

# G-Tree: An Efficient and Scalable Index for Spatial Search on Road Networks

Ruicheng Zhong, Guoliang Li, Kian-Lee Tan, Lizhu Zhou, and Zhiguo Gong

**Abstract**—In the recent decades, we have witnessed the rapidly growing popularity of location-based systems. Three types of location-based queries on road networks, single-pair shortest path query,  $k$  nearest neighbor ( $k$ NN) query, and keyword-based  $k$ NN query, are widely used in location-based systems. Inspired by R-tree, we propose a height-balanced and scalable index, namely G-tree, to efficiently support these queries. The space complexity of G-tree is  $\mathcal{O}(|V|\log|V|)$  where  $|V|$  is the number of vertices in the road network. Unlike previous works that support these queries separately, G-tree supports all these queries within one framework. The basis for this framework is an assembly-based method to calculate the shortest-path distances between two vertices. Based on the assembly-based method, efficient search algorithms to answer  $k$ NN queries and keyword-based  $k$ NN queries are developed. Experiment results show G-tree's theoretical and practical superiority over existing methods.

**Index Terms**—Single-pair shortest path, KNN search, keyword search, road network, index, spatial databases

## 1 INTRODUCTION

NOWADAYS, mobile devices have become more and more popular in our daily life. As existing mobile devices are usually equipped with global positioning system (GPS) chips, users can easily get their locations. To provide users with location-based services (LBS), many LBS systems (e.g., foursquare and Google Maps) have been deployed and widely accepted by mobile users. In our daily life, there are many LBS applications on road networks. For example, in navigation systems, a user may want to find the shortest path between two vertices on a road network. In tourist guide applications, a tourist may look for  $k$  nearest “seafood restaurants” while walking in a city. There are three widely-used types of queries on road networks, single-pair shortest path query (SPSP),  $k$  nearest neighbor ( $k$ NN) query, and keyword-based  $k$ NN ( $k^2N^2$ ) query.

There are many existing works to support these three types of queries. [1], [5], [9], [10], [20], [21] studied the problem of shortest path queries between two vertices, [2], [6], [7], [12], [13], [16], [19] addressed the  $k$ NN search problem, and [14], [17], [18] worked on the problem of keyword-based  $k$ NN query on road networks. Though they achieve good results for individual type of queries, there are still limitations to be addressed for them to be of practical use. First, some are not efficient enough, especially on  $k$ NN

query (e.g., [7], [12], [13], [16]). Second, some fail to scale up to very large datasets (e.g., [6], [13], [19]). Third, they cannot be adaptive to different types of queries. For example, if a LBS system needs to support shortest-path query,  $k$ NN query and  $k^2N^2$  query at the same time, it has to exploit three different methods or indices, and the overhead is very costly. To this end, we aim to design an efficient and scalable index on road networks, which can not only handle very large datasets effectively, but also support various spatial queries on road networks.

However, to our best knowledge, there is no such index on road networks. Inspired by R-tree on the metric space, we devise a similar tree-based index on road networks, namely G-tree. We recursively partition the road network into sub-networks and build a tree structure on top of sub-networks where each G-tree node corresponds to a sub-network. The space complexity of our G-tree structure is only  $\mathcal{O}(|V|\log|V|)$ , which is much better than the state-of-the-art method [19] with  $\mathcal{O}(|V|^{1.5})$  complexity, where  $|V|$  is the number of vertices in the road network.

Using the G-tree index, we devise effective algorithms to support three types of spatial queries on road networks, single-pair shortest path query,  $k$  nearest neighbor query and keyword-based  $k$ NN query ( $k^2N^2$ ), which are fundamental and popular for a variety of applications. For SPSP, we devise an assembly-based algorithm to compute the shortest-path distance between two vertices with  $\mathcal{O}(|V|)$  time complexity. For  $k$ NN, we develop an efficient top- $k$  search scheme and devise effective pruning techniques based on the minimum boundary between the query location and a G-tree node. We also combine the spatial and textual information together to handle  $k^2N^2$  query judiciously. Experiments on eight real-world datasets show that our method significantly outperforms state-of-the-art methods, even by 2-3 orders of magnitude. To summarize, we make the following contributions.

- We propose a balanced search tree index, G-tree, which has low space overhead but superior

- R. Zhong, G. Li, and L. Zhou are with the Department of Computer Science and Technology, Tsinghua National Laboratory for Information Science and Technology (TNList), Tsinghua University, Beijing 100084, China. E-mail: zrc1101001@gmail.com, [liguoliang, dcszlj]@tsinghua.edu.cn.
- K.-L. Tan is with the Department of Computer Science, National University of Singapore, Singapore. E-mail: tankl@comp.nus.edu.sg.
- Z. Gong is with the Department of Computer Science, University of Macau, Macau, China. E-mail: fstzgg@umac.mo.

Manuscript received 3 June 2014; revised 19 Dec. 2014; accepted 14 Jan. 2015. Date of publication 2 Feb. 2015; date of current version 2 July 2015.

Recommended for acceptance by W. Wang.

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.2015.2399306

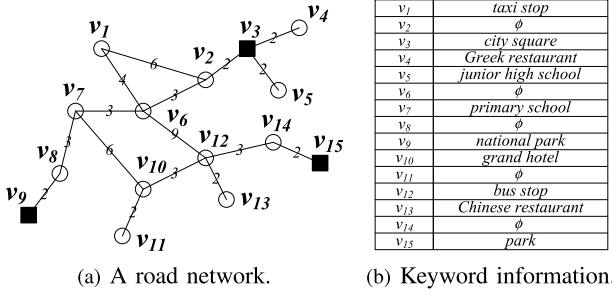


Fig. 1. An example of the road network.

performance and good scalability. The space complexity of our G-tree structure is only  $\mathcal{O}(|V|\log|V|)$ .

- We propose efficient search algorithms for three types of queries: SPSP,  $k$ NN, and  $k^2N^2$ . We devise an assembly-based algorithm to compute shortest-path distance for SPSP queries. Based on the algorithm and the G-tree, we devise a best-first search algorithm to support  $k$ NN queries on road networks. We also extend the  $k$ NN algorithm to support  $k^2N^2$  queries.
- Our method has theoretical and practical superiority over existing methods. Experimental results on real-world datasets show that G-tree significantly outperforms state-of-the-art methods and scales very well.

The structure of this paper is organized as follows. We define three types of queries and review related works in Section 2. We present our G-tree index and search algorithms in Sections 3 and 4 respectively. We discuss the path recovery issue in Section 5. Miscellaneous issues are discussed in Section 6. Experiments are reported in Section 7 and we conclude in Section 8.

## 2 PRELIMINARIES

### 2.1 Problem Formulation

**Data model.** We model a road network as an undirected weighted graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is a set of vertices and  $\mathcal{E}$  is a set of edges. Each edge  $(u, v)$  in  $\mathcal{E}$  has a weight  $w(u, v)$ , e.g., road distance, travel time, etc., which is a positive value. Given a path between vertex  $u$  and vertex  $v$ , the sum of weights of edges along the path is called the distance of the path. Among all paths between vertexes  $u$  and  $v$ , the one with the minimum distance is called the shortest path. Let  $SP(u, v)$  denote a shortest path between  $u$  and  $v$ , and  $SPDist(u, v)$  denote the shortest-path distance between  $u$  and  $v$ . We will discuss how to extend our method to support directed graphs in Section 6.2. Each vertex  $v$  contains a list of keywords, denoted by  $v.\mathcal{W}$ , which is the textual description of the vertex.

For example, Fig. 1a shows a road network. The weight of edge  $(v_2, v_6)$  is 3.  $SP(v_4, v_9) = v_4v_3v_2v_6v_7v_8v_9$  is a shortest path between  $v_4$  and  $v_9$  and  $SPDist(v_4, v_9) = 15$ . Fig. 1b displays the textual description of each vertex on the road network in Fig. 1a.

**Query model.** We consider three types of queries: single-pair shortest path query,  $k$  nearest neighbors query and keyword-based  $k$ NN query ( $k^2N^2$ ).

1) *Single-pair shortest path query.* Given a graph  $\mathcal{G}$  and a query  $q = \langle u, v \rangle$ , SPSP returns  $SPDist(u, v)$ . For example,

consider a query  $q = \langle v_4, v_9 \rangle$  in Fig. 1a. we can figure out that  $SP(v_4, v_9) = v_4v_3v_2v_6v_7v_8v_9$ , thus  $SPDist(v_4, v_9) = 15$ .

2)  *$k$  nearest neighbor query.* Given a graph  $\mathcal{G}$  and a query  $q = \langle v_q, \mathcal{C}, k \rangle$ , where  $v_q$  is the query vertex,  $\mathcal{C}$  is a set of object vertices, and  $k$  is an integer, the answer  $\mathcal{R}$  is a set of  $k$  nearest objects to the query location such that,

- 1) The size of  $\mathcal{R}$  is  $k$ , i.e.,  $|\mathcal{R}| = k$ ;
- 2) Each answer is an object, i.e.,  $\mathcal{R} \subseteq \mathcal{C}$ ;
- 3)  $\forall v \in \mathcal{R}, u \in \mathcal{C} - \mathcal{R}, SPDist(v_q, v) \leq SPDist(v_q, u)$ .

For simplicity, in the paper we assume that both the query location and the objects are at vertices. If the query location is on an edge, we use the two end vertices of the edge to do  $k$ NN search and merge the answer sets of the two vertices to generate the final results. If the query location is outside the graph, we first find its nearest edge (by euclidean distance) and then use the nearest location on the edge to find the final results. Similar variations can be applied to the object.

For example, in Fig. 1a, consider a query  $q = \langle v_4, \{v_3, v_9, v_{15}\}, 2 \rangle$ . There are three objects (which are denoted by solid rectangles). We have  $SPDist(v_4, v_3) = 2$ ,  $SPDist(v_4, v_9) = 15$  and  $SPDist(v_4, v_{15}) = 21$ . The top-2 answers are  $\mathcal{R} = \{v_3, v_9\}$ .

3) *Keyword-based  $k$ NN query.* Given a graph  $\mathcal{G}$  and a query  $q = \langle v_q, \mathcal{C}, \mathcal{W}, k \rangle$  and a ranking function  $\mathcal{F}$ , where  $\mathcal{W}$  is a set of keywords, the answer  $\mathcal{R}$  is a set of top- $k$  objects in  $\mathcal{C}$  based on the ranking function  $\mathcal{F}$  such that,

- 1) The size of  $\mathcal{R}$  is  $k$ , i.e.,  $|\mathcal{R}| = k$ ;
- 2) Each answer is an object, i.e.,  $\mathcal{R} \subseteq \mathcal{C}$ ;
- 3)  $\forall v \in \mathcal{R}, u \in \mathcal{C} - \mathcal{R}, \mathcal{F}(v_q, v) \geq \mathcal{F}(v_q, u)$ .

A general ranking function  $\mathcal{F}$  considers not only the spatial proximity but also textual relevancy between the queries and objects, which is defined as below

$$\mathcal{F}(q, v) = \alpha \left( 1 - \frac{SPDist(q, v)}{M_S} \right) + (1 - \alpha) \frac{TXDist(q, v)}{M_T}, \quad (1)$$

where  $M_S$  is the maximum spatial distance,  $M_T$  is the maximum textual distance, and  $TXDist(q, v)$  can be any textual relevancy and for simplicity, in the paper we take the following  $tf \cdot idf$  based textual relevance as an example,

$$TXDist(q, v) = \sum_{w \in q.\mathcal{W}} tf(w, v) * idf(w, \mathcal{V}), \quad (2)$$

where  $tf(w, v)$  is the term frequency of term  $w$  in object  $v$  and  $idf(w, \mathcal{V})$  is the inverse document frequency of term  $w$  in the vertex set  $\mathcal{V}$ .

For example, consider a query  $q = \langle v_4, \{v_3, v_9, v_{15}\}, \{ 'park' \}, 2 \rangle$  in Fig. 1. Let us denote  $\alpha = 0.5$ ,  $M_S = 21$  and  $M_T = 0.5$ . Additionally, we have  $TXDist(q, v_3) = 0$ ,  $TXDist(q, v_9) = 0.25$  and  $TXDist(q, v_{15}) = 0.5$ . Thus, the overall ranking scores are respectively  $\mathcal{F}(q, v_3) = \frac{1}{2} \cdot \frac{6}{21} + \frac{1}{2} \cdot \frac{0}{0.5} = 0.45$ ,  $\mathcal{F}(q, v_9) = \frac{1}{2} \cdot \frac{19}{21} + \frac{1}{2} \cdot \frac{0.25}{0.5} = 0.43$  and  $\mathcal{F}(q, v_{15}) = \frac{1}{2} \cdot \frac{0}{21} + \frac{1}{2} \cdot \frac{0.5}{0.5} = 0.5$ . Therefore, the top-2 answer  $\mathcal{R} = \{v_{15}, v_3\}$ .

It is noteworthy that G-tree can support other functions (e.g., Jaccard and BM25) for textual relevancy. In this paper, we focus on effective index and efficient algorithms.

## 2.2 Related Works

In this section, we introduce the related works for different types of spatial queries.

*Single-pair shortest path queries.* Many previous studies [1], [5], [9], [10], [20], [21] addressed the SPSP problem on road networks. *TNR* [1] and *CH* [5] are the state-of-the-art methods. The idea of *TNR* is to figure out a set of transit nodes  $t_{v_x}$  for every vertex  $v_x$  on the road network, and pre-compute all pairs of distances between these transit nodes and the distance between a vertex to its transit nodes. Then, the shortest path between  $u$  and  $v$  is the minimal one of  $\text{SPDist}(u, t_u) + \text{SPDist}(t_u, t_v) + \text{SPDist}(t_v, v)$ . However, *TNR* incurs a large amount of memory, and its performance is poor for local queries (where the two vertices are close). *CH* first pre-computes a “contracted” version of the original road network by appending additional edges, and then uses bidirectional Dijkstra algorithm to search for the shortest paths. Though *CH* incurs a lower memory overhead, it has to visit a large number of vertices when two vertices are remote. Note that, though *TNR* is very efficient on SPSP queries, it is hard to be extended to support  $k$ NN queries efficiently. The key of  $k$ NN search is to prune those unpromising objects, however *TNR* fails to provide any additional information (e.g., boundaries) to support this.

Although *HEPV* [9] and *HiTi* [10] also organize the road network into a hierarchical structure, they suffered from tremendous storage overhead for maintaining large numbers of path information, thus they can only afford at most three hierarchical levels. Furthermore, they still use ‘half-blind’ Dijkstra-like network expansion algorithm, which is totally different from the dynamic programming scheme of G-tree based on the assembly-based method.

*K nearest neighbor queries on road network.* The authors in [13], [16], [19] addressed the  $k$ NN search problem on road networks. *INE* [16] extended the Dijkstra algorithm by expanding neighbor vertices from the query location until  $k$ NN answers have been found. *IER* [16] improved *INE* by utilizing spatial pruning techniques, e.g., taking the euclidean distance as a bound, to prune unpromising expansions. *IER* and *INE* are ‘blind’ algorithms since they neither capture the global distance from objects to the query location nor prune unnecessary objects efficiently.

*ROAD* [13] extends the Dijkstra algorithm by using a hierarchical structure. *ROAD* recursively partitions a road network into sub-networks, pre-computes the shortest-path distances of “shortcuts” within a sub-network, and organizes them in a hierarchical manner. By using Dijkstra-like network expansion, *ROAD* can skip sub-networks which do not contain an object. However it cannot prune sub-networks with objects widely scattered. For example, if the objects are widely distributed (e.g., gas stations), *ROAD* will degenerate to the Dijkstra algorithm and have to traverse the whole network. Thus *ROAD* performs poorly, especially on large networks. We call *ROAD* a ‘half-blind’ algorithm as it partially captures global distance information.

Although *ROAD* adopts the network partition to form a hierarchical structure, pre-computes shortest distances between borders and uses bottom-up method to derive those distances, the idea and technique of G-tree differ greatly as follows. First, the index structure is different. *ROAD* implements “Route Overlay(RO)” and “Association

Directory(AD)” to organize the pre-computed shortcuts and objects in a  $B^+$ -Tree. However, its query cost is high since each node expansion will traverse a tree-like entry under RO, especially when objects are widely scattered. G-tree employs light-weight distance matrices, and each distance matrix on G-tree node is accessed only once for a single query. Second, the algorithms of finding the  $k$ NN answers are fundamentally different. *ROAD* employs an expansion-based method and cannot utilize the global distance information, e.g., the shortest-path distance from a query location to tree nodes, to do effective pruning. G-tree adopts a best-first search algorithm which only accesses tree branches containing objects and thus reduces the space complexity significantly. Thus our method significantly outperforms *ROAD* (see Section 7).

*SILC* [19] pre-computes the shortest paths between all vertex pairs and uses a quadtree-based encoding to store the shortest paths. It utilizes the materialized pairs to find  $k$  nearest neighbors by using euclidean distance and stores the shortest-path distance as a bound. However if there are large numbers of objects clustered in a small region, *SILC* is inefficient. Moreover, *SILC* consumes  $\mathcal{O}(|\mathcal{V}|^{1.5})$  storage space and incurs high pre-processing overhead, and thus it is impractical for large road networks.

There are some studies which assume that the object set is given [2], [6], [7], [12]. They pre-compute and materialize results of potential queries. However these approaches highly depend on the given object set. They also involve high pre-computation cost and large memory overhead.

This paper is a major-value added version of our previous work [23]. The new contributions include (1) We considered three types of spatial queries under the framework of G-tree and added how to support single-pair shortest path query in Section 4.1 and keyword-based  $k$ NN query in Section 4.3. We improved the  $k$ NN search algorithm in Section 4.2; (2) We conducted new experiments to evaluate G-tree in Section 7.3.1 and new experiments to evaluate keyword-based  $k$ NN query in Section 7.3.3; (3) We discussed the path recovery issue in Section 5.

*Keyword-based kNN query on road network.* The authors in [14], [17], [18] studied the problem of keyword-based  $k$ NN query on road networks. However these works are all different from ours, both in the problem definition and the query processing techniques. The authors in [17] focuses on keyword-based  $k$ NN search without combining textual relevancy and spatial distance. [14] studies keyword-based  $k$ NN search in the distributed environment. Though [18] considers both textual and spatial relevancy between the query location and the object, it assumes that  $\mathcal{C} = \mathcal{V}$ , thus the availability is low. Moreover, [18] supposes objects are on edges, which is different from our assumption that objects are on vertices. Furthermore, [18] still employs a Dijkstra-like expansion algorithm to locate objects, which is rather inefficient for objects that are far away from the query location.

IR-Tree [3] and [4] studied the problem of finding top- $k$  spatial objects on metric space. Though the textual inverted list of G-tree is inspired by the idea of pseudo documents in IR-Tree, the motivations are different. IR-Tree focused on efficient solutions to handle textual-spatial queries on metric space, while our goal is to devise a tree index on road



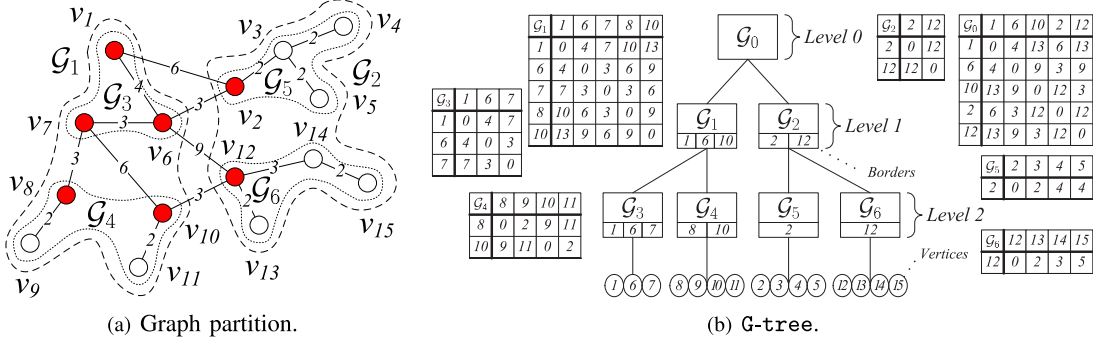


Fig. 2. A G-tree (The solid vertices in the graph are borders.  $f = 2$  and  $\tau = 4$ ).

network which provides efficiency, scalability and good extensibility for various types of queries.

*kNN queries for moving objects on road network.* There are a number of studies on  $kNN$  queries for moving objects monitoring [8], [15], [22], [24]. These works studied the problem of finding nearest moving objects (e.g., taxis) to a location and focused on dealing with frequent updates of moving objects. In our case, we emphasize on the efficiency of the  $kNN$  queries on static objects (e.g., gas stations), and the ideas and techniques are different.

### 3 THE G-tree INDEX

In this section, we first introduce the basic idea and the overview of our index in Section 3.1. Then, we formally define the G-tree in Section 3.2 and present how to construct the G-tree in Section 3.3. Finally, we discuss the space complexity of the G-tree in Section 3.4.

#### 3.1 Basic Idea and Overview

R-tree is a height-balanced index structure that has been widely adopted for spatial data because of its efficiency and scalability. On one hand, it can prune unpromising subtrees at query time to reduce the search space. On the other hand, it can adapt to different types of queries, e.g.,  $kNN$  query and spatial keyword query. Inspired by its salient features, our goal is to devise a similar index on road networks. To this end, we recursively partition the road network into sub-networks and construct a tree-structured index on top of the sub-networks where each G-tree node corresponds to a sub-network. Fig. 2 shows an example of G-tree. (We will formally define the G-tree in Section 3.2).

In addition, it is easy to calculate the euclidean distance between two vertices (or a vertex and a R-tree node) in  $\mathcal{O}(1)$  time and space complexity. However on road networks, it is challenging to compute the shortest-path distance between two vertices (or a vertex and a G-tree node). If we pre-compute all shortest-path distances between two vertices and materialize them, it will take  $\mathcal{O}(|V|^2)$  space storage and  $\mathcal{O}(|V|^3)$  pre-processing time cost, and the query time cost is  $\mathcal{O}(1)$ . If we calculate the distance online, the space complexity is  $\mathcal{O}(|V|)$  and the time complexity is  $\mathcal{O}(|V|\log|\mathcal{E}|)$ . Obviously, these two extreme paradigms are neither efficient nor scalable. Thus it calls for an efficient method to compute shortest-path distances with high efficiency and good scalability (i.e., low storage cost). To address this issue, we devise an *assembly-based method* to

efficiently compute the shortest-path distance, which only pre-computes and materializes essential pairs of shortest-path distances on G-tree (Fig. 2b), so that any shortest-path distance can be assembled piece by piece with this materialization. Furthermore, it costs only  $\mathcal{O}(|V|\log|V|)$  space overhead, and  $\mathcal{O}(|V|)$  time complexity. Therefore, G-tree can be scalable to very large datasets. We formally define the G-tree in the next section, and present the details of the *assembly-based method* in Section 4.1.

#### 3.2 Definition of G-tree

We first introduce some important concepts which will be used throughout the paper.

**Definition 1 (Graph Partition).** Given a graph  $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$ , where  $\mathcal{V}$  is the vertex set and  $\mathcal{E}$  is the edge set of  $\mathcal{G}$ , a partition of  $\mathcal{G}$  is a set of subgraphs, i.e.,  $\mathcal{G}_1 = \langle \mathcal{V}_1, \mathcal{E}_1 \rangle, \mathcal{G}_2 = \langle \mathcal{V}_2, \mathcal{E}_2 \rangle, \dots, \mathcal{G}_f = \langle \mathcal{V}_f, \mathcal{E}_f \rangle$  such that

- 1)  $\bigcup_{1 \leq i \leq f} \mathcal{V}_i = \mathcal{V}$ ,
- 2) For  $i \neq j$ ,  $\mathcal{V}_i \cap \mathcal{V}_j = \emptyset$ , and
- 3)  $\forall u, v \in \mathcal{V}_i$ , if  $(u, v) \in \mathcal{E}$ , then  $(u, v) \in \mathcal{E}_i$ .

To differentiate those *portal* vertices in each subgraph, we define a concept, called *borders*.

**Definition 2 (Borders).** Given a subgraph  $\mathcal{G}_i$  of  $\mathcal{G}$ , a vertex  $u \in \mathcal{V}_i$  is called a border if  $\exists (u, v) \in \mathcal{E}$  and  $v \notin \mathcal{V}_i$ . We use  $\mathcal{B}(\mathcal{G}_i)$  to denote the border set in graph  $\mathcal{G}_i$ .

A subgraph  $\mathcal{G}_i$  is called a super-graph of another subgraph  $\mathcal{G}_j$  if  $\mathcal{V}_i \supseteq \mathcal{V}_j$  and  $\mathcal{E}_i \supseteq \mathcal{E}_j$ .

Now, we formally define the G-tree as follows.

**Definition 3 (G-tree).** A G-tree is a balanced search tree that satisfies the following properties.

- 1) Each node represents a subgraph. The root node corresponds to the graph  $\mathcal{G}$ . The subgraph of a parent node is a super-graph of those of its child nodes.
- 2) Each non-leaf node has  $f(\geq 2)$  children.
- 3) Each leaf node contains at most  $\tau(\geq 1)$  vertices. All leaf nodes appear at the same level.
- 4) Each node maintains its border set and a distance matrix. In the distance matrix of a non-leaf node, the columns/rows are all borders in its children and the value of each entry is the shortest-path distance between the two borders. In the distance matrix of a leaf node, the rows are all borders in the node, columns

are all vertices in this node, and the value of each entry is the shortest-path distance between the border and the vertex.

- 5) For keyword queries, each node also stores additional information, e.g. inverted lists, boundaries, etc.

Properties (1)-(3) ensure that the G-tree has a balanced search tree structure. It is worth noting that for each node we do not maintain the physical subgraph. Instead, we only maintain a dummy subgraph ID. To support local Dijkstra and path recovery on leaf nodes, we store the original graph  $\mathcal{G} = \langle \mathcal{V}, \mathcal{E} \rangle$  individually. As there is a one-to-one correspondence between a node and a subgraph, for simplicity, “nodes” and “subgraphs” are interchangeably used if the context is clear. In the paper “nodes” refer to G-tree nodes and “vertices” refer to vertices in the graph

Property (4) is used to compute the shortest-path distances based on the assembly-based method, which will be discussed in Section 4.1. It is worth noting that these matrices are not actually stored on G-tree nodes but maintained on the continuous memory/disk storage space and only an offset pointer is maintained on each node.

Property (5) makes G-tree adaptive for various queries. Boundaries can be added to do pruning at query time, and keyword inverted list can be employed to support keyword search. We will discuss the details in Section 4.

**Example 1.** Fig. 2b shows the G-tree of the road network in Fig. 2a. The borders of each node are shown in the rectangle box under the node. The distance matrix of each node is listed around the tree node. For  $\mathcal{G}_1$ , its children  $\mathcal{G}_3$  and  $\mathcal{G}_4$  contain five borders  $\{v_1, v_6, v_7, v_8, v_{10}\}$ , thus the rows/columns of  $\mathcal{G}_1$ 's distance matrix are the five borders. The set of vertices of each leaf node are shown in the circled numbers. For instance, in  $\mathcal{G}_4$ 's distance matrix, the rows are borders  $\{v_8, v_{10}\}$  and the columns are vertices  $\{v_8, v_9, v_{10}, v_{11}\}$ . The entry  $(v_8, v_{11}) = 11$  since the shortest distance between border  $v_8$  and  $v_{11}$  is 11.

### 3.3 G-tree Construction

In this section, we present how to construct the G-tree. First, we discuss how to establish a tree structure from the road network with the help of graph partition algorithm. Then, we introduce how to calculate the distance matrices.

#### 3.3.1 Tree Construction

We use a graph partition based method to build the G-tree. Initially, we take the graph  $\mathcal{G}$  as the root. Then we partition  $\mathcal{G}$  into  $f$  equal-sized subgraphs (i.e.,  $|\mathcal{V}_1| \approx \dots \approx |\mathcal{V}_f|$ ) and take them as the root's children. Next we recursively partition the children and repeat this step until each leaf-node's subgraph has no more than  $\tau$  vertices. During the partitioning, for each subgraph, we will add its borders into the corresponding node. For example, in Fig. 2a, suppose  $f = 2$  and  $\tau = 4$ , the graph  $\mathcal{G}_0$  is partitioned into two subgraphs  $\mathcal{G}_1$  and  $\mathcal{G}_2$ .  $\mathcal{G}_1$  is further partitioned into  $\mathcal{G}_3$  and  $\mathcal{G}_4$ .  $\mathcal{G}_2$  is partitioned into  $\mathcal{G}_5$  and  $\mathcal{G}_6$ .

Graph partitioning is an important step in G-tree construction. The optimal one should not only generate

approximately equal-sized subgraphs, but also minimize the number of borders. However, it has been proven that optimal graph partitioning is NP-Hard. In this paper, we adopt a famous heuristics algorithm, called the multilevel partitioning algorithm [11]. It first reduces the graph size by coarsening the vertices and edges, and then partitions on the coarsened graph using traditional graph partitioning algorithms, e.g. Kernighan-Lin algorithm. Finally it uncoarsens the subgraphs to generate partitions of the original graph. Moreover, the multilevel partition algorithm can guarantee that each subgraph nearly has the same size and thus G-tree is a balanced search tree structure.

#### 3.3.2 Distance Matrices

For distance matrices of G-tree, we need to compute the shortest-path distance between a border and a border (non-leaf node)/vertex(leaf node). We can exploit a single source shortest-path algorithm, e.g., Dijkstra algorithm, which starts from each border/vertex within one G-tree node, and expands the edges until all borders of such a node have been reached. In Section 6.1, we will introduce an efficient bottom-up algorithm to speed up this procedure.

### 3.4 Space Complexity of the G-tree

**Height.** The height of the G-tree is  $\mathcal{H} = \log_f \frac{|\mathcal{V}|}{\tau} + 1$ .

**Number of nodes.** At level 0, there is one node (the root). In level  $i$ , there are  $f^i$  nodes. There are  $\frac{|\mathcal{V}|}{\tau}$  leaf nodes. Thus the total number of nodes is  $\mathcal{O}(\frac{f-1}{f-1} \frac{|\mathcal{V}|}{\tau}) = \mathcal{O}(\frac{|\mathcal{V}|}{\tau})$ .

**Number of borders.** A road network is usually modeled as a planar graph [19]. We also consider the planar graph in the space analysis. Consider a node on the  $i$ th level. Its borders are generated by its parent which has  $|\mathcal{V}|/f^{i-1}$  vertices. According to the Planar Separator Theorem, the  $f$ -partition on the parent totally contains  $\mathcal{O}(\log_2 f \cdot \sqrt{|\mathcal{V}|/f^{i-1}})$  borders. The  $f$  children of the parent share these borders. Thus each node at level  $i$  has  $\mathcal{O}(\log_2 f \cdot \sqrt{|\mathcal{V}|/f^{i+1}})$  borders on average. As there are  $f^i$  nodes, the number of borders in level  $i$  is  $\mathcal{O}(\log_2 f \cdot \sqrt{|\mathcal{V}|f^{i-1}})$ . The total number of borders in the G-tree is  $\mathcal{O}(\sum_{1 \leq i \leq \mathcal{H}} \log_2 f \cdot \sqrt{|\mathcal{V}|f^{i-1}}) = \mathcal{O}(\frac{\log_2 f}{\sqrt{f}} |\mathcal{V}|)$ .

**Distance matrices.** The average number of borders in a leaf node is  $\mathcal{O}(\log_2 f \cdot \sqrt{|\mathcal{V}|/f^{\mathcal{H}+1}}) = \mathcal{O}(\log_2 f \cdot \sqrt{\tau})$ . The number of vertices in a leaf node is  $\tau$ . Thus the distance-matrix size of a leaf node is  $\mathcal{O}(\log_2 f \cdot \tau^{1.5})$ . The total distance-matrix size of all leaf nodes is  $\mathcal{O}(\log_2 f \cdot \tau^{1.5} \frac{|\mathcal{V}|}{\tau}) = \mathcal{O}(\log_2 f \cdot \sqrt{\tau} |\mathcal{V}|)$ . For each non-leaf node, the rows/columns of its distance matrix are the union of the borders in its children. Each node on level  $i$  generates  $\mathcal{O}(\log_2 f \cdot \sqrt{|\mathcal{V}|/f^i})$  borders. Thus the matrix size of each node at level  $i$  is  $\mathcal{O}(\log_2^2 f \cdot |\mathcal{V}|/f^i)$ . As there are  $f^i$  nodes at level  $i$ , the distance-matrix size at level  $i$  is  $\mathcal{O}(\log_2^2 f \cdot |\mathcal{V}|)$ . Hence the total matrix size of non-leaf nodes is  $\mathcal{O}(\mathcal{H} \log_2^2 f \cdot |\mathcal{V}|) = \mathcal{O}(\log_2^2 f \cdot \log_f \frac{|\mathcal{V}|}{\tau} \cdot |\mathcal{V}|)$ .

**Road networks.** A hash table  $\text{leaf}(v)$  is maintained to map a vertex to the corresponding leaf node with cost  $\mathcal{O}(|\mathcal{V}|)$ . We also store the original graph for local Dijkstra search and path recovery (Section 5) with cost  $\mathcal{O}(|\mathcal{E}| + |\mathcal{V}|)$ .

**Overall space.** The overall size of the G-tree is  $\mathcal{O}(\frac{|\mathcal{V}|}{\tau} + \frac{\log_2 f}{\sqrt{\tau}} |\mathcal{V}| + \log_2 f \cdot \sqrt{\tau} |\mathcal{V}| + \log_2^2 f \cdot \log_f \frac{|\mathcal{V}|}{\tau} \cdot |\mathcal{V}| + 2|\mathcal{V}| + |\mathcal{E}|) = \mathcal{O}(\log_2 f \cdot \sqrt{\tau} |\mathcal{V}| + \log_2^2 f \cdot \log_f \frac{|\mathcal{V}|}{\tau} \cdot |\mathcal{V}| + |\mathcal{E}|)$ . As  $\log_2^2 f, \sqrt{\tau}$  and  $\log_f \frac{|\mathcal{V}|}{\tau}$  are small, G-tree is scalable.

## 4 SEARCH ALGORITHM

In this section, we first present the search algorithms for three types of queries: SPSP,  $k$ NN and  $k^2N^2$  in Sections 4.1, 4.2, and 4.3 respectively. Then we analyze the time and space complexity in Section 4.4.

### 4.1 Single-Pair Shortest Path Query

Given a query  $q = \langle u, v \rangle$ , the single-pair shortest-path query returns  $\text{SPDist}(u, v)$ . However, calculating shortest-path distance on road networks is either a time-consuming or space-consuming task. To address this problem, we propose an *assembly-based method* which pre-computes and materializes essential pairs of graph distances on G-tree, so that any shortest-path distance can be assembled piece by piece using the materialized pairs.

Formally we consider two general cases: (1)  $u$  and  $v$  are in different leaf nodes of the G-tree, and we use function  $\text{MINDIST-OUTSIDE-LEAF}$ ; (2)  $u$  and  $v$  are in the same leaf node, and we use function  $\text{MINDIST-INSIDE-LEAF}$ .

**MINDIST-OUTSIDE-LEAF.** For ease of presentation, we first introduce a closure-based method.

**Closure-based method.** We introduce the notion of closure as formalized in Lemma 1, and

**Lemma 1 (Closure).** *Given a subgraph  $G_i = \langle \mathcal{V}_i, \mathcal{E}_i \rangle$ , for any vertex  $u \notin \mathcal{V}_i$  and  $v \in \mathcal{V}_i$ , any shortest path between  $u$  and  $v$  must contain a border in  $\mathcal{B}(G_i)$ , i.e., for any shortest path  $\text{SP}(u, v)$ ,  $\exists w \in \mathcal{B}(G_i), w \in \text{SP}(u, v)$ .*

**Proof.** We prove it by contradiction. Let  $\text{SP}(u, v) = uv_1v_2 \cdots v_xv$ . Suppose  $v_1, v_2, \dots, v_x, v$  are not borders in  $\mathcal{B}(G')$ . If  $v_x \notin G'$ , then  $v$  is a border based on the definition which contradicts with the assumption. If  $v_x \in G'$ , we find the vertex with the largest subscript which is in  $G'$ , e.g.,  $v_i$ . Then  $v_{i-1} \notin G' (v_0 = u)$ . Based on the definition of border,  $v_i$  is a border which contradicts with the assumption. Thus the lemma is proved.  $\square$

As shown in Fig. 2a, consider  $v_9 \in G_1$  and  $v_4 \in G_5$ . As  $v_4$  and  $v_9$  are not within the same subgraph, any path from  $v_4$  to  $v_9$  must contain a border in  $G_1$ , e.g.,  $v_6$ . Similarly, any path must contain a border in  $G_4$ , e.g.,  $v_8$ .

Based on Lemma 1, we know that any path from  $u$  to  $v$  can be decomposed into two parts: the first part is from  $u$  to the border of a G-tree leaf node which  $v$  belongs to; and the second part is from such border to  $v$ . Thus,

$$\text{SPDist}(u, v) = \min_{b \in \mathcal{B}(\text{leaf}(v))} (\text{SPDist}(u, b) + \text{SPDist}(b, v)), \quad (3)$$

where  $\text{leaf}(v)$  denotes the leaf node that contains  $v$ . Basically,  $\text{leaf}(v)$  is implemented by a hash table which maps a vertex to the corresponding leaf node. Equation (3) suggests that we only need to materialize the distances between vertices and borders, and utilize these distances to compute any  $\text{SPDist}(u, v)$ . As the number of borders is less than the

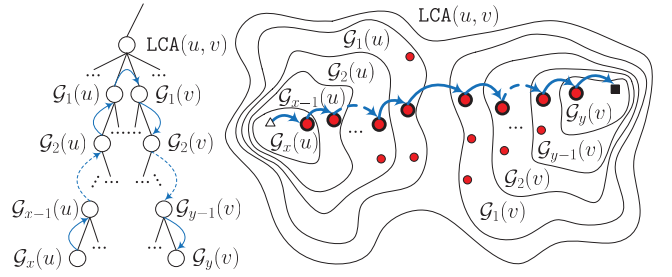


Fig. 3. A sketch of the assembly-based method.

total number of vertices, this method seems to work. However, in the worst case, the total number of borders is still  $\mathcal{O}(|\mathcal{V}|)$  (see Section 3.4). To address this issue, we propose an assembly-based method.

**Assembly-based method.** To overcome the problem of high space overhead in the closure-based method, we propose a novel way which only uses  $\mathcal{O}(|\mathcal{V}| \log |\mathcal{V}|)$  space overhead and  $\mathcal{O}(|\mathcal{V}|)$  time complexity.

The assembly-based method is the generalized version of the closure-based method. Rather than decomposing one path from  $u$  to  $v$  into two sub-paths by a single border, i.e.  $(u, b)$  and  $(b, v)$ , it decomposes one path into  $m + 1$  sub-paths by  $m$  borders, i.e.  $(u, b_1)(b_1, b_2), \dots, (b_m, v)$ , while these  $m$  borders belong to  $m$  distinct and consecutive G-tree nodes on the path from  $\text{leaf}(u)$  to  $\text{leaf}(v)$ . Fig. 3 illustrates this procedure.

Formally, let  $\text{LCA}(u, v)$  denote the least common ancestor of nodes  $\text{leaf}(u)$  and  $\text{leaf}(v)$ . Based on Lemma 1, any path from  $u$  to  $v$  must bypass a series of G-tree nodes,  $G_x(u) = \text{leaf}(u), G_{x-1}(u), \dots, G_1(u), G_1(v), G_2(v), \dots, G_y(v) = \text{leaf}(v)$ . To compute the shortest-path distance from  $u$  to  $v$ , we consider the following cases.

- 1) The shortest-path distance from  $u$  to  $G_x(u)$ . In this case,

$$\text{SPDist}(u, G_x(u)) = \min_{u_i \in \mathcal{B}(G_x(u))} (\text{SPDist}(u, u_i)), \quad (4)$$

where  $\text{SPDist}(u, u_i)$  is materialized in the distance matrix of the leaf node  $G_x(u)$ .

- 2) The shortest-path distance from  $u$  to  $G_{i-1}(u)$ . As the shortest path from  $u$  to  $G_{i-1}(u)$  must contain a border in  $G_i(u)$  based on the closure property, we have,

$$\text{SPDist}(u, G_{i-1}(u)) = \min_{u_i \in \mathcal{B}(G_i(u))} (\text{SPDist}(u, u_i) + \text{SPDist}(u_i, G_{i-1}(u))), \quad (5)$$

where  $\text{SPDist}(u, u_i)$  can be computed iteratively.

- 3) The shortest-path distance from  $u$  to  $G_1(v)$ . As the shortest path from  $u$  to  $G_1(v)$  must contain a border in  $G_1(u)$ , we have,

$$\text{SPDist}(u, G_1(v)) = \min_{u_1 \in \mathcal{B}(G_1(u))} (\text{SPDist}(u, u_1) + \text{SPDist}(u_1, v_1)). \quad (6)$$



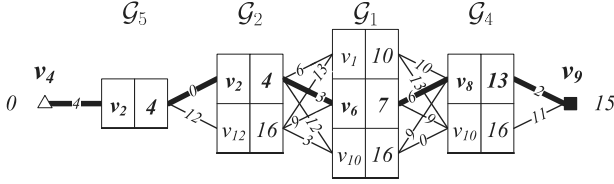


Fig. 4. Dynamic programming of MinDist-Outside-Leaf.

Based on Equations (4), (5) and (6), we have

$$\begin{aligned}
 \text{SPDist}(u, v) = & \min_{u_x \in \mathcal{B}(\mathcal{G}_x(u))} (\text{SPDist}(u, u_x) \\
 & + \min_{u_{x-1} \in \mathcal{B}(\mathcal{G}_{x-1}(u))} (\text{SPDist}(u_x, u_{x-1}) \\
 & + \cdots + \min_{u_1 \in \mathcal{B}(\mathcal{G}_1(u))} (\text{SPDist}(u_2, u_1) \\
 & + \min_{v_1 \in \mathcal{B}(\mathcal{G}_1(v))} (\text{SPDist}(u_1, v_1) \\
 & + \min_{v_2 \in \mathcal{B}(\mathcal{G}_2(v))} (\text{SPDist}(v_1, v_2) \\
 & + \cdots + \min_{v_x \in \mathcal{B}(\mathcal{G}_x(v))} (\text{SPDist}(v_x, v)). \quad (7)
 \end{aligned}$$

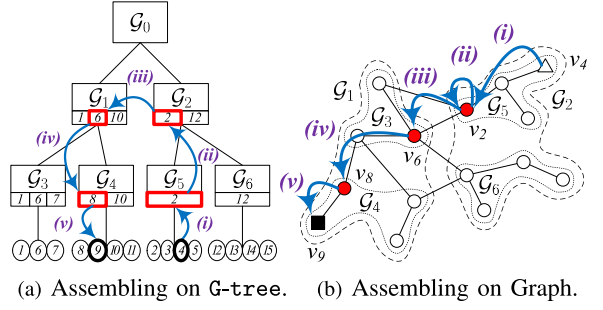
Equation (7) indicates that we can use a *dynamic-programming algorithm* to efficiently calculate  $\text{SPDist}(u, v)$ . Based on the distance matrices on G-tree, we first compute  $\text{SPDist}(u, u_x)$ ,  $u_x \in \mathcal{B}(\mathcal{G}_x(u))$ . Next, we move forward to the next level to compute  $\text{SPDist}(u, u_{x-1})$  for  $u_{x-1} \in \mathcal{B}(\mathcal{G}_{x-1}(u))$ , and iteratively can reach  $\mathcal{G}_1(u)$ . Then we cross from  $\mathcal{G}_1(u)$  to  $\mathcal{G}_1(v)$ , and move to the other branch. Iteratively, we can finally compute  $\text{SPDist}(u, v)$ .

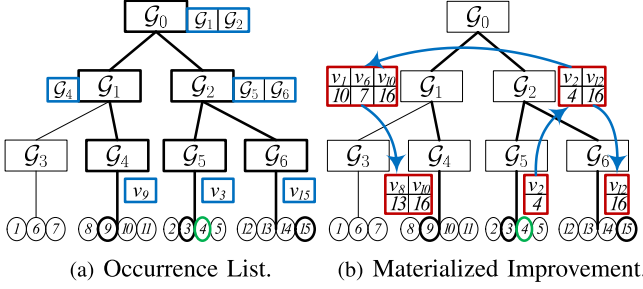
To summarize, given two vertices  $u$  and  $v$  where  $u \notin \text{leaf}(v)$ , to compute  $\text{SPDist}(u, v)$ , we first compute their least common ancestor and the nodes on the paths from  $\text{LCA}(u, v)$  to  $\text{leaf}(u)$  and  $\text{leaf}(v)$ . Then we use dynamic programming to compute  $\text{SPDist}(u, v)$ . Figs. 4 and 5 show how to implement MinDist-Outside-Leaf.

We had proved that this paradigm takes only  $\mathcal{O}(|\mathcal{V}| \log |\mathcal{V}|)$  space overhead in Section 3.4, and we will show its computation cost is  $\mathcal{O}(|\mathcal{V}|)$  in Section 4.4. Compared with previous work [19] (which takes  $\mathcal{O}(|\mathcal{V}|^{1.5})$  space), our method significantly reduces the time/space overhead and can scale up to very large datasets. Furthermore, we will show the advantages of our assembly-based method to support other queries, e.g.,  $k$ NN queries.

**MinDist-Inside-Leaf.** Given two vertices  $u$  and  $v$ , and  $u, v \in \text{leaf}(u)$ , consider a shortest path  $\text{SP}(u, v)$ . There are two cases: (1)  $\text{SP}(u, v)$  does not contain a vertex outside node  $\text{leaf}(u)$ . In this case, we use the Dijkstra algorithm to compute the shortest path in node  $\text{leaf}(u)$ . Let  $\text{DijkDist}(u, v)$  denote the distance inside  $\text{leaf}(u)$ . Since the subgraph w.r.t. the leaf node is not large, the Dijkstra algorithm is efficient enough. (2)  $\text{SP}(u, v)$  contains a vertex outside node  $\text{leaf}(u)$ . In this case,  $\text{SP}(u, v)$  must contain two borders  $b_1, b_2$  in  $\text{leaf}(u)$ . Let  $\text{BorderDist}(u, v)$  denote the shortest distance from  $u$  to  $v$  with outside vertices,

$$\begin{aligned}
 \text{BorderDist}(u, v) = & \min_{b_1, b_2 \in \mathcal{B}(\text{leaf}(u))} (\text{SPDist}(u, b_1) \\
 & + \text{SPDist}(b_1, b_2) + \text{SPDist}(b_2, v)). \quad (8)
 \end{aligned}$$



Fig. 6. Occurrence list and  $k$ NN Search.

$k$ NN query, therefore, such materialization-based method can avoid the duplicated computation. In Section 4.4, we will show it only takes  $O(\frac{\log_2 f}{\sqrt{\tau}} |\mathcal{V}|)$  space complexity.

The second issue is how to efficiently obtain the top- $k$  answers. Intuitively, the nearer those subgraphs (nodes) are to  $v_q$ , the greater likelihood that they may contain the top- $k$  objects. To this end we utilize a priority-queue to manage those traversed G-tree nodes. In this way, we can always dequeue the most promising (nearest) nodes each time. So the problem is to calculate the minimal distance between  $v_q$  and a G-tree node  $n$ , denoted by  $\text{SPDist}(v_q, n)$ . Based on the assembly-based method, we have

$$\text{SPDist}(v_q, n) = \min_{b \in B(n)} \text{SPDist}(v_q, b). \quad (10)$$

As  $\text{SPDist}(v_q, b)$  is calculated and materialized progressively at query time, hence, we only need  $O(1)$  time to figure out  $\text{SPDist}(v_q, n)$ . This implies that, during the process of the step-by-step node expansion on the G-tree, we can exploit the intermediate materialization to obtain  $\text{SPDist}(v_q, n)$ , which is critical and indispensable for top- $k$  finding and effective pruning. Thus G-tree is efficient enough to deal with very large datasets for  $k$ NN queries.

Notice that, based on the candidate object set  $\mathcal{C}$  specified by the user, we can filter out those unpromising G-tree nodes before enqueued into the priority-queue. To this end, we introduce the concept of *occurrence list*, i.e.,  $\mathcal{L}(n)$ . The occurrence list for a leaf node is the list of objects in this node; the occurrence list of a non-leaf node is a list of its children who contain objects. Fig. 6a illustrates an example for  $\mathcal{C} = \{v_3, v_9, v_{15}\}$ . For example, for  $G_5$ , its occurrence list is  $\{v_3\}$  as vertex  $v_3$  is an object. For  $G_2$ , its occurrence list is  $\{G_5, G_6\}$  as nodes  $G_5$  and  $G_6$  contain objects in  $\mathcal{C}$ . The benefit of the occurrence list is to decouple G-tree's structure with object information, thus G-tree can effectively support various object sets.

Algorithm 2 shows the pseudo-code of our algorithm, which works as follows. First of all, we locate the leaf node of query location  $v_q$  and objects  $\mathcal{C}$  by  $\text{leaf}(v)$ , and construct the occurrence lists in a bottom-up manner (line 1). Then we create a priority queue  $\mathcal{Q}$  and a result set  $\mathcal{R}$  (line 2). Initially we use MINDIST-INSIDE-LEAF as discussed in Section 4.1 to calculate  $\text{SPDist}(v_q, v \in \mathcal{L}(\text{leaf}(v_q)))$ , and put  $\langle v, \text{SPDist}(v_q, v) \rangle$  into  $\mathcal{Q}$  (line 4). Next, we iteratively dequeue the first element  $e$  of  $\mathcal{Q}$  and handle it separately according to whether  $e$  is an object, or a G-tree node (line 6 to line 15).

### Algorithm 2. KNNSEARCH ( $q = \langle v_q, \mathcal{C}, k \rangle, \mathcal{G}$ )

**Input:**  $q = \langle v_q, \mathcal{C}, k \rangle$ : A query;  $\mathcal{G}$ : A G-tree

**Output:**  $\mathcal{R}$ : The top- $k$  result list;

```

1 Compute the occurrence list  $\mathcal{L}$  based on  $\mathcal{C}$ ;
2 Initialize priority queue  $\mathcal{Q} = \emptyset$  and result set  $\mathcal{R} = \emptyset$ ;
3 foreach  $v \in \mathcal{L}(\text{leaf}(v_q))$  do
4    $\mathcal{Q}.$  Enqueue( $\langle v, \text{MINDIST} - \text{INSIDE} - \text{LEAF}(v_q, v) \rangle$ );
5 Initialize  $\mathcal{T}_n = \text{leaf}(v_q)$  and  $\mathcal{T}_{\min}$ ;
6 while  $\mathcal{R}.$  Size()  $< k$  and ( $\mathcal{Q} \neq \emptyset$  or  $\mathcal{T}_n \neq \text{root}$ ) do
7   if  $\mathcal{Q} = \emptyset$  then UpdateT( $\mathcal{T}_n, \mathcal{T}_{\min}, \mathcal{Q}$ );
8    $\langle e, dis \rangle \leftarrow \mathcal{Q}.$  Dequeue();
9   if  $dis > \mathcal{T}_{\min}$  then
10    UpdateT( $\mathcal{T}_n, \mathcal{T}_{\min}, \mathcal{Q}$ );
11     $\mathcal{Q}.$  Enqueue( $\langle e, dis \rangle$ );
12   else if  $e$  is an object then Insert  $e$  into  $\mathcal{R}$ ;
13   else if  $e$  is a node then
14     foreach  $c \in \mathcal{L}(e)$  do
15        $\mathcal{Q}.$  Enqueue( $\langle c, \text{SPDist}(v_q, c) \rangle$ );
16 Function UpdateT( $\mathcal{T}_n, \mathcal{T}_{\min}, \mathcal{Q}$ )
17    $\mathcal{T}_n \leftarrow \mathcal{T}_n.$  father and update  $\mathcal{T}_{\min}$ ;
18   foreach  $c \in \mathcal{L}(\mathcal{T}_n)$  do
19      $\mathcal{Q}.$  Enqueue( $\langle \text{SPDist}(v_q, c) \rangle$ );
```

It is worth noting that, for the purpose of effective pruning, we set up a global pointer  $\mathcal{T}_n$ , which keeps track of the current uppermost node ever visited, together with the minimum graph distance  $\mathcal{T}_{\min}$  within  $\mathcal{T}_n$  from  $v_q$ , i.e.,  $\mathcal{T}_{\min} = \min_{b \in B(\mathcal{T}_n)} \text{SPDist}(v_q, b)$  (line 3). Initially,  $\mathcal{T}_n$  is set to  $\text{leaf}(v_q)$ .  $\mathcal{T}_{\min}$  signifies the minimum graph distance to step off the current  $\mathcal{T}_n$ . Hence, once we dequeue an object  $e$  from  $\mathcal{Q}$  whose  $\text{SPDist}(v_q, e) \leq \mathcal{T}_{\min}$ , we can confirm that  $e$  is a top- $k$  answer (line 12). Otherwise, we should move the  $\mathcal{T}_n$  to its parent node, expand its sibling nodes into  $\mathcal{Q}$ , and update the  $\mathcal{T}_{\min}$  (line 10).

**Example 3.** Consider a query  $q = \{v_4, \{v_3, v_9, v_{15}\}, 2\}$  on the graph in Fig. 1a. We first construct the occurrence list and then compute the top-2 answers as follows.

**Step 1.** Enqueue the object in  $\mathcal{L}(G_5)$ , i.e.  $\langle v_3 \in \mathcal{L}(G_5), \text{SPDist}(v_4, v_3) = 2 \rangle$ . Initially,  $\mathcal{T}_n = G_5, \mathcal{T}_{\min} = 4$ .

Queue:  $\langle v_3, 2 \rangle$

**Step 2.** Dequeue  $\langle v_3, 2 \rangle$ . Since  $\text{SPDist}(v_4, v_3) = 2 < 4 = \mathcal{T}_{\min}$ , thus,  $\mathcal{R} = \{v_3\}$ . Now,  $\mathcal{Q} = \emptyset$ .

**Step 3.** As  $\mathcal{Q} = \emptyset$  but  $\mathcal{T}_n \neq G_0$ , update  $\mathcal{T}_n = G_2$  and  $\mathcal{T}_{\min} = 4$ . Enqueue  $\langle G_6 \in \mathcal{L}(G_2), \text{SPDist}(v_4, G_6) = 16 \rangle$ .

Queue:  $\langle G_6, 16 \rangle$

**Step 4.** Dequeue  $\langle G_6, 16 \rangle$ . Since  $\text{SPDist}(v_4, G_6) = 16 > 4 = \mathcal{T}_{\min}$ , update  $\mathcal{T}_n = G_0$  and  $\mathcal{T}_{\min} = \infty$ . Enqueue  $\langle G_1 \in \mathcal{L}(G_0), \text{SPDist}(v_4, G_1) = 7 \rangle$ . Enqueue  $\langle G_6, 16 \rangle$  back.

Queue:  $\langle G_1, 7 \rangle \mid \langle G_6, 16 \rangle$

**Step 5.** Dequeue  $\langle G_1, 7 \rangle$ . Find  $G_4 \in \mathcal{L}(G_1)$ . Get  $\text{SPDist}(v_4, G_4) = 13$  and enqueue  $\langle G_4, 13 \rangle$ .

Queue:  $\langle G_4, 13 \rangle \mid \langle G_6, 16 \rangle$

**Step 6.** Dequeue  $\langle G_4, 13 \rangle$ . Locate  $v_9 \in \mathcal{L}(G_4)$ . Calculate  $\text{SPDist}(v_4, v_9) = 15$  and enqueue  $\langle v_9, 15 \rangle$ .

Queue:  $\langle v_9, 15 \rangle \mid \langle G_6, 16 \rangle$

**Step 7.** Dequeue  $\langle v_9, 15 \rangle$ . Since  $\text{SPDist}(v_4, v_9) = 15 < \infty = \mathcal{T}_{\min}$ ,  $\mathcal{R} = \{v_3, v_9\}$ . Top-2 answers have been generated. Algorithm terminates.



### 4.3 Keyword-Based $k$ NN Query

Given a query  $q = \langle v_q, \mathcal{C}, \mathcal{W}, k \rangle$  and the ranking function  $\mathcal{F}$ , the  $k^2N^2$  query returns top- $k$  objects in  $\mathcal{C}$  based on the ranking function  $\mathcal{F}$ . Compared with  $k$ NN search,  $k^2N^2$  is more challenging since it considers not only spatial distance but also textual relevancy for ranking. In this section, we extend the  $k$ NN algorithm to efficiently answer  $k^2N^2$  queries. From previous sections, the most important component of  $k$ NN search is the distance calculation between the query location and a G-tree node, i.e.,  $\text{SPDist}(v_q, n)$  in Equation (10), which is used as the distance boundary for effective pruning. Similarly, we can define the textual boundary between a query  $q$  and a node  $n$ , denoted by  $\text{TXDist}(q, n)$ , and then aggregate the spatial and textual boundary to obtain the overall boundary  $\mathcal{F}(q, n)$  for pruning spatially and textually at the same time.

To this end, we add an inverted list of  $\langle kw, wt \rangle$  on each node of the G-tree, which represents the maximum  $\text{tf} \cdot \text{idf}$ , denoted  $wt$ , for keyword  $kw$  under the sub-tree of this node. Formally, given a G-tree node  $n$  and a keyword  $w$ ,  $wt$  is

$$wt(w, n) = \max_{v \in \text{SubTree}(n)} \text{tf}(w, v) * \text{idf}(w, \mathcal{V}), \quad (11)$$

where  $\text{SubTree}(n)$  denotes all the vertices in the sub-tree of node  $n$ . Thus, for a query with the keyword set  $q.\mathcal{W}$ ,

$$\text{TXDist}(q, n) = \sum_{w \in q.\mathcal{W}} wt(w, n). \quad (12)$$

Note that,  $wt(w, n)$  is only calculated once in a bottom-up manner offline. Thus,  $\text{TXDist}(q, n)$  can be aggregated online. Similar to the distance matrices, we maintain an offset pointer on each G-tree node, and  $wt(w, n)$  is physically stored in continuous memory/disk storage. Once we get both  $\text{SPDist}(v_q, n)$  and  $\text{TXDist}(q, n)$ , we can pack them together to obtain  $\mathcal{F}(q, n)$  by Equation (1).

Based on  $\mathcal{F}(q, n)$ , we can devise the  $k^2N^2$  algorithm by making minor modifications on Algorithm 2.

- 1) Pre-compute the  $wt(w, n)$  before the query;
- 2) Replace all  $\text{SPDist}(v_q, n)$  into  $\mathcal{F}(q, n)$ ;
- 3) Change the  $\mathcal{Q}$  into a max priority-queue;
- 4) Alter the intention of  $\mathcal{T}_{\min}$  into  $\mathcal{T}_{\max}$ . It is worth noting that the textual part of  $\mathcal{T}_{\max}$  is the global one, i.e.,  $\sum_{w \in q.\mathcal{W}} wt(w, \text{root})$ .

Algorithm 3 presents the algorithmic description for  $k^2N^2$ . Here, we use an example to show the process.

---

#### Algorithm 3. $K^2N^2$ ( $q = \langle v_q, \mathcal{C}, \mathcal{W}, k \rangle, \mathcal{F}, \mathcal{G}$ )

---

**Input:**  $q = \langle v_q, \mathcal{C}, \mathcal{W}, k \rangle$ : A query;  $\mathcal{F}$ : A ranking function;  $\mathcal{G}$ : A G-tree  
**Output:**  $\mathcal{R}$ : The top- $k$  result list;  
1 Line 1 to Line 4 in Algorithm 2;  
2 Compute global  $wt(w, n)$  based on occurrence list  $\mathcal{L}$ ;  
3 Initialize  $\mathcal{T}_n = \text{leaf}(v_q)$  and  $\mathcal{T}_{\max}$ ;  
4 Line 6 to Line 15 in Algorithm 2, replace all  $\text{SPDist}(v_q, n)$  into  $\mathcal{F}(q, n)$ , replace all  $\mathcal{T}_{\min}$  into  $\mathcal{T}_{\max}$ ;  
5 **Function** UpdateT( $\mathcal{T}_n, \mathcal{T}_{\max}, \mathcal{Q}$ )  
6  $\mathcal{T}_n \leftarrow \mathcal{T}_n.\text{father}$  and update  $\mathcal{T}_{\max}$ , where  $\text{TXDist}(v_q, \mathcal{T}_n) = \sum_{w \in q.\mathcal{W}} wt(w, \text{root})$ ;  
7 **foreach**  $c \in \mathcal{L}(\mathcal{T}_n)$  **do**  
8  $\mathcal{Q}.\text{Enqueue}(\langle \mathcal{F}(v_q, c) \rangle)$ ;

---

**Example 4.** Consider the query  $q = \langle v_4, \{v_3, v_9, v_{15}\}, \{\text{'park'}, 2\}$  in Fig. 1, where  $\alpha = 0.5$ ,  $M_S = 21$  and  $M_T = 0.5$ . To save space, we omit the process here. The top-2 answers are  $\langle v_{15}, 0.50 \rangle$  and  $\langle v_3, 0.45 \rangle$ .

**Theorem 1.** Algorithms 1(SPSP), 2( $k$ NN) and 3( $k^2N^2$ ) can correctly report all answers.

**Proof (SPSP).** According to Lemma 1, let  $b_1 \in \mathcal{G}_1, b_2 \in \mathcal{G}_2, \dots, b_p \in \mathcal{G}_p$  denote  $p$  consecutive borders on the shortest path of  $\text{SP}(u, v)$ . Suppose there exists another path  $\text{SP}'(u, v) = ub'_1b'_2 \dots b'_pv$  that  $|\text{SP}'(u, v)| < |\text{SP}(u, v)|$ , where  $b'_x \in \mathcal{G}_x$ . If  $\exists i$  such that  $b'_i = b_i$ , then  $|\text{SP}'(u, b'_i)| + |\text{SP}'(b'_i, v)| = |\text{SP}(u, b_i)| + |\text{SP}(b_i, v)|$ , thus,  $\forall i$  that  $b'_i \neq b_i$ . However, according to Equation (7),  $b'_1$  must be selected since  $|\text{SP}'(u, b'_1)| + |\text{SP}'(b'_1, v)| < |\text{SP}(u, b_1)| + |\text{SP}(b_1, v)|$ , i.e.,  $b'_1 = b_1$ , which contradicts that  $b'_1 \neq b_1$ . Thus SPSP can correctly find the shortest-path answer.

$k$ NN and  $k^2N^2$ . Suppose  $\exists v \in \mathcal{C} - \mathcal{R}$  that  $\text{SPDist}(v_q, v) < \text{SPDist}(v_q, \bar{v} \in \mathcal{R})$ , where  $\text{SPDist}(v_q, \forall v' \in \mathcal{R}) \leq \text{SPDist}(v_q, \bar{v})$ . Thus,  $\text{SPDist}(v_q, \text{leaf}(v)) < \text{SPDist}(v_q, \bar{v} \in \mathcal{R})$ . According to the best-first search,  $\text{leaf}(v)$  must dequeue before  $\bar{v}$ , hence,  $v$  must be one of top- $k$  answers, i.e.,  $v \in \mathcal{R}$ . This contradicts the assumption of  $v \in \mathcal{C} - \mathcal{R}$ . Therefore,  $k$ NN can correctly find top- $k$  answers. The correctness of  $k^2N^2$  can be proved similarly.  $\square$

### 4.4 Time and Space Complexity

#### 4.4.1 Single-Pair Shortest Path Query

**Time complexity.** If two vertices  $u$  and  $v$  are within the same leaf node, its time cost is  $\mathcal{O}(\tau \log \tau)$  for local Dijkstra search. Otherwise, it will at most scan a series of distance matrices from  $\text{leaf}(u)$  to  $\text{leaf}(v)$  bypassing the *root*, and the worst-case time complexity is the size of distance matrices scanned. According to analysis in Section 3.4, we can derive the time cost to be  $\mathcal{O}(2 \cdot \sum_{i=1 \dots \mathcal{H}} \log_2^2 f \cdot |\mathcal{V}| / f^i) = \mathcal{O}(\log_2^2 f \cdot |\mathcal{V}|)$ .

#### 4.4.2 $k$ NN Query and $k^2N^2$ Query

**Time Complexity.** The  $k$ NN search consists of two parts. The first one is the local Dijkstra search within MINDIST-INSIDE-LEAF. The time complexity is  $\mathcal{O}(\tau \log \tau)$ . The second one is MINDIST-OUTSIDE-LEAF. Since the materialized algorithm will only access each node of G-tree once, the worst-case time of MINDIST-OUTSIDE-LEAF is the total size of the distance matrices of G-tree, i.e.,  $\mathcal{O}(\log_2^2 f \cdot \log_f \frac{|\mathcal{V}|}{\tau} \cdot |\mathcal{V}|)$ . To sum up, the overall time cost is  $\mathcal{O}(\tau \log \tau + \log_2^2 f \cdot \log_f \frac{|\mathcal{V}|}{\tau} \cdot |\mathcal{V}|)$ . In practice, the complexity is much smaller than the worst-case complexity.

#### 4.4.3 Space Complexity

**Object occurrence list  $\mathcal{L}(n)$ .** The occurrence list  $\mathcal{L}(n)$  for leaf nodes of G-tree requires  $\mathcal{O}(|\mathcal{C}|)$  space. In the worst case,  $\mathcal{L}(n)$  for all non-leaf nodes is equal to  $\mathcal{O}(f \cdot \frac{|\mathcal{V}|}{\tau})$ . To sum up, the overall space cost is  $\mathcal{O}(|\mathcal{C}| + f \cdot \frac{|\mathcal{V}|}{\tau})$ . Compared with distance matrices of G-tree, the space cost of  $\mathcal{L}(n)$  is negligible.

**Inverted list for  $k^2N^2$ .** Let  $|\Psi|$  denotes  $\sum_{v \in \mathcal{V}} |v.\mathcal{W}|$ , i.e., the total number of keywords for all vertices. Suppose that each keyword on vertices is distinct (the worst case), thus, given a keyword  $w$  on  $v$ ,  $wt(w, n)$  is materialized in all nodes on the branch from  $\text{leaf}(v)$  to the root. Therefore, the overall space cost is  $\mathcal{O}(|\Psi| \mathcal{H})$ .

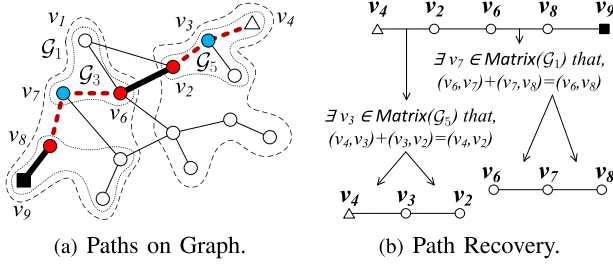


Fig. 7. An example of shortest-path recovery.

## 5 SHORTEST PATH RECOVERY

It is worth noting that all the algorithms in Section 4 returns shortest-path distance rather than vertex-by-vertex path. However, the latter is important and useful (e.g., in navigation system). In this section, we discuss how to recover the path from the query location  $v_q$  to an answer  $v_a \in \mathcal{R}$ .

### 5.1 Overview

We consider the case that  $v_q$  and  $v_a$  are in different nodes. Since we facilitate the graph distance calculation by the assembly-based method, we can only get a list of selected borders from  $v_q$  to  $v_a$ , i.e. the imperfect shortest path  $SP'(v_q, v_a) = v_q b_1 b_2 \dots b_m v_a$ . As there may be no direct edges between two adjacent borders, e.g.,  $\langle b_i, b_{i+1} \rangle$ , we need to add some other vertices between them to generate the real shortest path  $SP(v_q, v_a)$ .

The main idea is to apply the divide-and-conquer technique to iteratively add new vertices into the  $SP'(v_q, v_a)$ . Given two vertices  $v$  and  $v'$ , once we can find a  $v_x$  such that  $SPDist(v, v_x) + SPDist(v_x, v') = SPDist(v, v')$ , we can confirm that  $v_x$  is one vertex on  $SP(v, v')$ . For example, in Fig. 7, to compute the shortest-path distance from  $v_4$  to  $v_9$ , we get  $SP'(v_4, v_9) = v_4 v_2 v_6 v_8 v_9$ . As there is no edge between  $v_4$  and  $v_2$ , we need to find a vertex, i.e.,  $v_3$ , to add between them (since  $SPDist(v_4, v_2) = SPDist(v_4, v_3) + SPDist(v_3, v_2)$ ). Similarly we add vertex  $v_7$  between  $v_6$  and  $v_8$ . Thus the vertex-by-vertex shortest path is  $SP(v_4, v_9) = v_4 v_3 v_2 v_6 v_7 v_8 v_9$ .

### 5.2 Shortest Path Computation

Formally, given a section of a shortest path  $SP'(v_q, v_a) = v_q b_1 b_2 \dots b_m v_a$ , for each two adjacent vertices pair  $\langle v, v' \rangle$  in  $SP'(v_q, v_a)$ , there are generally three cases

Case 1:  $v, v'$  are borders in different leaf nodes:

**Lemma 2.** Given two borders  $v$  and  $v'$  in different leaf nodes, there must exist a border  $b$  in the distance matrix from  $LCA(v, v')$  to the root such that  $v, v', b$  appear in the matrix and  $SPDist(v, v') = SPDist(v, b) + SPDist(b, v')$ .

**Proof.** Suppose there is at least one vertex except  $v$  and  $v'$  on  $SP(v, v')$  which is in the subgraph of  $LCA(v, v')$ . According to Lemma 1, there must exist a border  $b$  on  $SP(v, v')$  such that  $v, v'$  and  $b$  are in the distance matrix of  $LCA(v, v')$  by Definition 3. Otherwise, all  $SP(v, v')$  vertices except  $v$  and  $v'$  are outside of  $LCA(v, v')$ . Since the root contains all vertices and paths, there must exist one ancestor node  $\mathcal{G}'$  of  $LCA(v, v')$  which contains at least one vertex except  $v$  and  $v'$  of  $SP(v, v')$ . Therefore, one border  $b$  must exist in the distance matrix of  $\mathcal{G}'$ , together with  $v, v'$ .  $\square$

For example, consider  $\langle v_6, v_8 \rangle$ , where  $v_6 \in \mathcal{G}_3$  and  $v_8 \in \mathcal{G}_4$ .  $LCA(v_6, v_8) = \mathcal{G}_1$ . We find a border  $v_7$  in the distance matrix of  $\mathcal{G}_1$  such that  $SPDist(v_6, v_8) = SPDist(v_6, v_7) + SPDist(v_7, v_8)$ . Then we add  $v_7$  into  $\langle v_6, v_8 \rangle$  and check  $\langle v_6, v_7 \rangle$  and  $\langle v_7, v_8 \rangle$ . As  $w(v_6, v_7) = SPDist(v_6, v_7)$  and  $w(v_7, v_8) = SPDist(v_7, v_8)$ , hence, the shortest path is  $v_6 v_7 v_8$ .

Case 2:  $v, v'$  are borders in the same leaf nodes:

**Lemma 3.** Given two borders  $v$  and  $v'$  in the same leaf node, if there does not exist a non-border vertex  $v_x$  such that  $SPDist(v, v') = SPDist(v, v_x) + SPDist(v_x, v')$ , then there must exist a border  $b$  in the distance matrix of ancestors of the node such that  $v, v', b$  appear in the matrix and  $SPDist(v, v') = SPDist(v, b) + SPDist(b, v')$ .

**Proof.** Lemma 3 can be proved similar to Lemma 2 on provided that  $LCA(v, v') = \text{leaf}(v) = \text{leaf}(v')$ .  $\square$

For example, consider  $\langle v_1, v_7 \rangle$  in  $\mathcal{G}_3$ . We find a border  $v_6$  such that  $SPDist(v_1, v_7) = SPDist(v_1, v_6) + SPDist(v_6, v_7)$ . We add  $v_6$  into  $\langle v_1, v_7 \rangle$  and check  $\langle v_1, v_6 \rangle$  and  $\langle v_6, v_7 \rangle$ . As  $w(v_1, v_6) = SPDist(v_1, v_6)$  and  $w(v_6, v_7) = SPDist(v_6, v_7)$ , the shortest path is  $v_1 v_6 v_7$ .

Case 3:  $v, v'$  are in the same leaf and  $v$  is not a border:

**Lemma 4.** Given a non-border vertex  $v$  and a border  $v'$  in the same leaf node, if there does not exist a border  $b$  such that  $SPDist(v, v') = SPDist(v, b) + SPDist(b, v')$ , then there must exist a shortest path from  $v$  to  $v'$  which only includes vertices in the leaf node.

**Proof.** If not all vertices on  $SP(v, v')$  are within  $\text{leaf}(v)$ , there must exist another border  $b$  on  $SP(v, v')$  according to Lemma 1. Otherwise,  $SP(v, v')$  only includes vertices in  $\text{leaf}(v)$ . Thus the lemma is proved.  $\square$

For example, consider  $\langle v_4, v_2 \rangle$  in  $\mathcal{G}_5$ . As there is only one border  $v_2$  in  $\mathcal{G}_5$ , we cannot add borders into the path. We visit  $v_2$ 's neighbor vertex and find vertex  $v_3$  such that  $SPDist(v_4, v_2) = w(v_4, v_3) + SPDist(v_3, v_2)$ . We add  $v_3$  into  $\langle v_4, v_2 \rangle$ . Since both  $(v_4, v_3)$  and  $(v_3, v_2)$  are edges in  $\mathcal{G}$ , thus the shortest path is  $v_4 v_3 v_2$ .

Based on these properties, we propose a recursive algorithm to recover the vertex-by-vertex path for a given  $\langle v, v' \rangle$  pair, which is shown in Algorithm 4.

#### Algorithm 4. PREC ( $\langle v, v' \rangle, \mathcal{G}$ )

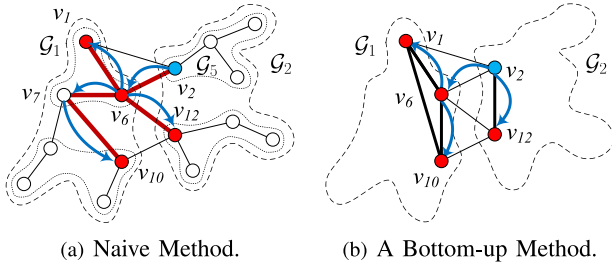
**Input:**  $\langle v, v' \rangle$  : Two adjacent vertices;  $\mathcal{G}$  : A G-tree

**Output:**  $\mathcal{R}$  : The vertex-by-vertex path from  $v$  to  $v'$  ;

```

1 if  $\langle v, v' \rangle \in \mathcal{E}$  then return  $\emptyset$  ;
2 if  $\text{leaf}(v) = \text{leaf}(v')$  then
2    $b \leftarrow \text{FIND\_BORDER}(LCA(v, v'), \text{root})$  ;
3   if  $b = \emptyset$  then
4      $u \leftarrow \text{FIND\_NEIGHBOR}(v)$  ;
5     return  $(v, u, \text{PREC}(u, v'), v')$ 
6   else return  $(v, \text{PREC}(v, b), b, \text{PREC}(b, v'), v')$ 
7 else if  $\text{leaf}(v) \neq \text{leaf}(v')$  then
8    $b \leftarrow \text{FIND\_BORDER}(LCA(v, v'), \text{root})$  ;
10 return  $(v, \text{PREC}(v, b), b, \text{PREC}(b, v'), v')$ 
```

*Space and time complexity.* The space complexity is the number of vertices in  $SP(v_q, v_a)$ , denoted by  $N_v$ . Let  $\mathcal{B}_{max}$  denote the maximal number of borders in each node. In

Fig. 8. Distance matrix computation (from  $v_2$ ).

each step, we add a border from the LCA node to the root with time complexity  $\mathcal{O}(\mathcal{H}\mathcal{B}_{max})$  or enumerate vertices in a leaf node with time complexity  $\mathcal{O}(\tau\mathcal{B}_{max})$ . Thus the overall time complexity is  $\mathcal{O}((\mathcal{H} + \tau)\mathcal{B}_{max}N_v)$ .

## 6 DISCUSSIONS

### 6.1 Computing Distance Matrix Efficiently

The naive way to generate the distance matrices is to calculate the graph distances node by node, and pair by pair. This is obviously ineffective for very large datasets.

To address this issue, we propose a bottom-up method to efficiently compute the distance matrices. The basic idea is that we first compute the distance matrices in the lower level and then use these distances to compute the distance matrices in the upper level (Fig. 8b). Given two borders  $u$  and  $v$  from two sibling nodes  $\mathcal{G}(u)$  and  $\mathcal{G}(v)$ , consider a shortest path between  $u$  and  $v$ , e.g.,  $SP(u, v) = uu_1u_2 \dots u_xw_1w_2 \dots w_zv_yv_{y-1}v_1v$ . Suppose  $u_x$  is the last vertex in  $\mathcal{G}(u)$  and  $v_y$  is the first vertex in  $\mathcal{G}(v)$ . Obviously  $u_x$  is a border of  $\mathcal{G}(u)$  and  $v_y$  is a border of  $\mathcal{G}(v)$ . As  $SPDist(u, u_x)$  has been computed in the lower level, we can skip vertexes  $u_1u_2 \dots u_{x-1}$  and directly use the “shortcut”  $\langle u, u_x \rangle$ . Our algorithm works as follows.

- 1) Initially, for each leaf node, we use the Dijkstra algorithm to compute the shortest-path distance between any two borders in the leaf node.
- 2) We remove all non-border vertices in the leaf node and add shortcuts between any two borders of the leaf node.
- 3) We move to the parent of leaf nodes and use the Dijkstra algorithm to compute the shortest-path distance between any two borders in the parent based on the updated graph.
- 4) We repeat steps 2 and 3 and terminate the algorithm if we have processed the root node.

### 6.2 Extension to Directed Graphs

G-tree can support directed graphs with a minor change. First, in the distance matrix, we keep the shortest distances of directed paths from a vertex to a border/vertex. We only need to slightly modify the pre-computation method in Section 6.1 to compute the directed distances. Second, our method relies on using the assembly-based method to implement the  $SPDist$  function. Nevertheless, the assembly based method still works for directed graphs based on the following reasons. (1) The closure property (Lemma 1) still holds for directed graphs, i.e., given a subgraph  $\mathcal{G}' = \langle \mathcal{V}', \mathcal{E}' \rangle$ , any directed path from  $u \in \mathcal{V}'$  to  $v \notin \mathcal{V}'$  must contain at

TABLE 1  
Datasets

Data	Description	# Vertices	# Edges
CAL	California(Undirected)	21,048	21,693
SF	San Francisco(Undirected)	174,956	223,001
COL	Colorado(Undirected)	435,666	528,533
FLA	Florida(Undirected)	1,070,376	1,356,399
E-USA	East USA(Undirected)	3,598,623	4,389,057
C-USA	Center USA(Undirected)	14,081,816	17,146,248
USA	USA(Undirected)	23,947,347	29,166,672
WA	Washington(Directed)	514,654	1,246,353

least one border in  $\mathcal{B}(\mathcal{G}')$ . (2) We can still use the algorithm in Section 4.1 to compute the shortest distance of a directed path. Third, we slightly modify the Dijkstra algorithm to support the directed graphs.

### 6.3 On Planarity of Network

Though we assume that the underlying road network is a planar graph, this is only for the purpose of analyzing the bounds of space complexity. The experimental result in Fig. 10 shows that the index sizes of G-tree on real networks actually follow the derived bound, and G-tree works very well with non-planar networks, and the bridges or tunnels will not affect the correctness of the algorithms.

## 7 EXPERIMENTS

**Datasets.** We used eight real-world datasets<sup>1,2,3</sup> with various sizes from 20,000 vertices to 24 million vertices, which are widely used in previous studies [13]. The statistics of these datasets are shown in Table 1.

**Settings.** For the G-tree, the default fanout is  $f = 4$  and  $\tau$  is set to 64(CAL), 128(SF), 128(COL), 256(FLA), 256(E-USA), 512(C-USA) and 512(USA). To evaluate the efficiency of SPSP query, we randomly chose 10,000 pairs of vertices for each dataset; to evaluate the performance of  $kNN$  and  $k^2N^2$  query, we randomly chose 100 vertices as the query location, and for each location we generated 100 groups of objects, thus we had 10,000 queries for each query set. For objects we uniformly selected 0.0001, 0.001, 0.01, 0.1, 1 of vertices from the dataset as objects (default is 0.01). For  $k$ , we used 1, 5, 10, 20, 50 (default is 10).

For SPSP query, we compared our G-tree with *TNR* [1], *CH* [5] and *A\** algorithm. For  $kNN$  query, we compared with *INE* [16], *SILC* [19] and *ROAD* [13].<sup>4</sup> For  $k^2N^2$  query, since no previous works addressed the same problem as ours, we implemented a heuristic method, denoted by *NKS*, which is based on *INE* but used the  $k$ th answer’s boundary to do pruning. All the algorithms were implemented in C++. All experiments were conducted on a Linux computer with Intel 2.50 GHz CPU and 16 GB memory.

1. <http://www.cs.fsu.edu/~lifeifei/SpatialDataset.htm>

2. <http://www.dis.uniroma1.it/challenge9/index.shtml>

3. <http://depts.washington.edu/giscup/roadnetwork>

4. We set “No. of levels( $l$ )” to 6, 7, 7 and 8 respectively on CAL, SF, COL and FLA, so that the height of ROAD is the same with G-tree.



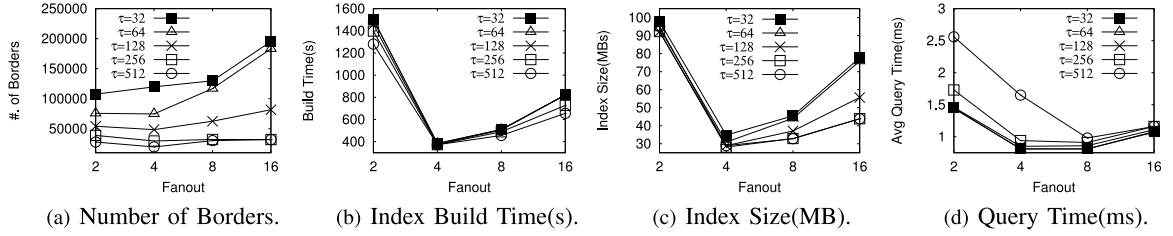


Fig. 9. Evaluation on parameters:  $f$  and  $\tau$  (COL dataset,  $k = 10$ , 1 percent uniform vertices as objects).

### 7.1 Evaluation on G-tree Parameters: $f$ and $\tau$

We evaluated the impact of  $f$  (the fanout) and  $\tau$  (the number of vertices in a leaf node) of the G-tree by investigating the number of borders, the index size, the index build time and the query time. We varied  $\tau$  in {32, 64, 128, 256, 512} and  $f$  in {2, 4, 8, 16}.

Fig. 9 shows the results. We made two observations. First, with the increase of  $f$ , the number of borders and the index size, the index build time and the query time first decreased and then increased. Our method achieved the best results when  $f = 4$ . The main reason is that larger  $f$  will generate more borders as there are more partitions, however, larger  $f$  will reduce the height of G-tree and the number of nodes that need to be partitioned. Second, with the increase of  $\tau$ , the number of borders, the index size and the build time decreased, since bigger  $\tau$  results in smaller tree size. However, larger  $\tau$  incurs more time on local Dijkstra search within one leaf node. To balance between the query efficiency and indexing size, we selected  $\tau = 128$  for COL dataset as a trade off.

### 7.2 Evaluation on G-tree Construction

We evaluated the time and space overhead of indexing, and compared with SILC, ROAD, TNR and CH. Fig. 10 illustrates the index sizes and index construction time. Note that, the index size of G-tree included tree structure, distance matrices, leaf( $u$ ) hashtable and the original graph.

G-tree outperformed ROAD, SILC and TNR both in construction time and index sizes, and achieved the same result as CH. For the construction time, G-tree was better than ROAD and TNR by an order of magnitude and three orders of magnitude than SILC. For index sizes, on COL, G-tree consumed 37.0 MB, ROAD used 145 MB and SILC required 1,543 MB; on USA, G-tree consumed 2.6 GB but TNR cost 9.1 GB. This is because SILC required  $\mathcal{O}(|V|^{1.5})$  space to compute all-pair shortest paths. ROAD incurred a large number of shortcuts pre-computation. TNR needed to calculate the distances between all pairs of transit nodes.

### 7.3 Evaluation on Query Efficiency

**Evaluation on SPSP search.** We compared the SPSP search efficiency of G-tree with TNR, CH and A\*. First, we tested the performance by randomly choosing two nodes on each dataset. Second, we tested the performance varying the distance between two nodes from near to far. We partitioned USA road network into  $128 \times 128$  grids, and generated nine test cases  $Q_i (i = 0 \dots 8)$  which indicated that two nodes were exactly  $2^{i-1}$  grids away (For  $Q_0$ , two nodes were within one grid). The result is shown in Fig. 11.

For the random dataset, G-tree is approximately five times faster than CH, and nearly three orders of magnitudes faster than A\*. Though TNR was faster than G-tree, it was actually trading space for time. Thus, it is not surprising that TNR took almost 20 times more space than G-tree (see Fig. 10b). When two nodes were within one grid, G-tree and CH were two orders of magnitudes faster than TNR, since TNR failed to handle local queries efficiently by Dijkstra algorithm.

**Evaluation on kNN search.** We evaluated the  $k$ NN search efficiency of G-tree, ROAD, SILC and INE by varying the number of answers  $k$ , the number of objects  $|C|$ , datasets, and distances from query location to top- $k$  answers. As default, we tested on COL dataset, set  $k = 10$ , and  $|C| = 0.01|V|$ . Note that, we used 1 percent vertices as objects, since in most cases the number of candidate objects were relatively few compared with the size of the road network. Even in other cases, we can see that it was much easier for top- $k$  search when  $|C| = |V|$  (Fig. 12b), since it would quickly find top- $k$  answers around the query location.

**kNN search by varying  $k$ .** The result is shown in Fig. 12a. We can see that G-tree outperformed the state-of-the-arts ROAD and SILC by 2-3 orders of magnitude, and outperformed INE by one order of magnitude. Since SILC had to search multiple quadrees to find distances between query location and objects, this operation was very costly and inefficient for larger  $k$ . As ROAD employs an expansion-based method, it can only prune the nodes which have no objects and cannot use distance-based pruning. Note that, INE was faster than ROAD and SILC on this test—since objects were

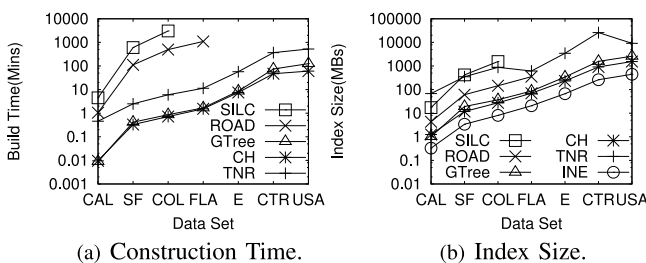


Fig. 10. Index comparison.

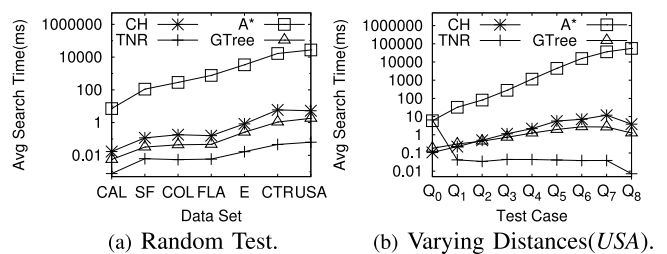


Fig. 11. Evaluation on SPSP Query.

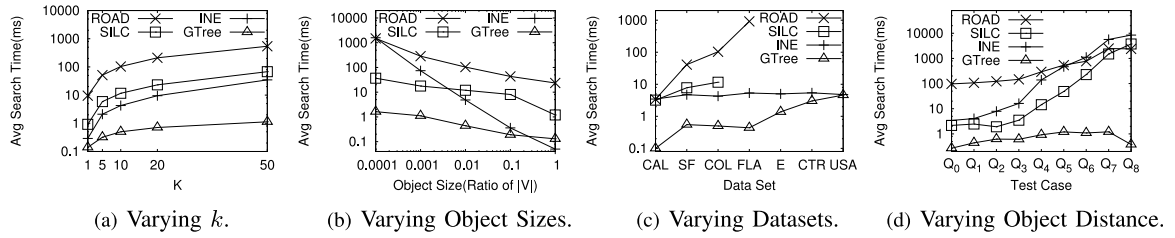


Fig. 12. Performance comparison on  $k$ NN search (Default: *COL* dataset,  $k = 10$ , 1 percent uniform vertices as objects).

uniformly distributed *INE* could quickly come up with top- $k$  answers with Dijkstra algorithm.

*k*NN search by varying object sizes. Fig. 12b shows the result. We made three observations. First, the efficiency of these four methods increased for larger  $|C|$ , as more objects implied less top- $k$  finding overhead from query location. Second, G-tree outperformed three baseline methods significantly for different  $|C|$  settings. Third, when  $|C| = |V|$ , G-tree achieved comparable performance with *INE* by using identical leaf-node Dijkstra search. However, *SILC* and *ROAD* still performed poorly, as *SILC* used morton block to compress/decompress path information and *ROAD* needed to traverse numbers of tree-like entries on RO index for every vertex expansion.

*k*NN search by varying datasets. Since *ROAD* and *SILC* took a mass of pre-processing time and consumed large amount of memory,<sup>5</sup> we only tested *ROAD* on *CAL*, *SF*, *COL* and *FLA*, and *SILC* on *CAL*, *SF* and *COL*. Fig. 12b shows the result. G-tree outperformed *ROAD* and *SILC* on every dataset by 1-2 orders of magnitude. However, *INE*'s performance degenerated slowly with the increase of the data size. This is because that the efficiency of *INE* actually depends on the minimum distances to top- $k$  objects. Since we had assumed that objects are uniform distributed and  $|C| = 0.01|V|$ , such distance did not change tremendously as the data size become larger. However, we will see that *INE* performed very poorly when objects were far away from the query location in the next test.

*k*NN search by varying object distances. We adopted the same setting which was used for SPSP queries in Fig. 11b to guarantee the objects of  $Q_i$  were at least  $2^{i-1}$  grids away. Fig. 12c shows the result. G-tree significantly outperformed *ROAD*, *SILC* and *INE*, even by 3-4 orders of magnitude. All baseline methods performed poorly since they had to traverse long distance paths before they accessed all top- $k$  answers. In contrast, G-tree can quickly locate top- $k$  objects within fewer tree nodes, no matter how far away objects were. This test proves G-tree's superiority and efficiency for  $k$ NN search.

*Evaluation on  $k^2N^2$  search.* We obtained 14 million POIs from a popular directory website,<sup>6</sup> extracted the location and keywords, and attached each of them to the nearest vertex on the road network. Table 2 shows the statistics.

We proposed a baseline method, namely *NKS*, which adopted the framework of *INE* and used the  $k$ th obtained answer  $v^k$  to prune unpromising Dijkstra expansion, i.e., stop the expansion at  $v_x$  when  $\text{SPDist}(q, v_x) + \max_{v_i \in C} \text{TXDist}(q, v_i) < \mathcal{F}(q, v^k)$ .

We compared G-tree with *NKS* by varying the number of answers  $k$ , datasets, the weight factor  $\alpha$  and the number of keywords  $|W|$ . As default, we tested on *COL* dataset, and set  $k = 10$ ,  $|C| = 0.01|V|$ ,  $\alpha = 0.5$  and  $|W| = 1$ .

*$k^2N^2$  search by varying  $k$ .* Fig. 13a shows the result. We can see that G-tree outperformed *NKS* by one order of magnitude. This is because that G-tree can fully utilize spatial and textual pruning at the query time, thus, even for  $k = 50$ , G-tree can answer a query within 1 ms on average.

*$k^2N^2$  search by varying datasets.* Fig. 13b shows that G-tree achieved excellent efficiency and scalability as the data size increased. Compared with *NKS*, G-tree was nearly 3 orders of magnitude faster on *USA*.

*$k^2N^2$  search by varying  $\alpha$ .* According to Equation (1),  $\alpha = 0$  considers only the textual part, while setting  $\alpha = 1$  considers only the spatial part. Fig. 13c shows the result. We can see that *NKS* improved its performance as  $\alpha$  increased. For  $\alpha = 0$ , *NKS* was extremely inefficient, while it was the opposite when  $\alpha = 1$ . The reason is that *NKS* was simply based on Dijkstra algorithm, thus, it was incapable to handle purely textual queries effectively. Note that G-tree is more efficient when  $\alpha = 0$ , as the spatial distance calculation was not required in this case.

*$k^2N^2$  search by varying  $|W|$ .* We can see that G-tree still outperformed *NKS* by at least 1 order of magnitude, which is shown in Fig. 13d. Note that, the variation on  $|W|$  did not affect G-tree's performance significantly, and had minor impact on *NKS*. The reason is that, according to Equations (1) and (2), the increase of  $|W|$  is actually augmenting the weight of the textual part. Hence, this is equal to decreasing the  $\alpha$  to some extent.

## 7.4 Evaluation on Path Recovery

For SPSP query, we compared G-tree with *TNR* and *CH*. For  $k$ NN query, we compared with *ROAD*, *SILC* and *INE*. We used the same setting as the existing experiments. The result is shown in Fig. 14.

For SPSP, G-tree outperformed *TNR* but five times slower than *CH*. However, G-tree was five times faster than

TABLE 2  
Textual Statistics on Seven Datasets

Size	CAL	SF	COL	FLA	E-USA	C-USA	USA
$ POIs / V $	1.01	2.62	0.53	0.84	1.03	0.42	0.54
Distinct Kwd	29,439	165,607	97,165	230,863	632,823	806,494	1,443,560
Average Kwd	6.58	6.43	6.71	6.93	6.56	6.83	6.73
Min-Max Kwd	[3,21]	[3,21]	[3,21]	[3,27]	[2,25]	[3,34]	[1,36]
Inv List(MB)	5.3	46.3	44.7	147.5	527.3	993.8	2035.5

5. We estimated 4.8 days for *ROAD*, and 36.5 GB for *SILC* on *USA*.

6. <http://www.factual.com>

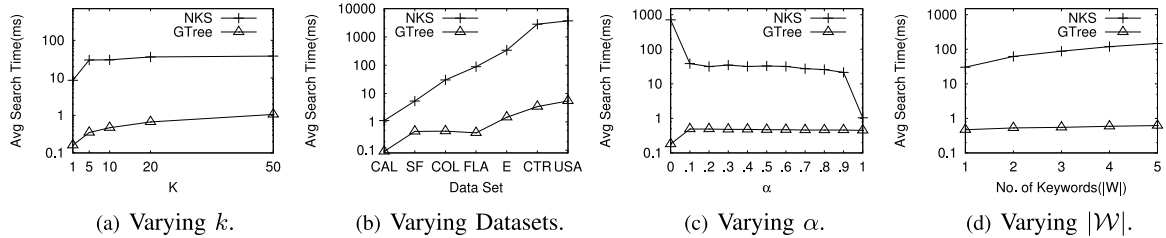


Fig. 13. Performance comparison on  $k^2 N^2$  search (Default: COL dataset,  $k = 10$ ,  $|C| = 0.01|V|$ ,  $|W| = 1$  and  $\alpha = 0.5$ ).

CH for distance query in Fig. 11a. For  $k$ NN, G-tree still beats ROAD, SILC and INE even by comparing the total time ( $1*k$ NN +  $10*Path$  Recovery).

## 7.5 Scalability

Fig. 10b indicated that the index size of G-tree increased linearly with the increases of the data size. Figs. 11a, 12c, and 13b had displayed G-tree's superior scalability for three different types of queries.

## 7.6 Evaluation on Directed Graph

We evaluated the search efficiency of G-tree on a directed graph WA. We set  $f = 4$  and  $\tau = 128$ . The G-tree index was 52.13 MB and the index building time was 16 s. We compared the  $k$ NN performance with SILC, ROAD and INE by varying  $k$  and object size  $|C|$ , as shown in Fig. 15. G-tree still significantly outperformed existing methods. The results are consistent with those on undirected graphs.

## 8 CONCLUSION

In this paper we have proposed an efficient and scalable index on road networks. We devised an assembly-based method to efficiently calculate shortest-path distance. We proposed efficient search algorithms to answer SPSP,  $k$ NN and  $k^2 N^2$  queries. We also discussed path recovery issue. Experimental results show G-tree's theoretical and practical superiority over state-of-the-art methods.

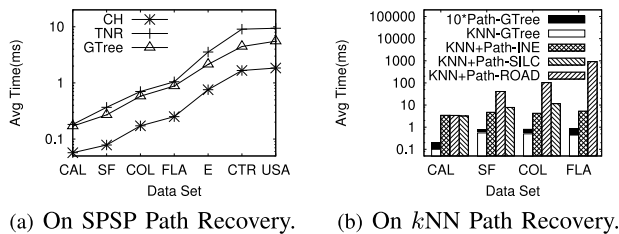


Fig. 14. Evaluation on path recovery.

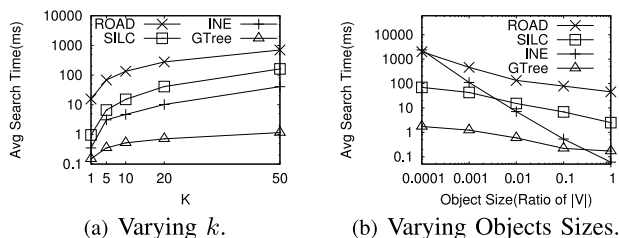


Fig. 15. Evaluation on directed graphs.

## ACKNOWLEDGMENTS

This work was partly supported by the 973 Program of China (2015CB358700 and 2011CB302206), and the NSFC project (61272090, 61373024 and 61422205), YETP0105, Tencent, Huawei, SAP, the "NEXt Research Center" (WBS: R-252-300-001-490), and the FDCT/106/2012/A3. Guoliang Li is the corresponding author.

## REFERENCES

- [1] H. Bast, S. Funke, and D. Matijevic, "Transit—Ultrafast shortest-path queries with linear-time preprocessing," in *Proc. 9th DIMACS Implementation Challenge*, 2006, pp. 175–192.
- [2] H.-J. Cho and C.-W. Chung, "An efficient and scalable approach to CNN queries in a road network," in *Proc. 31st Int. Conf. Very Large Data Bases*, 2005, pp. 865–876.
- [3] G. Cong, C. S. Jensen, and D. Wu, "Efficient retrieval of the top-k most relevant spatial web objects," *Proc. VLDB Endowment*, vol. 2, no. 1, pp. 337–348, 2009.
- [4] I. D. Felipe, V. Hristidis, and N. Rishe, "Keyword search on spatial databases," in *Proc. IEEE 24th Int. Conf. Data Eng.*, 2008, pp. 656–665.
- [5] R. Geisberger, P. Sanders, D. Schultes, and D. Delling, "Contraction hierarchies: Faster and simpler hierarchical routing in road networks," in *Proc. 7th Int. Conf. Experimental Algorithm*, 2008, pp. 319–333.
- [6] H. Hu, D. L. Lee, and V. C. S. Lee, "Distance indexing on road networks," in *Proc. 31st Int. Conf. Very Large Data Bases*, 2006, pp. 894–905.
- [7] H. Hu, D. L. Lee, and J. Xu, "Fast nearest neighbor search on road networks," in *Proc. 10th Int. Conf. Extending Database Technol.*, 2006, pp. 186–203.
- [8] C. S. Jensen, J. Kolár, T. B. Pedersen, and I. Timko, "Nearest neighbor queries in road networks," in *Proc. 11th ACM Int. Symp. Adv. Geographic Inf. Syst.*, 2003, pp. 1–8.
- [9] N. Jing, Y.-W. Huang, and E. A. Rundensteiner, "Hierarchical encoded path views for path query processing: An optimal model and its performance evaluation," *IEEE Trans. Knowl. Data Eng.*, vol. 10, no. 3, pp. 409–432, May/Jun. 1998.
- [10] S. Jung and S. Pramanik, "An efficient path computation model for hierarchically structured topographical road maps," *IEEE Trans. Knowl. Data Eng.*, vol. 14, no. 5, pp. 1029–1046, Sep/Oct. 2002.
- [11] G. Karypis and V. Kumar, "Analysis of multilevel graph partitioning," in *Proc. ACM/IEEE Conf. Supercomput.*, 1995, p. 29.
- [12] M. R. Kolahdouzan and C. Shahabi, "Voronoi-based k nearest neighbor search for spatial network databases," in *Proc. 13th Int. Conf. Very Large Data Bases*, 2004, pp. 840–851.
- [13] K. C. K. Lee, W.-C. Lee, B. Zheng, and Y. Tian, "Road: A new spatial object search framework for road networks," *IEEE Trans. Knowl. Data Eng.*, vol. 24, no. 3, pp. 547–560, Mar. 2012.
- [14] S. Luo, Y. Luo, S. Zhou, G. Cong, and J. Guan, "Disks: A system for distributed spatial group keyword search on road networks," *Proc. VLDB Endowment*, vol. 5, no. 12, pp. 1966–1969, 2012.
- [15] K. Mouratidis, M. L. Yiu, D. Papadias, and N. Mamoulis, "Continuous nearest neighbor monitoring in road networks," in *Proc. 13th Int. Conf. Very Large Data Bases*, 2006, pp. 43–54.
- [16] D. Papadias, J. Zhang, N. Mamoulis, and Y. Tao, "Query processing in spatial network databases," in *Proc. 13th Int. Conf. Very Large Data Bases*, 2003, pp. 802–813.
- [17] M. Qiao, L. Qin, H. Cheng, J. X. Yu, and W. Tian, "Top-k nearest keyword search on large graphs," *Proc. VLDB Endowment*, vol. 6, no. 10, pp. 901–912, 2013.



- [18] J. B. Rocha-Junior and K. Nøravåg, "Top-k spatial keyword queries on road networks," in *Proc. 15th Int. Conf. Extending Database Technol.*, 2012, pp. 168–179.
- [19] H. Samet, J. Sankaranarayanan, and H. Alborzi, "Scalable network distance browsing in spatial databases," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2008, pp. 43–54.
- [20] J. Sankaranarayanan, H. Samet, and H. Alborzi, "Path oracles for spatial networks," *Proc. VLDB Endowment*, vol. 2, no. 1, pp. 1210–1221, 2009.
- [21] J. Sankaranarayanan, H. Alborzi, and H. Samet, "Efficient query processing on spatial networks," in *Proc. 13th Annu. ACM Int. Workshop Geographic Inf. Syst.*, 2005, pp. 200–209.
- [22] H. Wang and R. Zimmermann, "Snapshot location-based query processing on moving objects in road networks," in *Proc. 16th ACM SIGSPATIAL Int. Conf. Geographic Inf. Syst.*, 2008, p. 50.
- [23] R. Zhong, G. Li, K.-L. Tan, and L. Zhou, "G-tree: An efficient index for knn search on road networks," in *Proc. 22nd Int. Conf. Inf. Knowl. Manage.*, 2013, pp. 39–48.
- [24] R. Zhong, G. Li, K.-L. Tan, and L. Zhou, "Processing of continuous location-based range queries on moving objects in road networks," *IEEE Trans. Knowl. Data Eng.*, vol. 23, no. 7, pp. 1065–1078, Jul. 2011.



**Ruicheng Zhong** is currently working toward the PhD degree in the Department of Computer Science, Tsinghua University, Beijing, China. His research interests mainly include spatio-textual data query, query processing on road network and spatial database.



**Guoliang Li** received the PhD degree from Tsinghua University, in 2009. He is currently an associate professor in the Department of Computer Science, Tsinghua University, Beijing. His research interests mainly include data cleaning and integration, spatial databases, and crowdsourcing.



**Kian-Lee Tan** received the BSc (Honors) and PhD degrees from the National University of Singapore (NUS) in 1989 and 1994. He is currently a professor in the Department of Computer Science, NUS. His major research interests include query processing and optimization, database security, and database performance. He is a member of the ACM.



**Lizhu Zhou** received the master's of science degree from the University of Toronto, in 1983. Currently, he is a full professor at the Tsinghua University, China. His major research interests include database systems, web data processing, and digital resource management. He is a member of the ACM.



**Zhiguo Gong** is currently an associate professor in the Department of Computer and Information Science, University of Macau, Macau, China. His research interests include databases, web information retrieval, and web mining.

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