



TASK: An Efficient Framework for Instant Error-tolerant Spatial Keyword Queries on Road Networks

Chengyang Luo
Zhejiang University
luocy1017@zju.edu.cn

Qing Liu
Zhejiang University
lqzju2010@gmail.com

Yunjun Gao
Zhejiang University
gaoyj@zju.edu.cn

Lu Chen
Zhejiang University
luchen@zju.edu.cn

Ziheng Wei
Huawei Cloud Computing
Technologies Co., Ltd
ziheng.wei@huawei.com

Congcong Ge
Huawei Cloud Computing
Technologies Co., Ltd
gecongcong1@huawei.com

ABSTRACT

Instant spatial keyword queries return the results as soon as users type in some characters instead of a complete keyword, which allow users to query the geo-textual data in a *type-as-you-search* manner. However, the existing methods of instant spatial keyword queries suffer from several limitations. For example, the existing methods do not consider the typographical errors of input keywords, and cannot be applied to the road networks. To overcome these limitations, in this paper, we propose a new query type, i.e., instant error-tolerant spatial keyword queries on road networks. To answer the queries efficiently, we present a framework, termed as TASK, which consists of index component, query component, and update component. In the index component, we design a novel index called reverse 2-hop label based trie, which seamlessly integrates spatial and textual information for each vertex of the road network. Based on our proposed index, we devise efficient algorithms to progressively return and update the query results in the query component and update component, respectively. Finally, we conduct extensive experiments on real-world road networks to evaluate the performance of our presented TASK. Empirical results show that our proposed index and algorithms are up to 1-2 orders of magnitude faster than the baseline.

PVLDB Reference Format:

Chengyang Luo, Qing Liu, Yunjun Gao, Lu Chen, Ziheng Wei, and Congcong Ge. TASK: An Efficient Framework for Instant Error-tolerant Spatial Keyword Queries on Road Networks. PVLDB, 16(10): 2418 - 2430, 2023.
doi:10.14778/3603581.3603584

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/ZJU-DAILY/TASK>.

1 INTRODUCTION

Geo-textual objects associated with both geographical and textual information are ubiquitous in daily life, such as restaurants, hotels,

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.
Proceedings of the VLDB Endowment, Vol. 16, No. 10 ISSN 2150-8097.
doi:10.14778/3603581.3603584

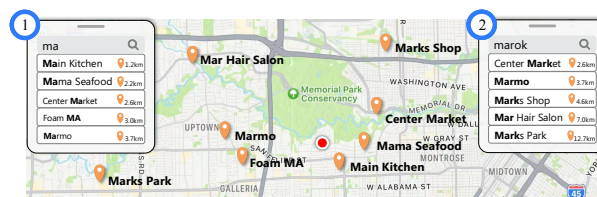


Figure 1: A motivating example

shopping malls, etc. In the literature, many types of spatial keyword queries have been studied [7, 11], e.g., top-*k* spatial keyword queries [12, 35], reverse spatial keyword queries [16, 27], why-not spatial keyword queries [6, 46], continuous spatial keyword queries [14, 37], to name but a few. In this paper, we focus on instant spatial keyword queries [19, 33, 48].

Instant queries, a.k.a., *search-as-you-type* or *type-ahead search*, return query results on-the-fly when users type in a query keyword character-by-character [3, 21–23]. They allow users to browse the results during typing characters. [19, 33, 48] investigate the instant queries for spatial databases. Specifically, with every character of keyword being typed in, the instant spatial keyword queries return geo-textual objects that are relevant to the query inputs. The instant spatial keyword queries save users from typing in complete keywords and thus are helpful for users in real-world applications. As an example, assume that a user wants to go to a restaurant called "Mama seafood" for dinner, and uses the location based service (LBS) to find the location. As shown in Figure 1, as soon as the user types in "Ma", the queries can return the results containing "Mama seafood" for her/him.

However, the state-of-the-art methods of instant spatial keyword queries [19, 33, 48] suffer from several limitations.

- The existing methods are mainly designed for the Euclidean space. For example, the bound-materialized trie proposed in [33] employs grid to index the space. Both the filtering-effective hybrid index [19] and prefix-region tree [48] leverage R-tree for the space indexing. All those techniques cannot be applied to road networks since they use different metrics to compute the distance between two objects. In real applications, it is straightforward and more practical to query the geo-textual data on road networks [1, 17, 20, 30, 32, 43]. Hence, it is necessary to develop techniques for instant spatial keyword queries on road networks.

Table 1: Taxonomy of representative related work and our work

Category	Index	Error tolerant	Multiple keywords	Search as you type	Character deletion	Non-tail operations	Road network
Instant spatial keyword queries	Bound-Materialized Trie [33]	✗	✗	✓	✗	✗	✗
	Filtering-effective hybrid index (FEH) [19]	✗	✗	✓	✗	✗	✗
	Prefix-region tree (PR-Tree) [48]	✗	✓	✓	✗	✗	✗
Spatial keywords queries over road networks	Map B-tree and inverted file [32]	✗	✗	✗	✗	✗	✓
	Compact tree [30]	✗	✗	✗	✗	✗	✓
	LB Index and KT Index [20]	✗	✓	✗	✗	✗	✓
	Keyword separated index (K-SPIN) [1]	✗	✓	✗	✗	✗	✓
Instant spatial keywords queries over road networks	Reverse 2-hop label based trie (Our work)	✓	✓	✓	✓	✓	✓

- The existing methods do not consider the typographical errors of input keywords. Sometimes, typing accurately is a tedious task, and the users' inputs tend to contain typographical errors, especially for mobile devices. Consequently, it is critical for the instant spatial keyword queries to tolerate typos [31, 34, 36, 38, 49], which can help users query in a friendly way. Thus, it motivates us to consider the error tolerance for the instant spatial keyword queries.
- The traditional instant spatial keyword queries mostly focus on the cases where users type in queries character-by-character. Nevertheless, some common yet important cases are ignored. For instance, users may (i) delete characters during the query, and (ii) add/delete characters anywhere in the query. Considering these cases can greatly enrich the instant spatial keyword queries.

To overcome these limitations, in this paper, we study a new problem, i.e., instant error-tolerant spatial keyword queries on road networks. Specifically, given a road network including geo-textual objects, a query location, and a query string that may have typographical errors, the instant error-tolerant spatial keyword queries return top- k geo-textual objects that are the most relevant to query location and query string. When the query string updates, e.g., the users input new characters to the query string or delete some characters from the query string, the queries should update the top- k geo-textual objects instantly. Our studied problem has wide applications in LBS. For example, in Figure 1, assume that the user wants to query "Marks shop", and starts typing in a string "Ma". Then, five geo-textual objects containing the prefix "Ma" are quickly returned. However, the desired "Marks shop" is not in the results. Next, the user proceeds to type in the string "Marok", where a typo is contained. The queries tolerate the typo, and return new results containing "Marks shop". It is worth mentioning that the traditional instant queries only deal with the cases of typing keywords character-by-character. Our work also considers the deletion of characters, which is beneficial in real applications.

In the literature, many indexes have been proposed to handle the instant spatial keyword queries and spatial keyword queries on road networks, as summarized in Table 1. Since those indexes do not take all the requirements of our problem into consideration, they cannot be applied to tackle our problem. To this end, we present a

novel index called `reverse 2-hop label based trie (R2T)` to answer the instant error-tolerant spatial keyword queries on road networks. R2T consists of two parts, i.e., `reverse 2-hop label` and `trie`. Specifically, the reverse 2-hop label is based on 2-hop labeling techniques [10], which enables efficient calculation of the distance between two vertices. With the help of the reverse 2-hop label, R2T is capable of efficient distance computation during the query processing. For the textual information, we employ trie that can efficiently support characters matching. The complete trie is complex and large, which is not efficient for queries since we need to traverse it multiple times. Thus, we design a novel structure called `node array` to store partial trie for each vertex with respect to the reverse 2-hop label. Moreover, we have discussed the index maintenance for dynamic road networks.

Based on R2T, we devise efficient algorithms to support instant error-tolerant spatial keyword queries on road networks, including `instant query algorithm` and `instant update algorithm`. The instant query algorithm aims to return the results when users type in characters for the first time. It first traverses the complete trie to get the active nodes, which contains the candidate results, and then, it visits the R2T to progressively find the geo-textual objects that are the most relevant to the query location and query string. The instant update algorithm is to update the query results according to the updated query string. To avoid querying from the scratch, we design `query information inheritance mechanism` to make full use of previous query processing information, which can dramatically improve the query performance. Furthermore, we extend our algorithms to support multiple query strings.

Our key contributions are summarized as follows:

- We identify and formalize the problem of instant error-tolerant spatial keyword queries on road networks, which is rooted in real-world applications. To the best of our knowledge, it is the first attempt to investigate this problem.
- We design a novel index called R2T, which seamlessly integrates the spatial and textual information for each vertex of the road network to facilitate queries. Efficient algorithms are also proposed to construct and maintain R2T.
- We present efficient algorithms to answer queries using R2T, which can return the results in a progressive way. In particular, the first type of algorithms focus on how to

retrieve results when users type in a query string for the first time. The second type of algorithms handle how to efficiently update results for the updated query string.

- We conduct extensive experiments on real-world road networks to demonstrate the efficiency of our proposed index and algorithms.

The rest of this paper is organized as follows. Section 2 reviews related work. Section 3 formally defines the problem. Section 4 introduces the framework TASK. Section 5 presents the structure, construction, and maintenance of R2T index. Section 6 proposes algorithms for queries. Experimental results are reported in Section 7. Finally, Section 8 concludes the paper.

2 RELATED WORK

Instant spatial keyword queries. Existing studies proposed different indexes to support instant spatial keyword queries [19, 33, 48]. Basu Roy and Chakrabarti [33] first introduced the instant queries to spatial database, and developed the index called materialized trie (MT). MT uses trie as the main index structure, and incorporates spatial information into the node of trie. Ji et al. [19] proposed an R-tree based method called Filtering-Effective Hybrid Indexing (FEH) for instant spatial keyword queries. FEH utilizes R-tree as the key index structure. In each R-tree node, FEH incorporates textual filters according to the geo-textual objects contained in this node. Zhong et al. [48] proposed the Prefix-region tree (PR-Tree), which considers spatial information and textual information in a balanced manner. In addition, as surveyed in [5], many indexes, such as IR-tree [26], have been proposed for spatial keyword queries. Nonetheless, all these indexes are designed for the Euclidean space, and cannot be directly used in our work.

The mentioned spatial keyword queries are related to location-aware autocompletion [18, 34]. Specifically, as users type in queries, the location-aware autocompletion returns the possible completion of the queries. However, location-aware autocompletion is not applicable to our problem for three reasons. Firstly, it uses Euclidean space, which is not suitable for road networks. Secondly, it considers spatial similarity and textual similarity separately, whereas our problem requires a comprehensive approach. Lastly, location-aware autocompletion requires querying from scratch whenever users input characters, resulting in time consumption.

Spatial keyword queries on road networks. Rocha-Junior and Nørnvåg [32] first explored the spatial keyword queries on road networks. An index consisting of map tree and inverted file is proposed to answer the queries. Then, Qiao et al. [30] devised an index structure based on distance oracles and compact trees of keywords. Jiang et al. [20] presented 2-hop label backward index (LB) and keyword-lookup tree index (KT). Abeywickrama et al. [1] designed the keyword separated index (K-SPIN), which includes a set of *on-demand inverted heap* of a keyword. Besides traditional spatial keyword queries on road networks, many variants have also been studied, such as aggregate queries [9], time-aware queries [43], continuous queries [47], why-not queries [44], diversified queries [39], and reverse queries [45]. All these studies need users to type in complete queries and hence cannot be applied in our work.

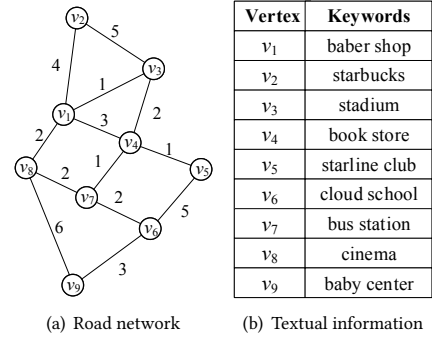


Figure 2: A running example

3 PROBLEM FORMULATION

In this section, we formalize the problem of instant error-tolerant spatial keyword queries on road networks.

The road network is denoted as a connected undirected weighted graph $G = (V, E)$. V and E are the sets of vertices and edges of G . The weight of an edge $e = (u, v)$, denoted by $w(e)$ or $w(u, v)$, is a positive integer, which denotes a metric, such as distance or travel time, between u and v . A path p is a sequence of vertices $p = (v_1, v_2, \dots, v_j)$, where $\forall 1 \leq i < j, (v_i, v_{i+1}) \in E$. The length of a path is the sum of weights of edges along the path. Given two vertices u and v of a road network G , the distance between u and v , denoted by $\text{dis}(u, v)$, is the shortest path between u and v in G . For example, in Figure 2, the length of path $p = (v_5, v_6, v_7)$ is $w(v_5, v_6) + w(v_6, v_7) = 5 + 2 = 7$; and the distance between v_5 and v_7 is $\text{dis}(v_5, v_7) = w(v_5, v_4) + w(v_4, v_7) = 2$.

The geo-textual objects and query location may appear on any point of the road network G . Given a geo-textual object and a query location on G , we can compute the distance between them by mapping each of them to an adjacent vertex with an offset. To make exposition simpler, we assume that the geo-textual objects and query location all appear on vertices, which follows previous studies such as [1, 20]. Under this assumption, in the road network G , a vertex v contains a set of keywords, which is denoted by $\text{doc}(v)$. We use the notation $kw \in \text{doc}(v)$ to denote that the vertex v includes the keyword kw . If a string str is the prefix of a keyword kw , we denote it by $str \leq kw$. For instance, in Figure 2, "bus" $\in \text{doc}(v_7)$ and "bu" \leq "bus". Next, we introduce the concepts of edit distance and prefix edit distance to measure the textual similarity.

DEFINITION 3.1. (Edit Distance, Prefix Edit Distance) Given a keyword kw , two strings str_1 and str_2 .

(1) The edit distance between str_1 and str_2 , denoted by $\text{ED}(str_1, str_2)$, is the minimum number of single-character edit operations, including insertion, deletion, and substitution, needed to transform str_1 to str_2 .

(2) The prefix edit distance between kw and str_1 , denoted by $\text{PED}(kw, str_1)$, is the minimum edit distance between kw 's prefix and str_1 , i.e., $\text{PED}(kw, str_1) = \min_{str' \leq kw} \text{ED}(str', str_1)$.

For example, $\text{ED}(\text{"school"}, \text{"scholar"}) = 3$, $\text{PED}(\text{"school"}, \text{"sco"}) = \text{ED}(\text{"sch"}, \text{"sco"}) = \text{ED}(\text{"sc"}, \text{"sco"}) = \text{ED}(\text{"scho"}, \text{"sco"}) = 1$. Based on the metric of textual similarity presented above, we formally define the error-tolerant spatial keyword queries on road networks.

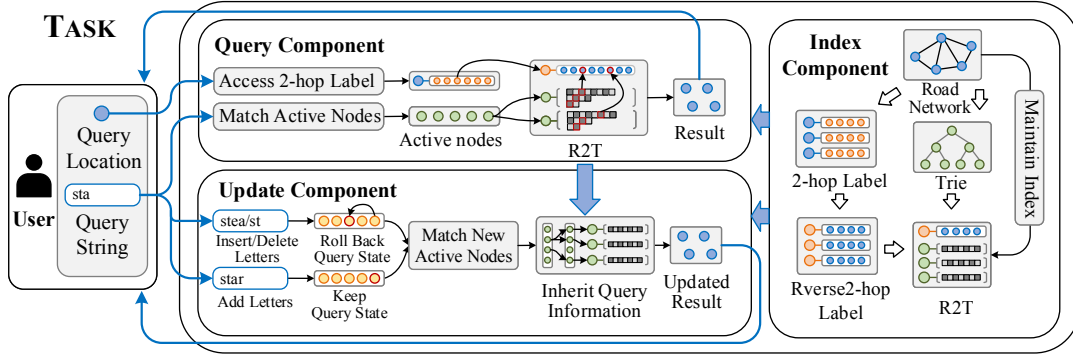


Figure 3: The workflow of TASK

DEFINITION 3.2. (Error-tolerant Spatial Keyword Queries on Road Networks) Given a road network $G = (V, E)$, a parameter k , an error threshold τ , a query $q = (q.loc, q.str)$, where $q.loc \in V$ and $q.str$ are the query location and query string, respectively. The error-tolerant spatial keyword queries on road networks return a set $\mathcal{R} \subseteq V$ such that:

- (1) $|\mathcal{R}| = k$;
- (2) $\forall v \in \mathcal{R}, \exists kw \in \text{doc}(v), \text{PED}(kw, q.str) \leq \tau$;
- (3) $\forall v \in \mathcal{R}$ and $\forall v' \in V - \mathcal{R}, \text{score}(q, v) \leq \text{score}(q, v')$, in which $\text{score}(q, v)$ computes the spatial and textual similarity between q and v . Specifically,

$$\text{score}(q, v) = \alpha \cdot \frac{\text{dis}(q.loc, v)}{\max_{u, u' \in V} \text{dis}(u, u')} + (1 - \alpha) \cdot \frac{\min_{kw \in \text{doc}(v)} \text{PED}(kw, q.str)}{\tau} \quad (1)$$

In Definition 3.2, (1) the first condition ensures that the returned results have at most k geo-textual objects; (2) the second condition gives the upper bound of the typographical error for the query string such that the keywords of each returned geo-textual object can match the query string; and (3) the third condition employs the normalized aggregated score of spatial similarity and textual similarity to rank the geo-textual objects such that the returned k geo-textual objects are optimal. Specifically, Equation 1 consists of two parts. The first part, i.e., $\frac{\text{dis}(q.loc, v)}{\max_{u, u' \in V} \text{dis}(u, u')}$, employs the shortest distance to measure the spatial similarity between the query location $q.loc$ and the geo-textual object v . The closer v is to $q.loc$, the more preferable v is to $q.loc$. The second part, i.e., $\frac{\min_{kw \in \text{doc}(v)} \text{PED}(kw, q.str)}{\tau}$, uses the prefix edit distance to measure the textual similarity between the query string $q.str$ and the keywords of v . The smaller prefix edit distance between $q.str$ and the keywords of v , the more similar they are. For equality, we employ the maximum shortest distance of the road network and the threshold τ to normalize these two parts, respectively. Moreover, the parameter $\alpha \in [0, 1]$ in Equation 1 is introduced to adjust the importance of spatial and textual similarities. In brief, the error-tolerant spatial keyword queries on road networks tolerate the input error, and return the most relevant geo-textual objects for users. Based on this, we formalize our studied problem below.

PROBLEM 1. Given a road network $G = (V, E)$, a parameter k , an error threshold τ , and a query $q = (q.loc, q.str)$, the problem of instant error-tolerant spatial keyword queries on road networks should (1) return the set \mathcal{R} satisfying the three conditions in Definition 3.2; and (2) update \mathcal{R} whenever the query string $q.str$ changes.

In a word, the goal of Problem 1 is two-folded. First, when a user types in a string, we should return the top- k geo-textual objects that best match user's location and the typed string. Second, when the typed string changes, e.g., a new character is inserted into the string or a character is deleted from the string, we should update the results instantly. For instance, in Figure 2, assume that $q.loc = v_1$, $k = 3$, $\tau = 1$, and $\alpha = 0.5$. When a user types in a query string "st", the results $\{v_3, v_4, v_2\}$ are returned. When the user proceeds to type in "a" after "st", i.e., the current query string is "sta", the results are updated to $\{v_3, v_2, v_5\}$ immediately.

4 FRAMEWORK OVERVIEW

Problem 1 returns k geo-textual objects with the minimum score, and the typo should be not larger than τ . A naive method is to traverse the road network from $q.loc$ in a breadth first manner. For the visited vertex, we check whether the prefix edit distance of the vertex's keywords is no larger than τ . If yes, we compute its score using Equation 1, and add it to the candidate set. Finally, k vertices with the minimum score are returned. When the query string changes, we can use the above naive method to compute the results from the scratch. Obviously, the simple combination of road network traversal and text examination leads to poor performance of the naive method.

Motivated by this, we propose an efficient framework, termed as TASK, to tackle the instant error-tolerant spatial keyword queries on road networks. As illustrated in Figure 3, TASK consists of an index component, a query component, and an update component. The index component is responsible for (1) constructing the R2T index for the road network and (2) maintaining R2T when the road network changes. The index component provides support for queries, and is the cornerstone of the framework. When a user types in a query string at the first time (i.e., when the current query string is empty), TASK uses the query component to return results to the user. Whenever the user makes change for the query string, TASK calls the update component, which employs the information of previous query processing to quickly update the results for user.

Table 2: 2-hop label and reverse 2-hop label of the road network in Figure 2

Vertex	2-hop label	Reverse 2-hop label
v_1	$(v_1, 0)$	$(v_1, 0), (v_3, 1), (v_8, 2), (v_4, 3), (v_2, 4), (v_5, 4), (v_7, 4), (v_6, 6), (v_9, 8)$
v_2	$(v_2, 0), (v_1, 4)$	$(v_2, 0)$
v_3	$(v_3, 0), (v_1, 1), (v_4, 2)$	$(v_3, 0)$
v_4	$(v_4, 0), (v_1, 3)$	$(v_4, 0), (v_5, 1), (v_7, 1), (v_3, 2), (v_6, 3), (v_8, 3), (v_9, 6)$
v_5	$(v_5, 0), (v_4, 1), (v_1, 4)$	$(v_5, 0)$
v_6	$(v_6, 0), (v_4, 3), (v_1, 6)$	$(v_6, 0), (v_7, 2), (v_9, 3), (v_8, 4)$
v_7	$(v_7, 0), (v_4, 1), (v_6, 2), (v_1, 4)$	$(v_7, 0), (v_8, 2)$
v_8	$(v_8, 0), (v_1, 2), (v_7, 2), (v_4, 3), (v_6, 4)$	$(v_8, 0), (v_9, 1)$
v_9	$(v_9, 0), (v_6, 3), (v_4, 6), (v_8, 6), (v_1, 8)$	$(v_9, 0)$

The user proceeds to type in the query string until the desirable results are found. In the following two sections, we will detail these three components. Note that, due to space limitation, some pseudocodes, examples, and proofs of Section 5 and Section 6 are moved to technical report [28].

5 R2T INDEX

In this section, we propose a novel index called reverse 2-hop label based trie (R2T for short) for Problem 1. We first introduce the structure of R2T, and then present the construction and maintenance algorithms of R2T.

5.1 R2T Structure

R2T integrates both road network information and textual information. Specifically, we employ the reverse 2-hop label and trie techniques to store spatial and textual information, respectively. First, we introduce the concept of 2-hop label [2, 8, 10, 15, 29].

Given a road network G , the 2-hop labeling technique assigns each vertex $v \in V$ a label $L(v)$ containing a set of pairs $(u, \text{dis}(u, v))$, i.e., $L(v) = \{(u, \text{dis}(u, v)) | u \in V\}$. The 2-hop labels of all vertices have the following property. Given any two vertices v and v' of the road network G , the distance between v and v' is $\text{dis}(v, v') = \min_{u \in L(v) \cap L(v')} \text{dis}(v, u) + \text{dis}(v', u)$. In other words, the distance between any two vertices can be computed via an intermediate hop. For example, Table 2 shows the 2-hop labels for all vertices of the road network G in Figure 2. We can observe that $L(v_2) = \{(v_2, 0), (v_1, 4)\}$ and $L(v_7) = \{(v_7, 0), (v_4, 1), (v_6, 2), (v_1, 4)\}$. The vertex v_1 is the only common label vertex in $L(v_2)$ and $L(v_7)$. Thus, $\text{dis}(v_2, v_7) = \text{dis}(v_2, v_1) + \text{dis}(v_1, v_7) = 4 + 4 = 8$. Based on the 2-hop label, the reverse 2-hop label is defined as follows.

DEFINITION 5.1. (Reverse 2-hop Label) *Given a road network G and 2-hop labels of all vertices in G , the reverse 2-hop label of a vertex v , denoted by $\tilde{L}(v)$, consists of pairs $(u, \text{dis}(u, v))$ with $(v, \text{dis}(u, v)) \in L(u)$, i.e., $\tilde{L}(v) = \{(u, \text{dis}(u, v)) | \forall u, (v, \text{dis}(u, v)) \in L(u)\}$.*

If a vertex v is included in the 2-hop label of u , the reverse 2-hop label of v contains u . For instance, in Table 2, $(v_1, 4) \in L(v_2)$ and $(v_2, 4) \in \tilde{L}(v_1)$. The reverse 2-hop label is mainly used for the

distance computation, which can benefit Problem 1 a lot. The 2-hop label only supports distance computation between two vertices. Given 2-hop labels of a road network, if we want to find the k nearest neighbors of a vertex v , we should compute the common 2-hop label vertices of v and every other vertex, which is time consuming. With the help of reverse 2-hop label, we only need to compute the common 2-hop label vertices of v and v 's k nearest neighbors, and thus improving the query efficiency. Note that, we assume that the pairs $(v, \text{dis}(u, v))$ in both 2-hop label and reverse 2-hop label are in *ascending order* w.r.t., $\text{dis}(u, v)$. Besides, when the context is clear, we use the notations $v \in L(u)$ and $v \in \tilde{L}(u)$ to denote $(v, \text{dis}(u, v)) \in L(u)$ and $(v, \text{dis}(u, v)) \in \tilde{L}(u)$ for simplicity.

Next, we introduce another technique used in R2T, i.e., trie. The trie is an ordered tree to represent a set of keywords. Its root node is an empty node. Each non-leaf node is labeled with one character, and each leaf node contains the last character of a keyword. The path from the root node to a leaf/intermediate node represents a keyword/prefix in the trie. As an example, Figure 4 depicts a trie representing the keywords set in Figure 2(b). The nodes n_4 and n_7 denote the prefix "bab" and keyword "baby", respectively.

Based on the reverse 2-hop label and trie, we formally define our proposed index as follows.

DEFINITION 5.2. (Reverse 2-hop Label Based Trie) *Given a road network G and a vertex $v \in V$, a reverse 2-hop label based trie of v , denoted by $\text{R2T}(v)$, consists of the reverse 2-hop label of v and the trie of v w.r.t. $\tilde{L}(v)$, i.e., $\text{R2T}(v) = (\tilde{L}(v), \mathcal{T}(v))$. In particular, $\mathcal{T}(v)$ is a trie of v to represent all the keywords contained in the vertices in $\tilde{L}(v)$, i.e., $\forall v' \in \tilde{L}(v), \forall kw \in \text{doc}(v')$.*

For example, in Table 2, $\tilde{L}(v_6) = \{(v_6, 0), (v_7, 2), (v_9, 3), (v_8, 4)\}$ and the keywords set w.r.t. $\tilde{L}(v_6)$ is {"cloud", "school", "bus", "station", "cinema", "baby", "center"}. Correspondingly, Figure 5 depicts $\mathcal{T}(v_6)$.

Combining Figures 4 and 5, we can find that each $\mathcal{T}(v)$ is a part of the complete trie. However, storing the entire $\mathcal{T}(v)$ directly in $\text{R2T}(v)$ has two drawbacks: (1) redundant storage and (2) repeated trie traversal. To address these issues, we use a node array to store $\mathcal{T}(v)$. The node array consists of a set of $(\text{ID}(n), \text{B}(n))$ for every leaf node n of $\mathcal{T}(v)$, where $\text{ID}(n)$ is the ID of n in the complete trie of the road network and $\text{B}(n)$ consists of bitmaps of trie nodes along the path from root node to n . Note that, the bitmaps of root node, leaf nodes, and the nodes whose bitmaps are the same as leaf nodes are not stored. In $\mathcal{T}(v)$, a node n 's bitmap is defined below: (1) Each bitmap has $|\tilde{L}(v)|$ bits, and each bit represents a vertex in $\tilde{L}(v)$. (2) A bit is set to "1" if the represented vertex has a keyword containing the prefix denoted by n , and it is set to "0" otherwise; the bits of the root node are all "1". We take the node n'_2 of $\mathcal{T}(v_6)$ in Figure 5 as an example. Since $|\tilde{L}(v_6)| = 4$, each bitmap of $\mathcal{T}(v_6)$'s node has 4 bits, representing v_6, v_7, v_9 , and v_8 , respectively. n'_2 denotes the prefix "b". According to Figure 2(b), "b" \leq "bus" $\in \text{doc}(v_7)$ and "b" \leq "baby" $\in \text{doc}(v_9)$. Thus, the bitmap of n'_2 is "0110".

In Figure 5, we can observe that there are a lot of "0" bits in bitmaps, which can be further compressed. Let b_1 and b_2 be the bitmaps of nodes n_1 and n_2 , respectively, and n_2 be the child node of n_1 . For the same bit of b_1 and b_2 , if both bits are "0", we delete the corresponding "0" bit from b_2 . Back to Figure 5, the bitmaps of n'_2 and n'_3 are "0110" and "0010", respectively. Since the first and last bits of n'_2 and n'_3 are all "0", we can delete the first and last bits

Algorithm 1: Instant Query Algorithm (IQA)

Input: a query $q = (q.loc, q.str)$, parameters k and τ , a trie T of G
Output: a set \mathcal{R} of k geo-textual objects

```
1  $\mathcal{R} \leftarrow \emptyset; C \leftarrow \emptyset;$   
2  $AN \leftarrow$  find active nodes of  $T$  for  $q.str$  [13];  
3 if  $AN = \emptyset$  then  
4   return  $\mathcal{R};$   
5 for  $\forall v \in L(q)$  do  
6   if  $T(v)$  contains the leaf nodes of  $AN$  then  
7      $v' \leftarrow \text{Min}\tilde{L}(v);$  // Computing the vertex with  
8     the minimal score using Algorithm 2  
9      $C \leftarrow C \cup (v', v, \text{score}(q, v'));$   
10  while  $|\mathcal{R}| < k \wedge C \neq \emptyset$  do  
11     $(v', v, \text{score}(q, v')) \leftarrow$   
12     $\arg \min_{(v', v, \text{score}(q, v')) \in C} \text{score}(q, v');$   
13     $C \leftarrow C - (v', v, \text{score}(q, v'));$   
14    if  $v' \notin \mathcal{R}$  then  
15       $\mathcal{R} \leftarrow \mathcal{R} \cup v';$   
16     $v'' \leftarrow \text{Min}\tilde{L}(v);$  // Computing the next vertex  
17    with the minimal score using Algorithm 2  
18     $C \leftarrow C \cup (v'', v, \text{score}(q, v''));$   
19 return  $\mathcal{R};$ 
```

first case, we first insert the keywords of u into the trie of v one by one. Correspondingly, we should insert the nodes representing the inserted keyword into $T(v)$. Then, we update the bitmaps of $T(v)$ by traversing the trie of v in a depth-first manner. If the visited trie node is newly inserted, we compute its bitmap and insert it into $T(v)$. Otherwise, we only need to update the existing bitmap as follows. (i) If u 's keywords contain the prefix represented by the visited trie node n , we add a "1" bit into n 's bitmap to represent u , and then visit n 's child node. (ii) If u 's keywords do not contain the prefix represented by the visited trie node n , we insert a "0" bit into n 's bitmap to represent u , and stop visiting n 's child node. After traversing the trie, all the bitmaps can be updated. (2) A pair $(u, \text{dis}(u, v))$ is deleted from $\tilde{L}(v)$. For the second case, we first delete all the nodes, which denote the unique keywords of u , and its bitmaps from $T(v)$. Similarly, by traversing the trie of v in a depth-first manner, we update the bitmaps of $T(v)$. (i) If the prefix represented by the visited trie node n is contained in u 's keywords, we should delete the "1" bit of u from n 's bitmap, and then visit n 's child node. (ii) If the prefix represented by the visited trie node n is not contained in u 's keywords, we should delete the "0" bit of u from n 's bitmap, and stop visiting n 's child node.

6 QUERY PROCESSING ALGORITHMS

Using R2T index, we propose the query processing algorithms.

6.1 Instant Query Algorithm

First, in this section, we present an algorithm, called Instant Query Algorithm (IQA), to handle the first case of Problem 1, i.e., when a user types in a query string for the first time, the query returns the top- k geo-textual objects with the minimal scores. With the help of the R2T index introduced in the previous section, the query

can be efficiently handled. First, we propose some lemmas, which establish the solid base to design IQA.

LEMMA 6.1. Given a vertex $v \in G$,

$$\bigcup_{v' \in \tilde{L}(v)} \tilde{L}(v') = V \quad (2)$$

Lemma 6.1 shows that, for a vertex v , all the reverse 2-hop labels $\tilde{L}(v')$ of $v' \in L(v)$ constitute the vertex set of the road network. For example, in Table 2, $L(v_2) = \{(v_2, 0), (v_1, 4)\}$ and $\tilde{L}(v_2) \cup \tilde{L}(v_1) = V$. Based on Lemma 6.1, we can compute the distance from a vertex to all the other vertices by using the reverse 2-hop labels. Moreover, the reverse 2-hop label also enable the computation of the nearest neighbors in a progressive way as shown in the following lemma.

LEMMA 6.2. Given a vertex $v \in G$, for a vertex $v' \in L(v)$, let $v'' = \arg \min_{v'' \in \tilde{L}(v')} \text{dis}(v', v'')$. Note that $v'' \neq v$. Then, the nearest neighbor of v is

$$v'' = \arg \min_{v'' \in \tilde{L}(v'), v' \in L(v)} \text{dis}(v, v') + \text{dis}(v', v'') \quad (3)$$

According to Lemma 6.2, in each reverse 2-hop label of $v' \in L(v)$, we can find the nearest neighbor of v' . Then, among all reverse 2-hop labels, the vertex with the minimal distance to v is the nearest neighbor of v . For instance, in Table 2, $L(v_4) = \{v_4, v_1\}$. In $\tilde{L}(v_4)$, the closest vertex to v_4 is v_5 with $\text{dis}(v_4, v_5) = 1$. In $\tilde{L}(v_1)$, the nearest vertex to v_1 is v_1 with $\text{dis}(v_1, v_1) = 0$. Since $\text{dis}(v_4, v_5) = 1$ and $\text{dis}(v_4, v_1) = 3$, v_5 is the nearest neighbor of v_4 . Lemma 6.2 can also be extended to our problem as follows.

COROLLARY 6.1. Given a query $q = (q.loc, q.str)$, for a vertex $v \in L(q)$, let $v' = \arg \min_{v' \in \tilde{L}(v)} \text{score}(q, v')$. Then, the vertex with the minimal score w.r.t. q is

$$v' = \arg \min_{v' \in \tilde{L}(v), v \in L(q)} \text{score}(q, v') \quad (4)$$

In other words, the minimal score vertex is among the vertices, which have the minimal score in the reverse 2-hop label of $v' \in L(v)$. Motivated by Corollary 6.1, we develop IQA to find the k vertices with the minimal score. The basic idea of IQA is to progressively return the vertices with the minimal score among all reverse 2-hop labels of $v \in L(q)$. To this end, we should find the vertex with the minimal score for each reverse 2-hop label of $v \in L(q)$. Then, in each round, we select the vertex with the minimal score among all reverse 2-hop labels of $v \in L(q)$, and add it to the results. Assume the vertex with the minimal score is in $\tilde{L}(v)$. After this, we should compute the next vertex with the minimal score in $\tilde{L}(v)$ for the next round processing. IQA repeats the selection of vertices having the minimal scores until all results are found.

Algorithm 1 shows the pseudo-code of IQA. First, IQA computes the active nodes in the complete trie of the road network (line 2). Here, the active node is the node whose represented string/keyword's edit distance to $q.str$ is not larger than τ . It is worth mentioning that, the active nodes are mainly used to compute the score of the vertex, which will be illustrated later. If there is no such active node, it means that all vertices do not match the $q.str$. IQA returns an empty result (lines 3-4). Otherwise, IQA finds the vertex with the

Algorithm 2: $\text{Min}\tilde{L}(v)$

Input: a set AN of active nodes, $\text{R2T}(v)$ index of a vertex v
Output: a vertex v' with the minimal score in $\tilde{L}(v)$

```
1  $v' \leftarrow \emptyset$ ;  $\text{MinScore} \leftarrow +\infty$ ;  
2 if  $\text{BTag}[v] = \emptyset$  then  
3    $\_$  initialize  $\text{BTag}[v]$ ;  
4 for each active node  $an \in AN$  do  
5   find a leaf node  $n$  of  $an$  using binary search;  
6   for  $i \leftarrow an.depth$  to 1 do  
7     if  $i = an.depth$  then  
8        $\_$   $\text{BTag}[v][n][i].y = \text{BTag}[v][n][i].y + 1$ ;  
9     else  
10       $\_$   $\text{BTag}[v][n][i].y = \text{BTag}[v][n][i+1].x$ ;  
11      $\text{BTag}[v][n][i].x \leftarrow$  the position of  
12      $\text{BTag}[v][n][i].y$ -th "1" bit in bitmap  $\text{B}[v][n][i]$ ;  
13     if  $\text{BTag}[v][n][i].x = \text{null}$  then  
14        $\_$  break;  
15      $i - -$ ;  
16  $v'' \leftarrow (\text{BTag}[v][n][i].x)$ -th vertex in  $\tilde{L}(v)$ ;  
17 if  $\text{score}(q, v'') < \text{MinScore}$  then  
18    $v' \leftarrow v''$ ;  
19    $\text{MinScore} \leftarrow \text{score}(q, v'')$ ;  
20 return  $v'$ ;
```

minimal score for each reverse 2-hop label of $v \in L(q)$ (lines 5-8). Then, IQA repeatedly selects the vertex with the minimal score among all reverse 2-hop labels until all results are found (lines 9-15). Finally, IQA returns the top- k geo-textual objects (line 16).

Next, we introduce how to find the vertex with the minimal score in a reverse 2-hop label using R2T index. Note that, we have to keep in mind that the vertex with the minimal score should satisfy the prefix edit distance constraint in Definition 3.2. Recall that, at the beginning of Algorithm 1, we have computed the active nodes in the complete trie. The edit distances between the prefixes represented by those active nodes and $q.str$ are no larger than τ . Hence, if a vertex's keyword contains the prefix denoted by an active node, the vertex is a candidate of the minimal score vertex. Thus, a naive method is to compute the score for all candidates in the reverse 2-hop label. However, this naive method is not progressive. In the sequel, we devise a method to progressively return the minimal score vertex in a reverse 2-hop label.

In the bitmap stored in $\text{T}(v)$ of $\text{R2T}(v)$, each bit represents whether the keyword of a corresponding vertex contains the prefix represented by a node in the trie. (1) All vertices with the same prefix in the trie node have the same textual similarity. (2) In the reverse 2-hop label, the vertices are in ascending order of the distances. Based on the above two facts, for an active node, the vertex with the minimal score is just the vertex represented by the first "1" bit in the corresponding bitmap of the active node. In the same way, the vertex denoted by the second "1" bit in the bitmap of the active node is the second minimal score vertex. Using this property of the active node, for a reverse 2-hop label, we can compute the minimal score vertex for each active node. The vertex having the smallest score among all active nodes is the minimal score vertex in the reverse 2-hop label.

Next, we introduce how to find the minimal score vertex for an active node. Finding the minimal score vertex for an active node is only to find the first "1" bit in the bitmap of the active node. A straightforward way is to convert the compressed bitmap to the full bitmap, and we can quickly find the vertex of the first "1" bit. However, this method is inefficient since it has to convert the compressed bitmaps from the active to the root node. In view of this, we devise a novel method via traversing the compressed bitmaps from the active node to the root node with the help of an auxiliary structure BTag , which is defined as follows.

DEFINITION 6.1. A BTag of a bitmap is a pair of $\langle x, y \rangle$, which denotes that the position of y -th "1" bit is x .

$\langle x, y \rangle$ of a bitmap means that from the first bit to the x -th bit, there are y "1" bits. For example, let a bitmap be "110011", the $\text{BTag}(5, 3)$ indicates that the fifth bit is the 3-th "1" bit in "110011". The BTag has the following property.

LEMMA 6.3. Given two trie nodes n_1 and n_2 , two pairs $\langle x_1, y_1 \rangle$ and $\langle x_2, y_2 \rangle$, which are the BTags of the bitmaps of n_1 and n_2 , respectively. Assume that n_1 is a child node of n_2 . If $y_2 = x_1$, the vertex represented by the x_1 -th bit in the bitmap of n_1 and the vertex denoted by the x_2 -th bit in the bitmap of n_2 are the same.

For instance, in Figure 5, the compressed bitmaps of n'_2 and n'_6 are "0110" and "10", respectively. Both $\langle 1, 1 \rangle$ of "10" and $\langle 2, 1 \rangle$ of "0110" represent the second vertex in $\tilde{L}(v_6)$, i.e., v_7 . Lemma 6.3 not only enables us to traverse the compressed bitmaps in a bottom-up manner to get the minimal score vertex for an active node, but also allows us to get the next minimal score vertex in the same way.

Based on the above discussion, we propose an algorithm to find the vertex with the minimal score in a reverse 2-hop label, whose pseudo-code is depicted in Algorithm 2. Initially, when Algorithm 2 computes the minimal score vertex for the first time, it initializes BTag (lines 2-3), i.e., to set BTag as $\langle 0, 0 \rangle$. Then, Algorithm 2 computes the minimal score vertex for each active node (lines 4-18). For each active node, it first gets a leaf node containing the bitmap of the active node using binary search (line 5). Note that, there may be multiple such leaf nodes. We only select any one of them. Next, Algorithm 2 traverses the bitmaps from the active node to the child of the root node for getting the position of the minimal score vertex of the active node (lines 6-14). Finally, the minimal score vertices of all active nodes are found, and the one having the smallest score among all active nodes is returned (lines 15-19).

Optimizations. We present two optimizations to speed up Algorithms 1 and 2.

Optimization 1. Algorithm 1 computes the minimal score vertex for each reverse 2-hop label (lines 5-8) for initialization. To improve the efficiency, we can first compute the lower bound of the score for a reverse 2-hop label. If the lower bound is larger than the current minimal score, the reverse 2-hop label definitely does not exist the minimum score vertex among all reverse 2-hop labels. Hence, we can skip the minimal score vertex computation for this reverse 2-hop label. Specifically, we can use Equation 1 to compute the lower bound of the score for the reverse 2-hop label $\tilde{L}(v)$. The spatial similarity is the minimal distance between q and the vertex in $\tilde{L}(v)$, and the textual similarity is the minimal edit distance between $q.str$ and the strings represented by active nodes.

Optimizations 2. Algorithm 2 needs to compute the minimal score vertex for every active node (lines 4-18). If the number of active nodes is large, it is costly. Actually, many active nodes have a parent-child relationship, and the vertices denoted by the bitmap of a child active node is a subset of that of its father active node. Thus, we need to only compute the minimal score vertex for father active nodes.

The time and space complexities of IQA are showed below.

THEOREM 6.4. *The time and space complexities of IQA are $O((k + w \cdot \log |V|) \cdot |AN| \cdot \sqrt{w \cdot \log |V|})$ and $O(w \cdot \log |V| \cdot |AN| \cdot |q.str|)$, respectively. $|AN|$ denotes the number of active nodes.*

6.2 Instant Update Algorithms

In this section, we present how to update the query results when the query string changes. A straightforward method is to query from the scratch. However, this method is not efficient, especially when updates are frequent. To this end, we extend Algorithms 1 and 2 to update the query results. We only need to modify three places in Algorithms 1 and 2. First, at the initialization step, Algorithm 1 should compute the first minimal score vertex for each reverse 2-hop label (lines 5-8). For the update algorithm, we can reuse the minimal score vertex found in the previous query, and the initialization can be omitted. Second, Algorithm 2 computes the minimal score vertex by traversing the bitmaps of active nodes and their ancestor nodes from the first bit, during which BTags are used. For the update algorithm, we can reuse the visited vertices, and traverse the bitmap from the current bit. Third, we have to recompute the score of previous query results based on the new query string to prune the unqualified vertices. Next, we detail the second modification for three cases as follows.

Case I: Inserting character(s) at the end of the query string. If the query string changes, the active nodes change accordingly. Let AN and AN' be the active nodes before and after inserting character(s) at the end of the query string. Then, it satisfies that $\forall an' \in AN'$, we can find an active $an \in AN$ such that $an = an'$ or an' is a descender node of an . If $an = an'$, we can still use the bitmap and corresponding BTag of an for an' to find the minimal score vertex. If an' is a child node of an , we should initialize the BTag of an' through an . Specifically, let $\langle x', y \rangle$ and $\langle x, y \rangle$ be the BTag of an' and an , respectively. Then, x' can be set to y . In the same way, we can iteratively derive the BTag for the bitmap of a descender node.

Case II: Deleting character(s) at the end of the query string. In order to achieve a fast update of deleting character(s) at the end, BTags for each insertion (Case I) need to be saved. BTags will be rolled back to the state corresponding to the string after the letter is deleted, and AN will also be rolled back. Note that if the BTags needed are not saved, then roll back to further back. The letters that pass at the end can be treated as Case I.

Case III: Inserting/deleting character(s) at random position of the query string. Let AN and AN' be the active nodes before and after deleting character(s) at random position of the query string. We rolled back BTags and AN like Case II to the position. Then, for each active nodes in AN' , it satisfies that $\forall an' \in AN'$, we can find an active $an \in AN$ such that $an = an'$ or an' is a descender node of an . This relationship is consistent with that in Case I, so we can update BTags with the method in Case I.

Table 3: Statistics of road networks

Dataset	Region	$ V $	$ E $	$ \text{doc}(V) $	$ W $
NY	New York City	264,346	733,846	157,100	6,556
FLA	Florida	1,070,376	2,712,798	343,452	16,656
CAL	California	1,890,815	4,657,742	401,258	20,319
LKS	Great Lakes	2,758,119	6,885,658	615,168	25,807
EU	Eastern USA	3,598,623	8,778,114	780,749	33,522
WU	Western USA	6,262,104	15,248,146	1,580,430	52,316
CTR	Central USA	14,081,816	34,292,496	2,782,249	78,658
USA	Full USA	23,947,347	58,333,344	4,118,452	112,353

6.3 Multiple Query Strings

Sections 6.1 and 6.2 mainly aim at a single query string. In this section, we discuss how to handle multiple query strings. After the user types in a query string, if the result does not meet his/her expectation, he/she may proceed to type in more query strings. The methods proposed in Sections 6.1 and 6.2 can also be extended to tackle multiple query strings. In particular, we discuss the extension of Algorithms 1 and 2 for the case of multiple query strings.

When a new query string $q.str'$ is added, Algorithm 1 should find new active nodes for $q.str'$ (line 2). Then, it iteratively computes the vertex with the minimal score to find the top- k results by using Algorithm 2 (lines 5-15), which is the same as the case of a single query string. Specifically, Algorithm 2 computes the vertex with the minimal score for each active node (lines 4-18). For a single query string, the first "1" bit in the bitmap of active node is the minimal score vertex. But, for multiple query strings, we should compute the scores of vertices denoted by "1" bits in the bitmap of active node until the lower bound of vertex's score is larger than the current minimum score. The lower bound of vertex's score can be computed using Equation 1, where the textual similarity is the sum of the minimal prefix edit distance for all query strings. In addition, if a vertex's prefix edit distance w.r.t. previously entered query strings is larger than τ , it cannot become the final result and thus can be skipped the score computation.

7 PERFORMANCE STUDY

This section evaluates the performance of our proposed index and algorithms. All algorithms were implemented in C++, and compiled by GCC 7.5.0 with $-O3$ optimization. The experiments were conducted on a machine running on Ubuntu server 18.04.5 LTS version with two Intel Xeon 2.40GHZ processors and 512G main memory.

Datasets: We employ eight real-world road networks in experiments. The road networks are obtained from DIMACS¹, which do not contain POIs. For each road network, we get POIs in the corresponding area from OpenStreetMap (OSM)² and map keywords of each POI to the closest vertex. Table 3 lists the statistics of the road networks, where $|\text{doc}(V)| = \sum_{v \in V} |\text{doc}(v)|$, and $|W|$ represents the number of distinct keywords in the road network.

Compared Methods: In experiments, the competitors include the naive method, the keyword query method, the instant query algorithm (IQA), and the instant update algorithm (IUA). Recall that, the first step of the naive method is to traverse the road network

¹<http://www.diag.uniroma1.it/~challenge9/download.shtml>

²<https://www.openstreetmap.org/>

from $q.loc$. We implemented three versions of the naive method by using different techniques to traverse the road network, including Dijkstra, G^* -tree, and 2-hop label. Note that, for G^* -tree, we employ the default parameter settings as reported in [25]. Also, we extend two state-of-the-art spatial keywords query methods, i.e., K-SPIN [1] and KT [20], to handle our problem. Specifically, first, we find all the complete keywords, whose prefix edit distance is not larger than τ . These keywords are used as candidate keywords. Then, we use K-SPIN and KT to find the k geo-textual objects with the minimum score.

Metrics: We report the query time, index construction/update time, and index size in our experiments. For the evaluation of query efficiency, we randomly generate 5000 queries and report the average query time. Due to space limitations and similar trends across different datasets, we present empirical results of partial datasets in this paper. The complete empirical results can be found in Appendix O of technical report [28].

7.1 Evaluation of Instant Query Algorithm

In this section, we study the efficiency of instant query algorithm.

Exp-1: Effect of $|q.str|$. We first verify the effect of the query string length $|q.str|$. We varied $|q.str|$ from 1 to 7 while keeping other parameters at their default values. Figure 7(a) shows the experimental results. We can observe that, with the growth of $|q.str|$, the query time of K-SPIN and KT drops; the query time of G^* -tree, 2-hop, and Dijkstra increases; and the query time of IQA ascends as $|q.str| \leq 3$. When $|q.str| > 3$, the query time of IQA almost does not change. For K-SPIN and KT, the longer $q.str$ is, the less the candidate keywords, resulting less query time. For G^* -tree, 2-hop, and Dijkstra, if $q.str$ becomes longer, there are less vertices satisfying the textual constraints. Thus, the search spaces of G^* -tree, 2-hop, and Dijkstra become larger, incurring more query time. For IQA, the number of active nodes increases with the growth of $|q.str|$ when $|q.str| \leq 3$. Since we set the default value of error threshold to 2, the number of active nodes goes down gradually when $|q.str| > 3$. Hence, the impact of the increased search space is offset, the running time of IQA almost keep the same as $|q.str| > 3$. Moreover, IQA outperforms other algorithms by 1-2 orders of magnitude.

Exp-2: Effect of $q.str$'s frequency. Next, we evaluate the effect of $q.str$'s frequency. Here, the $q.str$'s frequency is $\frac{|V(q.str)|}{|doc(V)|}$, where $V(q.str) = \{v | \forall v \in V, \exists kw \in doc(v), q.str \leq kw\}$. The query time are shown in Figure 7(b). The query time of all algorithms decrease with the increase of $q.str$'s frequency. The reason behind is that if $q.str$ is frequent in the road network, there are more vertices containing $q.str$. Thus, the results are easier to be found, resulting in less query time.

Exp-3: Effect of # of $q.str$. In this experiment, we explore the effect of the number of query strings, whose value is varied from 1 to 5. The results in Figure 7(c) show that the query time of IQA remained stable, while the query time of other algorithms increased as the number of query strings grew. IQA efficiently finds minimal score vertices by traversing active node bitmaps. Despite the increase in active nodes with more query strings, IQA can still find results by accessing only a few active nodes, resulting in consistent query times. Conversely, the other algorithms require traversing more

vertices to find results as the number of query strings increases, leading to longer query times.

Exp-4: Effect of $q.str$'s edit distance. Recall that our proposed algorithms can tolerate the typos in $q.str$. Thus, we study the effect of $q.str$'s edit distance. To this end, we assume that $q.str$ includes several typos. The $q.str$'s edit distance is the edit distance between $q.str$ and the correct string. Note that, to ensure non-empty query result, we set $\tau = 4$ in this experiment. Figure 7(d) shows the empirical results. We can observe that when the $q.str$'s edit distance increases, the query time of IQA, K-SPIN, and KT gradually drops while the query time of other algorithms gradually grows. This is because the larger the $q.str$'s edit distance, the less vertices in the road network satisfy the prefix edit distance constraint, meaning that G^* -tree, 2-hop, and Dijkstra need to search more vertices. For K-SPIN and KT, the larger the editing distance, the fewer candidate keywords, and therefore the shorter the query time. IQA can quickly skip invalid vertices since the larger $q.str$'s edit distance leads to less active nodes.

Exp-5: Effect of k . We investigate the effect of k by varying k from 1 to 64. Figure 7(e) plots the query time of four algorithms. As k becomes larger, the query time of all algorithms increase. The reason behind is that the larger k indicates more results. Hence, all algorithms take more time to query. However, the performance of IQA is still much better than that of other algorithms.

Exp-6: Effect of τ . Then, we evaluate the effect of error threshold τ on algorithms. The empirical results are reported in Figure 7(f). We can observe that the query time of G^* -tree, 2-hop, and Dijkstra decrease while the query time of IQA, K-SPIN, and KT increase with the growth of τ . If τ becomes larger, (1) more vertices satisfy the textual constraint, and (2) more active nodes and candidate keywords will be found. Thus, the G^* -tree, 2-hop, and Dijkstra spend less time while IQA, K-SPIN, and KT take more time.

Exp-7: Effect of α . α represents the user preference for score computation in Equation 1. This set of experiment test the effect of α on algorithms and the results are shown in Figure 7(g). We can observe that when $\alpha = 0$, other algorithms take more time. This is because the text pruning abilities of other algorithms are weak. For other values of α , the performance of all algorithms are stable, meaning that α has little effect on algorithms.

Exp-8: Scalability. In this set of experiments, we verify the scalability of the algorithms. In view of this, we vary the vertex cardinality $|V|$, the number of distinct keywords $|W|$, and the occurrences of keywords $|doc(V)|$. Figures 7(h), 7(i), and 7(j) plot the query time by changing $|V|$, $|W|$, and $|doc(V)|$, respectively. In Figure 7(h), with the growth of vertex cardinality, the performance of all algorithms degrade since the larger road network needs more time to find the results. In Figure 7(i), as the number of distinct keywords increases, the performance of four algorithms degrade as well. This is because when the number of distinct keywords grows, the keywords will become less frequent. Hence, all algorithms take more time with the growth of $|W|$, which is consistent with the empirical results of $q.str$'s frequency depicted in Figure 7(b). In Figure 7(j), as the occurrences of keywords increase, the performance of IQA keeps stable while that of other algorithms all degrade. If the occurrences of keywords grow, the vertices contain more keywords. Thus, the search

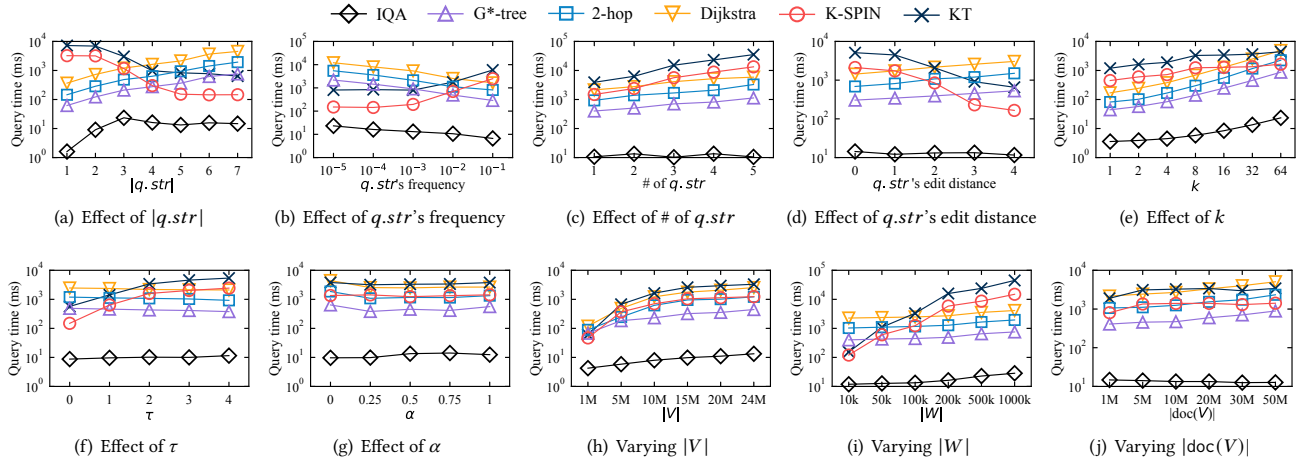


Figure 7: Evaluation of instant query algorithm on USA

spaces of other algorithms become larger, incurring more query time. On the other hand, the growth of the occurrences of keywords does not affect the complete trie of the road network. Hence, the active nodes found by IQA do not change as well. Therefore, the query time of IQA keeps stable.

7.2 Evaluation of Instant Update Algorithm

This section evaluates the performance of instant update algorithm. IQA, G*-tree, Dijkstra, 2-hop label, K-SPIN, and KT are instant query algorithms used for results update, i.e., to query from the scratch. IUA is the instant update algorithm presented in Section 6.2.

Exp-9: Effect of the position to insert characters. First, we explore the effect of inserting characters into different positions of $q.str$. Here, if we insert the characters into the position i of $q.str$, it means that we insert the characters after the i -th character of $q.str$. Note that, the $|q.str| \geq 7$, and we vary the position from 1 to 7. We insert one character in this experiment. Figure 8(a) depicts the empirical results. Specifically, the query time of all algorithms almost remain the same, except for IUA. Moreover, we can observe that if the insertion position is closer to the end of $q.str$, IUA has better performance. This is because the closer to the end of the query string, the less BTags will be updated. Thus, IUA needs less time to update. In addition, IUA has better performance compared with other algorithms. Specifically, IUA returns results in an average time of 2.1ms while IQA takes 10ms.

Exp-10: Effect of # of inserted characters. Next, we investigate the effect of inserting different number of characters into $q.str$. To this end, we extract characters from a complete keyword, and take the remaining string as $q.str$. Then, in experiments, we insert the extracted characters into $q.str$. In this way, we ensure the query result is non-empty. The query time of all algorithms are shown in Figure 8(b). When the number of inserted characters increases, IQA, G*-tree, Dijkstra, and 2-hop label take more query time. This is because when we insert more characters, $q.str$ becomes more accurate. They should traverse more vertices to find the results, and thus need more time to query. For IUA, if we insert more characters,

the difference between the original query string and new query string is greater. Hence, IUA has to need more time to update BTags. Although IUA becomes less efficient when more characters are inserted, it is still better than IQA, and is 2 orders of magnitude faster than other algorithms.

Exp-11: Effect of the position to delete characters. In this experiment, we study the instant update algorithm by deleting characters from different positions of the query string. We vary the deletion position from 1 to 7. Here, the deletion position i means that the i -th character of $q.str$ is deleted. Figure 8(c) depicts the empirical results. As observed, the closer the deletion position is to the end of $q.str$, the better performance of IUA. The reason behind is similar with that of Exp-9, i.e., when deleting characters near the end of $q.str$, less BTags need to be updated.

Exp-12: Effect of # of deleted characters. Then, we verify the effect of deleting different number of characters from $q.str$. In view of this, we randomly choose a deletion position in $q.str$, and delete a certain number of characters, which varies from 1 to 7. Empirical results are plotted in Figure 8(d). We have the following observations. With the growth of the number of deleted characters, (1) the query time of G*-tree, Dijkstra, and 2-hop label drops while that of K-SPIN and KT increases, (2) the query time of IQA remains unchanged, and (3) the query time of IUA increases slowly. However, IUA is still able to handle the deletion of characters efficiently, and outperforms other algorithms by 1-2 orders of magnitude.

Exp-13: Instant query simulation. This experiment simulates the users' instant queries, i.e., to simulate users type in the query string character-by-character. As soon as a character is typed in, we perform the query/update algorithms for the new query string, and return the results. Finally, after the users type in the complete query string, we report the total query/update time, as depicted in Figure 8(e). Note that, the length of query string for this experiment changes from 1 to 7. As expected, the total query time of all algorithms increase if the length of query string becomes larger. Nevertheless, the total query time of other algorithms ascends much faster than that of both IUA and IQA. This is because

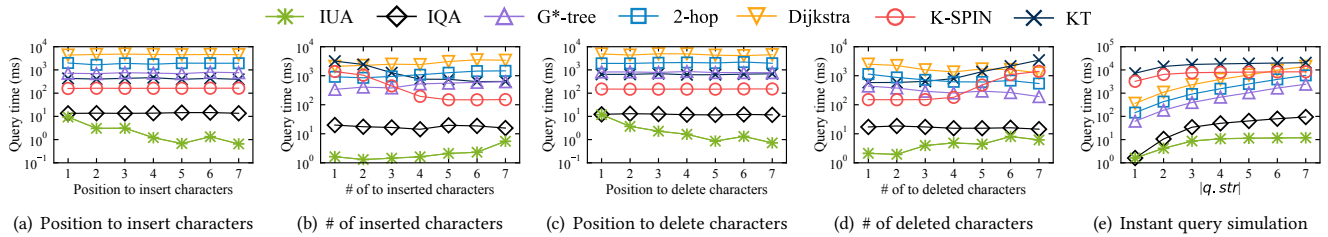
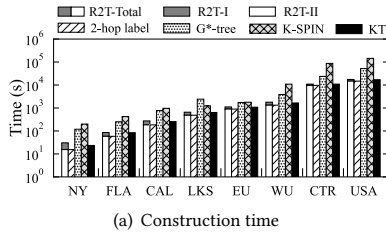
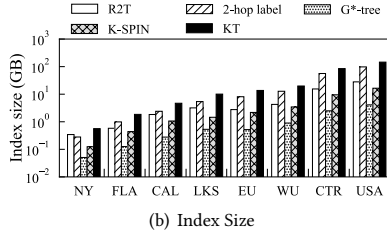


Figure 8: Evaluation of instant update algorithm on USA

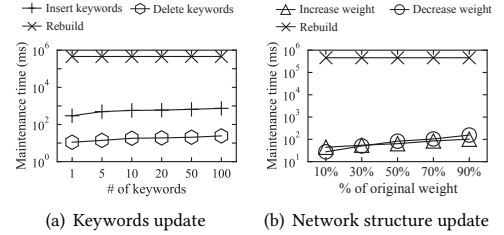


(a) Construction time



(b) Index Size

Figure 9: Index construction



(a) Keywords update

(b) Network structure update

Figure 10: Index maintenance (WU)

the query information inheritance mechanism of IUA makes it be able to perform update incrementally. In addition, the performance of IUA is much better than other algorithms. This is because IUA leverages a query information inheritance mechanism that enables the efficient utilization of information from previous queries. As a result, IUA can incrementally update the results, and significantly reduce query time compared with other algorithms.

7.3 Index Evaluation

In this section, we evaluate the performance of R2T index, including the index construction and maintenance.

Exp-14: Index Construction. First, we investigate the performance of index construction. In this experiment, we take G*-tree, 2-hop label, K-SPIN, and KT, which are the indexes of baselines, as competitors. Figures 9(a) and 9(b) show the index construction time and index size, respectively. Note that, in Figure 9(a), we split the construction time of R2T into two parts, i.e., R2T-I and R2T-II. Specifically, R2T-I represents the time of computing 2-hop label for every vertex, and R2T-II denotes the time of constructing $\tilde{L}(v)$ and $T(v)$ for every vertex after getting 2-hop labels. In Figure 9(a), we can observe that R2T-I takes up most of the construction time. For example, on the road network CTR, the time of R2T-I and R2T-II are 1289s and 9626.5s, respectively. R2T-I is 11.8% of the total construction time. Overall, the construction time of R2T is less than or comparable with other indexes except the 2-hop label. It is obvious since the construction time of R2T including the $\tilde{L}(v)$ and $T(v)$ construction time. In Figure 9(b), G*-tree has the smallest size and KT has the largest size. The size of R2T is smaller than that of the 2-hop label in most cases. The reason is two-fold. (1) In the road network, the vertices without keywords cannot contribute to the final results. We do not store these vertices in R2T. (2) The compression technique significantly reduces the bitmap size.

Exp-15: Index maintenance. Next, we explore the performance of index maintenance. We mainly consider two categories of the road network update, including keywords update and road network structure update. The keywords update is to insert/delete a certain number of keywords into/from $doc(v)$ of a vertex v . Figure 10(a) plots the index maintenance time for the keywords update. The road network structure update contains the update of edges' weight and the insertion/deletion of edges/vertices. Since the insertion/deletion of edges/vertices can be reduced to the update of edges' weight [42], we mainly consider the update of edges' weight in this experiment, which include the cases of increasing and decreasing the weight of an edge. For each road network, we select 1000 edges at random and change their weights. The maintenance time of road network structure update is shown in Figure 10(b). Overall, the maintenance time is much smaller than the time of rebuilding the index.

8 CONCLUSIONS

In this paper, we study the problem of instant error-tolerant spatial keyword queries on road networks. To efficiently answer the queries, we propose a new index called R2T. R2T employs reverse 2-hop label and tries to seamlessly integrate the spatial and textual information for each vertex of the road network. Based on R2T, we present a suite of algorithms to answer the queries. Both theoretical analysis and empirical evaluation demonstrate the efficiency of our proposed index and algorithms. This work is our first step towards the studied problem. In the future, we would like to investigate the instant error-tolerant spatial keyword queries for moving query object or in a distributed environment.

ACKNOWLEDGMENTS

This work was supported by the NSFC under Grants No. (61972338, 62025206, and 62102351). Qing Liu is the corresponding author of the work.

REFERENCES

- [1] Tenindra Abeywickrama, Muhammad Aamir Cheema, and Arijit Khan. 2020. K-SPIN: Efficiently Processing Spatial Keyword Queries on Road Networks. *IEEE Trans. Knowl. Data Eng.* 32, 5 (2020), 983–997.
- [2] Takuya Akiba, Yoichi Iwata, and Yuichi Yoshida. 2013. Fast Exact Shortest-path Distance Queries on Large Networks by Pruned Landmark Labeling. In *SIGMOD*. 349–360.
- [3] Hannah Bast and Ingmar Weber. 2006. Type Less, Find More: Fast Autocompletion Search with A Succinct Index. In *SIGIR*. 364–371.
- [4] Ramadhana Bramandia, Byron Choi, and Wee Keong Ng. 2010. Incremental Maintenance of 2-Hop Labeling of Large Graphs. *IEEE Trans. Knowl. Data Eng.* 22, 5 (2010), 682–698.
- [5] Lisi Chen, Gao Cong, Christian S. Jensen, and Dingming Wu. 2013. Spatial Keyword Query Processing: An Experimental Evaluation. *Proc. VLDB Endow.* 6, 3 (2013), 217–228.
- [6] Lei Chen, Jianliang Xu, Xin Lin, Christian S. Jensen, and Haibo Hu. 2016. Answering Why-not Spatial Keyword Top- k Queries via Keyword Adaption. In *ICDE*. 697–708.
- [7] Zhida Chen, Lisi Chen, Gao Cong, and Christian S. Jensen. 2021. Location- and Keyword-based Querying of Geo-textual Data: A Survey. *VLDB J.* 30, 4 (2021), 603–640.
- [8] Zitong Chen, Ada Wai-Chee Fu, Minhao Jiang, Eric Lo, and Pengfei Zhang. 2021. P2H: Efficient Distance Querying on Road Networks by Projected Vertex Separators. In *SIGMOD*. 313–325.
- [9] Zhongpu Chen, Bin Yao, Zhi-Jie Wang, Xiaofeng Gao, Shuo Shang, Shuai Ma, and Minyi Guo. 2021. Flexible Aggregate Nearest Neighbor Queries and its Keyword-aware Variant on Road Networks. *IEEE Trans. Knowl. Data Eng.* 33, 12 (2021), 3701–3715.
- [10] Edith Cohen, Eran Halperin, Haim Kaplan, and Uri Zwick. 2002. Reachability and Distance Queries via 2-Hop Labels. In *SODA*. 937–946.
- [11] Gao Cong and Christian S. Jensen. 2016. Querying Geo-textual Data: Spatial Keyword Queries and Beyond. In *SIGMOD*. 2207–2212.
- [12] Gao Cong, Christian S. Jensen, and Dingming Wu. 2009. Efficient Retrieval of the Top- k Most Relevant Spatial Web Objects. *Proc. VLDB Endow.* 2, 1 (2009), 337–348.
- [13] Dong Deng, Guoliang Li, He Wen, H. V. Jagadish, and Jianhua Feng. 2016. META: An Efficient Matching-based Method for Error-tolerant Autocompletion. *Proc. VLDB Endow.* 9, 10 (2016), 828–839.
- [14] Yuyang Dong, Chuan Xiao, Hanxiong Chen, Jeffrey Xu Yu, Kunihiko Takeoka, Masafumi Oyamada, and Hiroyuki Kitagawa. 2021. Continuous Top- k Spatial-keyword Search on Dynamic Objects. *VLDB J.* 30, 2 (2021), 141–161.
- [15] Ada Wai-Chee Fu, Huanhuan Wu, James Cheng, and Raymond Chi-Wing Wong. 2013. IS-LABEL: An Independent-set Based Labeling Scheme for Point-to-point Distance Querying. *Proc. VLDB Endow.* 6, 6 (2013), 457–468.
- [16] Yunjun Gao, Xu Qin, Baihua Zheng, and Gang Chen. 2015. Efficient Reverse Top- k Boolean Spatial Keyword Queries on Road Networks. *IEEE Trans. Knowl. Data Eng.* 27, 5 (2015), 1205–1218.
- [17] Yunjun Gao, Jingwen Zhao, Baihua Zheng, and Gang Chen. 2016. Efficient Collective Spatial Keyword Query Processing on Road Networks. *IEEE Trans. Intell. Transp. Syst.* 17, 2 (2016), 469–480.
- [18] Sheng Hu, Chuan Xiao, and Yoshiharu Ishikawa. 2018. An Efficient Algorithm for Location-aware Query Autocompletion. *IEICE Trans. Inf. Syst.* 101-D, 1 (2018), 181–192.
- [19] Shengyue Ji and Chen Li. 2011. Location-based Instant Search. In *SSDBM*, Vol. 6809. 17–36.
- [20] Minhao Jiang, Ada Wai-Chee Fu, and Raymond Chi-Wing Wong. 2015. Exact Top- k Nearest Keyword Search in Large Networks. In *SIGMOD*. 393–404.
- [21] Guoliang Li, Jianhua Feng, and Chen Li. 2013. Supporting Search-As-You-Type Using SQL in Databases. *IEEE Trans. Knowl. Data Eng.* 25, 2 (2013), 461–475.
- [22] Guoliang Li, Shengyue Ji, Chen Li, and Jianhua Feng. 2009. Efficient Type-ahead Search on Relational Data: A TASTIER approach. In *SIGMOD*. 695–706.
- [23] Guoliang Li, Shengyue Ji, Chen Li, and Jianhua Feng. 2011. Efficient Fuzzy Full-text Type-ahead Search. *VLDB J.* 20, 4 (2011), 617–640.
- [24] Ye Li, Leong Hou U, Man Lung Yiu, and Ngai Meng Kou. 2017. An Experimental Study on Hub Labeling based Shortest Path Algorithms. *Proc. VLDB Endow.* 11, 4 (2017), 445–457.
- [25] Zijian Li, Lei Chen, and Yue Wang. 2019. G*-Tree: An Efficient Spatial Index on Road Networks. In *ICDE*. 268–279.
- [26] Zhisheng Li, Ken C. K. Lee, Baihua Zheng, Wang-Chien Lee, Dik Lun Lee, and Xufa Wang. 2011. IR-Tree: An Efficient Index for Geographic Document Search. *IEEE Trans. Knowl. Data Eng.* 23, 4 (2011), 585–599.
- [27] Ying Lu, Jiaheng Lu, Gao Cong, Wei Wu, and Cyrus Shahabi. 2014. Efficient Algorithms and Cost Models for Reverse Spatial-keyword k -nearest Neighbor Search. *ACM Trans. Database Syst.* 39, 2 (2014), 13:1–13:46.
- [28] Chengyang Luo, Qing Liu, Yunjun Gao, Lu Chen, Ziheng Wei, and Congcong Ge. 2023. TASK: An Efficient Framework for Instant Error-tolerant Spatial Keyword Queries on Road Networks. <https://github.com/ZJU-DAILY/TASK/blob/main/paper/PVLDB2023.pdf>
- [29] Dian Ouyang, Lu Qin, Lijun Chang, Xuemin Lin, Ying Zhang, and Qing Zhu. 2018. When Hierarchy Meets 2-Hop-labeling: Efficient Shortest Distance Queries on Road Networks. In *SIGMOD*. 709–724.
- [30] Miao Qiao, Lu Qin, Hong Cheng, Jeffrey Xu Yu, and Wentao Tian. 2013. Top- k Nearest Keyword Search on Large Graphs. *Proc. VLDB Endow.* 6, 10 (2013), 901–912.
- [31] Jianbin Qin, Chuan Xiao, Sheng Hu, Jie Zhang, Wei Wang, Yoshiharu Ishikawa, Koji Tsuda, and Kunihiko Sadakane. 2020. Efficient Query Autocompletion with Edit Distance-based Error Tolerance. *VLDB J.* 29, 4 (2020), 919–943.
- [32] João B. Rocha-Junior and Kjetil Nørveg. 2012. Top- k Spatial Keyword Queries on Road Networks. In *EDBT*. 168–179.
- [33] Senjuti Basu Roy and Kaushik Chakrabarti. 2011. Location-aware Type Ahead Search on Spatial Databases: Semantics and Efficiency. In *SIGMOD*. 361–372.
- [34] Jin Wang and Chunbin Lin. 2020. Fast Error-tolerant Location-aware Query Autocompletion. In *ICDE*. 1998–2001.
- [35] Xiang Wang, Wenjie Zhang, Ying Zhang, Xuemin Lin, and Zengfeng Huang. 2017. Top- k Spatial-keyword Publish/Subscribe over Sliding Window. *VLDB J.* 26, 3 (2017), 301–326.
- [36] Chuan Xiao, Jianbin Qin, Wei Wang, Yoshiharu Ishikawa, Koji Tsuda, and Kunihiko Sadakane. 2013. Efficient Error-tolerant Query Autocompletion. *Proc. VLDB Endow.* 6, 6 (2013), 373–384.
- [37] Hongfei Xu, Yu Gu, Yu Sun, Jianzhong Qi, Ge Yu, and Rui Zhang. 2020. Efficient Processing of Moving Collective Spatial Keyword Queries. *VLDB J.* 29, 4 (2020), 841–865.
- [38] Junye Yang, Yong Zhang, Xiaofang Zhou, Jin Wang, Huiqi Hu, and Chunxiao Xing. 2019. A Hierarchical Framework for Top- k Location-aware Error-tolerant Keyword Search. In *ICDE*. 986–997.
- [39] Chengyuan Zhang, Ying Zhang, Wenjie Zhang, Xuemin Lin, Muhammad Aamir Cheema, and Xiaoyang Wang. 2014. Diversified Spatial Keyword Search on Road Networks. In *EDBT*. 367–378.
- [40] Mengxuan Zhang, Lei Li, Wen Hua, Rui Mao, Pingfu Chao, and Xiaofang Zhou. 2021. Dynamic Hub Labeling for Road Networks. In *ICDE*. 336–347.
- [41] Mengxuan Zhang, Lei Li, Wen Hua, and Xiaofang Zhou. 2021. Efficient 2-Hop Labeling Maintenance in Dynamic Small-world Networks. In *ICDE*. 133–144.
- [42] Mengxuan Zhang, Lei Li, and Xiaofang Zhou. 2021. An Experimental Evaluation and Guideline for Path Finding in Weighted Dynamic Network. *Proc. VLDB Endow.* 14, 11 (2021), 2127–2140.
- [43] Jingwen Zhao, Yunjun Gao, Gang Chen, and Rui Chen. 2018. Towards Efficient Framework for Time-aware Spatial Keyword Queries on Road Networks. *ACM Trans. Inf. Syst.* 36, 3 (2018), 24:1–24:48.
- [44] Jingwen Zhao, Yunjun Gao, Gang Chen, and Rui Chen. 2018. Why-not Questions on Top- k Geo-social Keyword Queries in Road Networks. In *ICDE*. 965–976.
- [45] Jingwen Zhao, Yunjun Gao, Gang Chen, Christian S. Jensen, Rui Chen, and Deng Cai. 2017. Reverse Top- k Geo-social Keyword Queries in Road Networks. In *ICDE*. 387–398.
- [46] Bolong Zheng, Kai Zheng, Christian S. Jensen, Nguyen Quoc Viet Hung, Han Su, Guohui Li, and Xiaofang Zhou. 2020. Answering Why-not Group Spatial Keyword Queries. *IEEE Trans. Knowl. Data Eng.* 32, 1 (2020), 26–39.
- [47] Bolong Zheng, Kai Zheng, Xiaokui Xiao, Han Su, Hongzhi Yin, Xiaofang Zhou, and Guohui Li. 2016. Keyword-aware Continuous k NN Query on Road Networks. In *ICDE*. 871–882.
- [48] Ruicheng Zhong, Ju Fan, Guoliang Li, Kian-Lee Tan, and Lizhu Zhou. 2012. Location-aware Instant Search. In *CIKM*. 385–394.
- [49] Xiaoling Zhou, Jianbin Qin, Chuan Xiao, Wei Wang, Xuemin Lin, and Yoshiharu Ishikawa. 2016. BEVA: An Efficient Query Processing Algorithm for Error-tolerant Autocompletion. *ACM Trans. Database Syst.* 41, 1 (2016), 5:1–5:44.