worst-case analysis of a new heuristic for the travelling salesman problem," No. RR-388. Carnegie-Mellon Univ. Pittsburgh Pa Management Sciences Research Group, 1976.
- [31] H. Chen, W.-S. Ku, M.-T. Sun, and R. Zimmermann, "The multi-rule partial sequenced route query," in *Proc. 16th ACM SIGSPATIAL Int. Conf. Advances Geographic Inf. Syst.*, 2008, Art. no. 10.
- [32] Y. Kanza, E. Safra, Y. Sagiv, and Y. Doytsher, "Heuristic algorithms for route-search queries over geographical data," in *Proc. ACM SIGSPATIAL Int. Conf. Advances Geographic Inf. Syst.*, 2008, Art. no. 11.
- [33] Y. Kanza, R. Levin, E. Safra, and Y. Sagiv, "Interactive route search in the presence of order constraints," *Proc. VLDB Endowment*, vol. 3, no. 1/2, pp. 117–128, 2010.
- [34] B. Yao, M. Tang, and F. Li, "Multi-approximate-keyword routing in GIS data," in *Proc. ACM SIGSPATIAL Int. Conf. Advances Geographic Inf. Syst.*, 2011, pp. 201–210.
- [35] J. Li, Y. D. Yang, and N. Mamoulis, "Optimal route queries with arbitrary order constraints," *IEEE Trans. Knowl. Data Eng.*, vol. 25, no. 5, pp. 1097–1110, May 2013.
- [36] C. Zhang, H. Liang, K. Wang, and J. Sun, "Personalized trip recommendation with poi availability and uncertain traveling time," in *Proc. 24th ACM Int. Conf. Inf. Knowl. Manage.*, 2015, pp. 911–920.



Bolong Zheng received the bachelor's and master's degrees in computer science from the Huazhong University of Science and Technology, Wuhan, China, in 2011 and 2013, respectively, and the PhD degree from the University of Queensland, in 2017. He is currently a postdoctoral research fellow with the University of Queensland. His research interests include spatial database, spatial keyword search, and graph database.



Han Su received the BS degree in software engineering from Nanjing University, in 2011 and the PhD degree in computer science from the University of Queensland, in 2015. She is currently an associate professor in the Big Data Research Center, University of Electronic Science and Technology of China. Her research interests include trajectory querying and mining.



Wen Hua received the bachelor's and PhD degrees in computer science from the Renmin University of China, in 2010 and 2015, respectively. She is currently a lecturer with the University of Queensland. Her research interests include natural language processing, information retrieval and extraction, data mining, social media analysis, and graph management.



Kai Zheng received the PhD degree in computer science from the University of Queensland, in 2012. He is a full professor in the School of Computer Science, Soochow University. He has been working in the area of social-media analysis, spatial-temporal databases, uncertain databases, data mining, and bioinformatics. He has published more than 70 papers in the most prestigious journals and conferences such as SIGMOD, ICDE, the *VLDB Journal*, ACM Transactions and IEEE Transactions. He was the program committee co-chair of the 18th APWeb Conference and general co-chair of the 22nd DASFAA Conference. He is a PC member of several prestigious conferences such as ACM SIGMOD, ACM CIKM, ACM SIGSPATIAL, and DASFAA. He is the guest editor of the *Journal of WWW*, *Frontier of Computer Science*, and *Geoinformatica*. He was the recipient of the Australian Discovery Early Career Research Award in 2013, and Best Paper Award of the International Conference on Data Engineering (2015). He is a member of the IEEE.



Xiaofang Zhou received the bachelor's and master's degrees in computer science from Nanjing University, in 1984 and 1987, respectively, and the PhD degree in computer science from the University of Queensland in 1994. He is a professor of computer science with the University of Queensland. He is the head of the Data and Knowledge Engineering Research Division, School of Information Technology and Electrical Engineering. He is also a specially appointed adjunct professor with Soochow University, China. His research is focused on finding effective and efficient solutions to managing integrating, and analyzing very large amounts of complex data for business and scientific applications. His research interests include spatial and multimedia databases, high performance query processing, web information systems, data mining, and data quality management. He is a senior member of the IEEE.



Guohui Li is a full professor in the School of Computer Science and Technology, Huazhong University of Science and Technology (HUST). His research interests include big data processing, data mining, social media data processing, and real-time computing. He is widely published in the leading scholarly journals and conferences in computer science, such as RTSS, AAAI, the *IEEE Transactions on Computers*, and the *IEEE Transactions on Mobile Computing*. He is a member of the IEEE.

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