

Effective Caching of Shortest Paths for Location-Based Services

ABSTRACT

Web search is ubiquitous in our daily lives. Caching have been extensively used to reduce the computation time of the search engine and reduce the network traffic beyond a proxy server. Another form of web search, known as online shortest path search, are popular due to advances in geo-positioning ~~capability of mobile devices~~. However, existing caching techniques are ineffective for shortest path queries. ~~There are~~ several crucial differences between web search results and shortest path results, in ~~terms of~~ query matching, cache item overlapping, and ~~variation of~~ query cost.

Motivated by this, we ~~first~~ identify several ~~issues~~ that are essential to the success of ~~an~~ effective cache for shortest path search. Our cache exploits the optimal subpath property, which allows a cached shortest path to answer any query with source and target nodes on the path. We utilize statistics from query log to estimate the benefit of caching a specific shortest path, and employ a greedy algorithm for placing beneficial paths in the cache. Also, we design a compact cache structure that supports efficient query matching at run-time. Empirical results on real datasets confirm the effectiveness of our proposed techniques.

1. INTRODUCTION

Web search is ubiquitous and extensively used in our daily lives. The scenario for typical web search is illustrated in Figure 1. A user may submit a query “Paris Eiffel Tower” to the search engine, which then computes relevant results and reports them back to the user. Usually, a *cache* is used to keep the results of frequent queries in order to achieve a high hit ratio [1, 2, 15, 16]. The cache can be employed at the search engine to save its computation time, e.g., when the query (result) can be found in the cache. To improve ~~the~~ ~~user~~ response time, a cache can be placed at a proxy, ~~which~~ resides in the same sub-network as the user. A query result that is available at the proxy can be reported immediately, without contacting the search engine.

The scenario of Figure 1 is applicable to *online shortest path search* on a map as well, which are popular due to advances in geo-positioning capability of mobile devices (e.g., PDA, smartphone). It has various applications for mobile users, e.g., a tourist asking

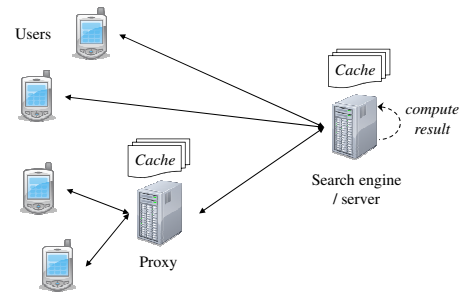


Figure 1: Scenario for Web Search

the way to a museum, a driver finding the route to a gas station, and a customer searching the path to a restaurant. When compared to commercial navigation software, online search (e.g., Google Maps, MapQuest) provide ~~these~~ benefits to mobile users: (i) they are free to use, and (ii) they do not require any installation and storage space at mobile devices. Figure 2 shows a road network in which a node v_i is a road junction and an edge (v_i, v_j) models a road segment with its distance shown as a number. The shortest path from v_1 to v_7 is the path $\langle v_1, v_3, v_4, v_5, v_7 \rangle$ and its path distance $3+6+9+5 = 23$ ~~is the minimum~~. Again, caching can be utilized at a proxy to reduce the ~~user~~ response time, or used at the server to reduce ~~its~~ computation time.

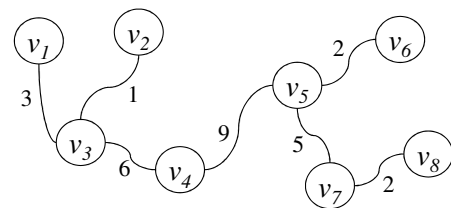


Figure 2: An Example Road Network

~~In this paper, we~~ study the caching of path search results in a scenario like Figure 1. Nevertheless, there are crucial differences between web search results and shortest path results, rendering existing caching techniques for web results ineffective in our context.

- **Exact matching vs. subpath matching:** The result of a web query (e.g., “Paris Eiffel Tower”) seldom matches with that of another query (e.g., “Paris Louvre Palace”). In contrast, a shortest path result ~~may~~ contain ~~several~~ paths that can be used for answering other queries. For example, the shortest path from v_1 to v_7 ($\langle v_1, v_3, v_4, v_5, v_7 \rangle$) contains the shortest

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$10.00.

path from v_3 to v_5 , the shortest path from v_4 to v_7 , etc. We need to re-visit the model for defining the benefit of a path.

- **Cache structure:** Web search caching may employ a hash table to check efficiently whether a query can be found in the cache. However, hash table cannot support subpath matching in our problem. A new structure is required to organize the cache content in an effective way for supporting subpath matching. Furthermore, this problem is complicated by the overlapping of paths. For example, the shortest path $\langle v_1, v_3, v_4, v_5, v_7 \rangle$ and the shortest path $\langle v_2, v_3, v_4, v_5, v_6 \rangle$ have significant overlap, although none of them contains the other path. It is challenging to develop a compact structure for storing these paths.
- **Query processing cost:** At the server side, when there is a cache miss, an algorithm will be invoked to compute the query result. In fact, some results are more expensive to obtain than other results. To optimize the server performance, some cost model has been developed to estimate the cost of a web query [1], in order to decide the importance of a result. However, to our knowledge, there is no work on estimating the cost of a shortest path query with respect to a given shortest path algorithm.

In order to tackle the above challenges, we contribute the followings:

- We formulate a systematic model for quantifying the benefit of caching a specific shortest path.
- We design techniques to extract statistics from query log and benchmark the cost of a shortest path call.
- We propose an algorithm for selecting paths to be placed in the cache.
- We develop a compact and efficient cache structure for storing shortest paths.

The rest of the paper is organized as follows. Section 2 studies the work related to our shortest path caching problem. Section 3 defines our problem formally and Section 4 formulates a model for capturing the benefit of a cached shortest path, examines the query frequency and the cost of shortest path call, and presents an algorithm for selecting appropriate paths to be placed in the cache. Section 5 investigates how to design a compact and efficient cache structure for storing shortest paths. Our proposed methods are then evaluated on real data in Section 6, followed by a conclusion in Section 7.

2. RELATED WORK

Web Search Caching: A web search query asks for the top- K relevant documents (i.e., web pages) that contain all query keywords, e.g., “Paris Eiffel Tower”. The typical value of K is the number of results (10) to be displayed on a result page [15]; a request for the next result page is interpreted as another query.

Web search caching have been used to improve the performance of the search engine. When a query can be answered by the cache, the cost of computing the query result can be saved. Markatos et al. [14] is a pioneering work on evaluating two caching approaches (dynamic caching and static caching) on real query logs. *Dynamic caching* [6, 13, 14] aims to cache the results of the most recently accessed queries. For example, in the Least-Recently-Used (LRU) method, when a query causes a cache miss, the least recently used result in the cache will be replaced by the current query result. This approach keeps the cache up-to-date and adapts quickly to

the distribution of the incoming queries; however, it incurs overhead on updating the cache frequently. On the other hand, *static caching* [1–3, 14–16] aims to cache the results of the most popular queries. It exploits a query log, which contains users’ queries in the past, to determine the most frequent queries. The above studies have shown that the frequency of queries follow the Zipfian distribution, i.e., a small number of queries have very high frequency and they remain popular for a period of time. Although the cache content is not the most up-to-date, it is able to answer the majority of frequent queries. A static cache can be updated periodically (e.g., every day) based on the latest query log. Static caching has an advantage that it incurs very low overhead at query time.

Early works on web search caching adopt the *cache hit ratio* as the performance metric. This metric reflects the number of queries that do not require computation cost. Recent work [1, 2, 15] on web search caching use the *server processing time* as the performance metric. They show that different queries have different query processing times. For example, a query that involves terms with large posting lists is expected to incur high cost. Refs. [1, 15] propose models for estimating the cost of processing queries based on the sizes of inverted lists, and exploit both frequency and cost information in their static caching methods.

None of the above work have studied the caching of shortest paths. In this paper, we adopt static caching approach because it performs well on query logs that are typically skewed in nature; and it incurs very low overhead at query time. The cost models of [1, 15] are specific to web search queries and inapplicable to our problem. In our caching problem, different shortest path queries also have different processing times. Thus, we will study a cost-oriented model for quantifying the benefit of placing a specific path in the cache.

Semantic Caching: In a client-server system, a cache may be employed at the client-side in order to reduce the communication cost and improve response latency of the client. The cache located at a client can only serve queries from the client itself but not from other clients. Such a cache is only beneficial for a query-intensive user. All techniques in this category adopt the dynamic caching approach.

Semantic caching [5] is a client-side caching model that associates cached results with valid ranges. Upon receiving a query Q , the relevant results from the cache are reported. A subquery Q' is constructed from Q such that Q' covers the query region that cannot be answered by the cache. The subquery Q' is then forwarded to the server in order to obtain the remaining results of Q . The work on [5] focuses on semantic caching of relational datasets. As example, assume that the dataset stores the age of each employee; and the cache contains the result of the query “find employees with age below 30”. Suppose that the client now issues a query Q “find employees with age between 20 and 40”. First, the employees with age between 20 and 30 can be obtained from the cache. Then, a subquery Q' “find employees with age between 30 and 40” is submitted to server for retrieving the remaining results.

Semantic caching has also studied for spatial data [9, 12, 20]. Zheng et al. [20] define the semantic region of a spatial object as its Voronoi cell, which can then be utilized to answer nearest neighbor queries for a moving client user. Hu et al. [9] study semantic caching of tree nodes in an R-tree index and examine how to process generic spatial queries on the cached tree nodes. Lee et al. [12] build generic semantic regions for spatial objects so that they are able to support generic spatial queries. However, there has not been any semantic caching technique for graphs or shortest path

results.

Shortest Path Computation: Existing shortest path indexes can be categorized into three types, which have trade-offs between their precomputation effort and query performance. A basic structure is the adjacency list, in which each node v_i is assigned a list for storing the adjacent nodes of v_i . It does not store any pre-computed information. Uninformed search (e.g., Dijkstra’s algorithm, bidirectional search) can be used to compute the shortest path; however, it incurs high query cost. *Fully-precomputed* indexes, e.g., distance index [8], shortest path quadtree [18], require precomputing the shortest paths between any two nodes in the graph. Although they support efficient querying, they incur huge precomputation time $O(|V|^3)$ and storage space $O(|V|^2)$, where $|V|$ is the number of nodes in the graph. *Partially-precomputed* indexes, e.g., landmark [11], HiTi [10], and TEDI [19] attempt to materialize some distances / paths in order to accelerate the processing of shortest path queries. They employ some parameter to control the trade-off between the query performance, precomputation overhead, and storage space.

As a possible approach to our caching problem, one could assume that a specific shortest path index is being used at the server. A portion of the index may be cached so that it can be used to answer certain queries rapidly. Unfortunately, this approach is tightly coupled to the assumed index. It is inapplicable to servers that employ other indexes different from the assumed index. Also, it cannot be directly applied to any new index developed in future.

In this paper, we view the shortest path index/algorithm as a black-box and decouple it from the cache. The main advantage is that our approach is applicable to any shortest path index, without knowing their implementation. Any new shortest path method can be integrated with our proposed cache seamlessly.

3. PROBLEM SETTING

We first introduce some background definitions and properties. Then we present our problem and objectives. Table 1 summarizes the notations to be used in this paper.

Notation	Meaning
$G(V, E)$	a graph with node set V and edge set E
v_i	a node in V
(v_i, v_j)	an edge in E
$W(v_i, v_j)$	the edge weight of $W(v_i, v_j)$
$Q_{s,t}$	shortest path query from node v_s to node v_t
$P_{s,t}$	the shortest path result of $Q_{s,t}$
$ P_{s,t} $	the size of $P_{s,t}$ (in number of nodes)
$E_{s,t}$	the expense of executing query $Q_{s,t}$
$\chi_{s,t}$	The frequency of a SP
Ψ	the cache
$\mathcal{U}(P_{s,t})$	the set of all subpaths in $P_{s,t}$
$\mathcal{U}(\Psi)$	the set of all subpaths of paths in Ψ
$\gamma(\Psi)$	the total benefit of the content in the cache
\mathcal{QL}	query log

Table 1: Summary of Notations

3.1 Definitions and Properties

We first provide the definitions of graph and shortest path.

DEFINITION 1. Graph model.

Let $G(V, E)$ be a graph with a set V of nodes and a set E of edges. Each node $v_i \in V$ models a road junction. Each edge $(v_i, v_j) \in E$ models a road segment, and its distance is denoted as $W(v_i, v_j)$.

DEFINITION 2. Shortest path: query and result.

A shortest path query, denoted by $Q_{s,t}$, consists of a source node v_s and a target node v_t .

The result of $Q_{s,t}$, denoted by $P_{s,t}$, is the path from v_s to v_t (on graph G) with the minimum sum of edge weights along the path. We can represent $P_{s,t}$ as a list of nodes: $\langle v_{x_0}, v_{x_1}, v_{x_2}, \dots, v_{x_m} \rangle$, where $v_{x_0} = s$, $v_{x_m} = t$, and the path distance is: $\sum_{i=0}^{m-1} W(v_{x_i}, v_{x_{i+1}})$.

We consider only undirected graphs in our examples. Our techniques can be easily applied to directed graphs. In the example graph of Figure 2, the shortest path from v_1 to v_7 is the path $P_{1,7} = \langle v_1, v_3, v_4, v_5, v_7 \rangle$ with its path distance $3+6+9+5 = 23$.

Shortest path exhibits the optimal subpath property (see Lemma 1). It states that every subpath of a shortest path is also a shortest path. For example, in Figure 2, the shortest path $P_{1,7} = \langle v_1, v_3, v_4, v_5, v_7 \rangle$ contains these shortest paths: $P_{1,3}, P_{1,4}, P_{1,5}, P_{1,7}, P_{3,4}, P_{3,5}, P_{3,7}, P_{4,5}, P_{4,7}, P_{5,7}$. As we will discuss shortly, this property can be exploited for caching shortest paths.

LEMMA 1. Optimal subpath property (from [4]).

The shortest path $P_{a,b}$ contains the shortest path $P_{s,t}$ if $s \in P_{a,b}$ and $t \in P_{a,b}$. Specifically, let $P_{a,b} = \langle v_{x_0}, v_{x_1}, v_{x_2}, \dots, v_{x_m} \rangle$. We have $P_{s,t} = \langle v_{x_i}, v_{x_{i+1}}, \dots, v_{x_j} \rangle$ if $s = v_{x_i}$ and $t = v_{x_j}$.

3.2 Problem and Objectives

We adopt the architecture of Figure 1 to our caching problem. Using mobile devices with positioning capabilities, users issue shortest path queries to an online server. The cache, as defined below, can be placed at either a proxy or the server. It helps improve the computation cost and the communication cost at the server/proxy, as well as reduce the response time of shortest path queries.

DEFINITION 3. Cache and budget.

Given a cache budget \mathcal{B} , a cache Ψ is allowed to store a collection of shortest path results such that $|\Psi| \leq \mathcal{B}$, where the cache size $|\Psi|$ is taken as the total number of nodes of shortest paths in Ψ .

As discussed in Section 2, recent literature on web search caching [1, 2, 15, 16] advocates the use of a static cache rather than a dynamic cache. Static caching has very low runtime overhead, while it only sacrifices the hit ratio a little. Thus, in this paper, we also adopt the static caching paradigm and exploit a query log to build a cache.

DEFINITION 4. Query log.

The query log \mathcal{QL} is a collection of timestamped queries that have been issued by users in the past.

We identify essential components in our static caching system in Figure 3. It involves: (i) a shortest path API, (ii) a cache, (iii) online module for cache lookup, and (iv) offline/periodic modules for collecting query log, benchmarking the cost of API, and building the cache.

Observe the shortest path API (in gray) is given by the system so we are not allowed to modify its implementation. For the server scenario, the shortest path API is linked to a typical shortest path algorithm (e.g., Dijkstra, A* search). For the proxy scenario, the shortest path API triggers the issue of a query message to the server. In either case, calling the shortest path API incurs expensive computation/communication cost, as defined below. Different queries may have different expenses. In general, a long-range query incurs more expense than a short-range query.

DEFINITION 5. Expense of executing query.

We define $E_{s,t}$ as the expense of the shortest path API to process query $Q_{s,t}$.

We employ a *cache* to reduce the overall cost of invoking the shortest path API. Upon receiving a query (at runtime), the server/proxy checks whether the cache contains the query result. If this is a hit, then the result from the cache is reported to the user immediately. This saves the expensive cost from calling the shortest path API. Otherwise, the result must be obtained by calling the shortest path API.

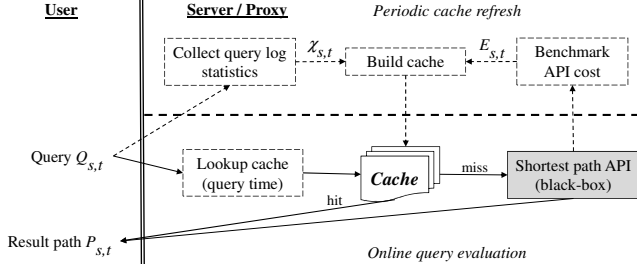


Figure 3: Components in a Static Caching System

Unlike web search caching, we observe that minimizing the cache hit ratio does not necessarily mean that the overall cost is reduced significantly. In the server scenario, the cost of calling the shortest path API (e.g., shortest path algorithm) is not fixed and heavily depends on the distance of the shortest path. As a case study, we randomly generate 500 shortest path queries on two real road networks (their descriptions available in Section 6), execute a shortest path algorithm (Dijkstra) for them. In Figure 4, each subfigure (for each road network) shows the shortest path distance (x-axis) and the cost (y-axis) of each individual query. Observe that there is a strong correlation between the cost and the shortest path distance. Caching a short-range path may only provide negligible improvement, even if the path is frequently queried. Therefore, we plan to design a benchmark on the cost of calling the API.

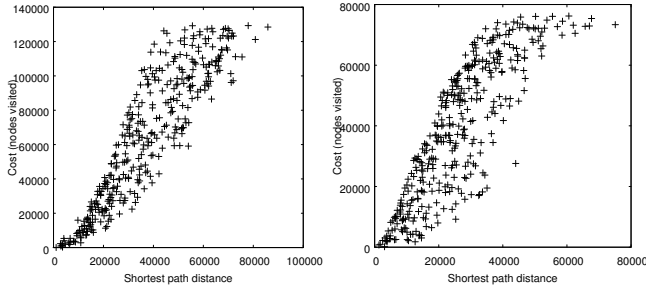


Figure 4: Cost vs. Distance of a Shortest Path Algorithm

Adopting the static caching paradigm, the server/proxy collects the query log and builds the cache periodically (e.g., every day). By extracting the distribution of the query log, we are able to estimate the probability of a specific shortest path being queried in future. Combining such information with the benchmark model, we can pick promising paths into the cache in order to minimize the overall cost of the system. We will also investigate the structure of the cache; it should be compact in order to accommodate as many paths as possible, and it should support efficient result retrieval.

Our main objective is to reduce the overall cost of calling the shortest path API. We define this as the problem below. In Sec-

tion 4, we will formulate the cache benefit notion $\gamma(\Psi)$, extract statistics to compute $\gamma(\Psi)$, and present an algorithm to this cache benefit maximization problem.

PROBLEM-1. Static cache benefit maximization problem.

Given a cache budget B and a query log QL , this problem is to build a cache Ψ with the maximum cache benefit $\gamma(\Psi)$ subject to the budget constraint B , where Ψ contains result paths $P_{s,t}$ whose queries $Q_{s,t}$ belong to QL .

Our secondary objectives are to: (i) develop a compact cache structure to maximize the accommodation of shortest paths, and (ii) provide efficient means for retrieving result from the cache. We will focus on these issues in Section 5.

3.3 Existing Solutions for Caching Results

We now revisit existing solutions for caching web search results [3, 14] and then explain why they are inadequate for caching shortest path results.

Dynamic Caching — LRU: A typical dynamic caching method for web search is the Least-Recently-Used (LRU) method [14]. When a new query is submitted, its result is inserted into the cache. The least-recently-used result would be evicted from the cache when it becomes full.

We proceed to illustrate the running steps of LRU on the map from Figure 2. Suppose that the cache budget B is 10 (i.e., space for 10 nodes). Table 2 shows the query and the cache content at each timestamp T_i . Each cached path is associated with the last timestamp being used. At times T_1 and T_2 , both queries produce cache misses and their results ($P_{3,6}$ and $P_{1,6}$) are inserted into the cache (as they fit). At time T_3 , query $Q_{2,7}$ causes a cache miss as it cannot be answered by any cached path. Before inserting its result $P_{2,7}$ into the cache, the least recently used path $P_{3,6}$ is evicted from the cache. At time T_4 , query $Q_{1,4}$ contributes a cache hit; it can be answered by the cached path $P_{1,6}$ because the source and target nodes v_1, v_4 fall on $P_{1,6}$ (see Lemma 1). The running steps at subsequent timestamps are shown in Table 2. In total, the LRU cache has only 2 hits.

Time	$Q_{s,t}$	$P_{s,t}$	Paths in LRU cache	event
T_1	$Q_{3,6}$	$\langle v_3, v_4, v_5, v_6 \rangle$	$P_{3,6} : T_1$	miss
T_2	$Q_{1,6}$	$\langle v_1, v_3, v_4, v_5, v_6 \rangle$	$P_{1,6} : T_2, P_{3,6} : T_1$	miss
T_3	$Q_{2,7}$	$\langle v_2, v_3, v_4, v_5, v_7 \rangle$	$P_{2,7} : T_3, P_{1,6} : T_2$	miss
T_4	$Q_{1,4}$	$\langle v_1, v_3, v_4 \rangle$	$P_{1,6} : T_4, P_{2,7} : T_3$	hit
T_5	$Q_{4,8}$	$\langle v_4, v_5, v_7, v_8 \rangle$	$P_{4,8} : T_5, P_{1,6} : T_4$	miss
T_6	$Q_{2,5}$	$\langle v_2, v_3, v_4, v_5 \rangle$	$P_{2,5} : T_6, P_{4,8} : T_5$	miss
T_7	$Q_{3,6}$	$\langle v_3, v_4, v_5, v_6 \rangle$	$P_{3,6} : T_7, P_{2,5} : T_6$	miss
T_8	$Q_{3,6}$	$\langle v_3, v_4, v_5, v_6 \rangle$	$P_{3,6} : T_8, P_{2,5} : T_6$	hit

Table 2: Example of LRU on a Sequence of Queries

Observe that LRU cannot determine the benefit of a path effectively. For example, the paths $P_{1,6}$ and $P_{2,7}$ (obtained at times T_2 and T_3) can answer many queries in subsequent timestamps, e.g., $Q_{1,4}, Q_{2,5}, Q_{3,6}, Q_{3,6}$. If they have been kept in the cache, there would be 4 cache hits. However, LRU mistakenly evicts them before they can be used to answer other queries.

Another limitation of LRU is that it is not designed to support subpath matching efficiently. Upon receiving a query $Q_{s,t}$, every path in the cache needs to be scanned in order to check whether the path contains s and t . This incurs significant runtime overhead, outweighing the advantages of the cache.

Static Caching — HQF: A typical static caching method for web search is the Highest-Query-Frequency (HQF) method [3]. In an

offline phase, the most frequent queries are selected from the query log \mathcal{QL} and then their results are inserted into the cache. The cache content remains unchanged during runtime.

According to Ref. [3], the frequency of a query is counted as the number of queries in \mathcal{QL} identical to it. Let's consider an example query log: $\mathcal{QL} = \{Q_{3,6}, Q_{1,6}, Q_{2,7}, Q_{1,4}, Q_{4,8}, Q_{2,5}, Q_{3,6}, Q_{3,6}\}$. Since $Q_{3,6}$ has the highest frequency (3), so HQF would pick the corresponding result path $P_{3,6}$. It fails to pick $P_{1,6}$ simply because its query $Q_{1,6}$ has a low frequency (1). In fact, the path $P_{1,6}$ is more promising than $P_{3,6}$ because $P_{1,6}$ can be used to answer more queries than $P_{3,6}$. This arises a problem in HQF because the query frequency definition does not capture characteristics specific to our problem—shortest paths may overlap, and the result of a query may be used to answer multiple other queries.

Other common limitations of LRU and HQF: Furthermore, both LRU and HQF have not considered the variations in the expense of obtaining shortest paths. Consider the cache in the server scenario as an example. Intuitively, it is more expensive to process a long-range query rather than a short-range query. Caching an expensive-to-obtain path could lead to great savings in future. An appropriate choice of cached path should take such an expense into account.

Also, they have not studied the utilization of the cache space for shortest paths. For example, in Table 2, the paths in the cache have overlapping and cause wasted space on storing duplicated nodes among the paths. It is important to design a compact cache structure that exploits the sharing among paths to avoid storing duplicated nodes.

4. BENEFIT-DRIVEN CACHING

In this section, we propose a benefit-driven approach to decide which shortest paths should be placed in the cache. Section 4.1 formulates a systematic model for capturing the *benefit* of a cache on potential queries. This model requires the knowledge of: (i) frequency $\chi_{s,t}$ of a query, and (ii) expense $E_{s,t}$ of processing a query. Thus, we investigate how to extract query frequency from a query log in Section 4.2, and benchmark the expense of processing a query in Section 4.3. Finally, in Section 4.4, we present an algorithm for selecting promising paths to be placed in the cache.

4.1 Benefit Model

~~This section presents a benefit model of a cache on shortest path queries.~~ We first study the benefit of a cached shortest path and then examine the benefit of a cache.

First, we consider a cache Ψ that contains one shortest path $P_{a,b}$ only. Recall from Figure 3 that, when a query $Q_{s,t}$ can be answered by a cached path $P_{a,b}$, this produces a cache hit and avoids the cost of invoking the shortest path API. In order to model the benefit of $P_{a,b}$, we must address the questions below:

1. Which queries $Q_{s,t}$ can be answered by the path $P_{a,b}$?
2. For query $Q_{s,t}$, how much cost can we save?

The first question is answered by Lemma 1. The path $P_{a,b}$ contains the path $P_{s,t}$ when both nodes v_s and v_t appear in $P_{a,b}$. Thus, we define the *answerable query set* of the path $P_{a,b}$ as:

$$\mathcal{U}(P_{a,b}) = \{P_{s,t} : s \in P_{a,b}, t \in P_{a,b}, s \neq t\} \quad (1)$$

It contains the set of queries that can be answered by $P_{a,b}$. Taking Figure 2 as the example graph, the answerable query set of path $P_{1,6}$ is: $\mathcal{U}(P_{1,6}) = \{P_{1,3}, P_{1,4}, P_{1,5}, P_{1,6}, P_{3,4}, P_{3,5}, P_{3,6}, P_{4,5}, P_{4,6}, P_{5,6}\}$. Table 3 shows the answerable query sets of other paths.

$P_{a,b}$	$\mathcal{U}(P_{a,b})$
$P_{1,4}$	$P_{1,3}, P_{1,4}, P_{3,4}$
$P_{1,6}$	$P_{1,3}, P_{1,4}, P_{1,5}, P_{1,6}, P_{3,4}, P_{3,5}, P_{3,6}, P_{4,5}, P_{4,6}, P_{5,6}$
$P_{2,5}$	$P_{2,3}, P_{2,4}, P_{2,5}, P_{3,4}, P_{3,5}, P_{4,5}$
$P_{2,7}$	$P_{2,3}, P_{2,4}, P_{2,5}, P_{2,7}, P_{3,4}, P_{3,5}, P_{3,7}, P_{4,5}, P_{4,7}, P_{5,7}$
$P_{3,6}$	$P_{3,4}, P_{3,5}, P_{3,6}, P_{4,5}, P_{4,6}, P_{5,6}$
$P_{4,8}$	$P_{4,5}, P_{4,7}, P_{4,8}, P_{5,7}, P_{5,8}, P_{7,8}$
$\mathcal{U}(\Psi)$, when $\Psi = \{P_{1,6}, P_{3,6}\}$	
$P_{1,3}, P_{1,4}, P_{1,5}, P_{1,6}, P_{3,4}, P_{3,5}, P_{3,6}, P_{4,5}, P_{4,6}, P_{5,6}$	

Table 3: Example of $\mathcal{U}(P_{s,t})$ and $\mathcal{U}(\Psi)$

Regarding the second question, the expected cost saving for query $Q_{s,t}$ depends on (i) its query frequency $\chi_{s,t}$, and (ii) the expense $E_{s,t}$ for the shortest path API to process it. Since the path $P_{a,b}$ can answer query $Q_{s,t}$, we save cost $E_{s,t}$ for $\chi_{s,t}$ times, i.e., $\chi_{s,t} \cdot E_{s,t}$ in total.¹

Combining the answers of both questions, we define the *benefit* of path $P_{a,b}$ as:

$$\gamma(P_{a,b}) = \sum_{P_{s,t} \in \mathcal{U}(P_{a,b})} \chi_{s,t} \cdot E_{s,t} \quad (2)$$

The path benefit $\gamma(P_{a,b})$ answers the question: “If path $P_{a,b}$ is in the cache, how much cost can we save in total?”

Let's assume that we are given the values of $\chi_{s,t}$ and $E_{s,t}$ for all pairs of v_s and v_t , as shown in Figure 5. We will study how to derive them in subsequent sections. To compute $\gamma(P_{1,6})$ of path $P_{1,6}$, we first find its answerable query set $\mathcal{U}(P_{1,6})$ (see Table 3). Since $\mathcal{U}(P_{1,6})$ contains the path $P_{1,4}$, it contributes a benefit of $\chi_{1,4} \cdot E_{1,4} = 1 \cdot 2$ (by lookup of Figure 5). Summing up the benefits of all paths in $\mathcal{U}(P_{1,6})$, we thus obtain: $\gamma(P_{1,6}) = 0 + 1 \cdot 2 + 0 + 1 \cdot 4 + 0 + 0 + 3 \cdot 3 + 0 + 0 + 0 = 15$. Similarly, we can compute the benefit of path $P_{3,6}$: $\gamma(P_{3,6}) = 0 + 0 + 3 \cdot 3 + 0 + 0 + 0 = 9$.

We then extend our equations to the general case—a cache with multiple shortest paths. Observe that a query can be answered by the cache Ψ if it can be answered by any path $P_{a,b}$ in Ψ . Thus, we define the answerable query set of Ψ as the union of all $\mathcal{U}(P_{a,b})$, and define the benefit of Ψ accordingly.

$$\mathcal{U}(\Psi) = \bigcup_{P_{a,b} \in \Psi} \mathcal{U}(P_{a,b}) \quad (3)$$

$$\gamma(\Psi) = \sum_{P_{s,t} \in \mathcal{U}(\Psi)} \chi_{s,t} \cdot E_{s,t} \quad (4)$$

The cache benefit $\gamma(\Psi)$ answers the question: “Using this cache Ψ , how much cost can we save in total?”

Suppose that the cache Ψ contains two paths $P_{1,6}$ and $P_{3,6}$. The answerable query set $\mathcal{U}(\Psi)$ of Ψ is shown in Table 3. By Equation 4, we compute the cache benefit as: $\gamma(\Psi) = 1 \cdot 2 + 1 \cdot 4 + 3 \cdot 3 = 15$.

Note that $\gamma(\Psi)$ is not a distributive function. For example, $\gamma(P_{1,6}) + \gamma(P_{3,6}) = 15 + 9 = 24 \neq \gamma(\Psi)$. Since the path $P_{3,6}$ appears in both answerable query sets $\mathcal{U}(P_{1,6})$ and $\mathcal{U}(P_{3,6})$, the benefit contributed by $P_{3,6}$ is double-counted in the sum $\gamma(P_{1,6}) + \gamma(P_{3,6})$. On the other hand, the value of $\gamma(\Psi)$ is correct because the path $P_{3,6}$ appears exactly once in the answerable query set $\mathcal{U}(\Psi)$ of the cache.

Benefit per size: The benefit model does not consider the size

¹We ignore the overhead of cache lookup as it is negligible compared to the expense $E_{s,t}$ of processing a query $Q_{s,t}$. Efficient cache structures will be studied in Section 5.

$|P_{a,b}|$ of a path $P_{a,b}$ into account. Suppose that we are given two paths $P_{a,b}$ and $P_{a',b'}$ such that they have the same benefit (i.e., $\gamma(P_{a,b}) = \gamma(P_{a',b'})$) and $P_{a',b'}$ has a smaller size than $P_{a,b}$. Intuitively, we prefer path $P_{a',b'}$ rather than path $P_{a,b}$ because $P_{a',b'}$ occupies less space, leaving more cache space for placing other paths. Thus, we define the *benefit-per-size* of a path $P_{a,b}$ as:

$$\bar{\gamma}(P_{a,b}) = \frac{\gamma(P_{a,b})}{|P_{a,b}|} \quad (5)$$

We will utilize this notion in Section 4.4.

Recall from Section 3.2 that our main problem is to build a cache Ψ such that its benefit $\gamma(\Psi)$ is maximized. This requires the values of $\chi_{s,t}$ and $E_{s,t}$. We proceed to discuss how to obtain these values in subsequent sections.

$\chi_{s,t}$	v_1	v_2	v_3	v_4	v_5	v_6	v_7	v_8
v_1	/	0	0	1	0	1	0	0
v_2	0	/	0	0	1	0	1	0
v_3	0	0	/	0	0	3	0	0
v_4	1	0	0	/	0	0	0	1
v_5	0	1	0	0	/	0	0	0
v_6	1	0	3	0	0	/	0	0
v_7	0	1	0	0	0	0	/	0
v_8	0	0	0	1	0	0	0	/

(a) $\chi_{s,t}$ values

$E_{s,t}$	v_1	v_2	v_3	v_4	v_5	v_6	v_7	v_8
v_1	/	2	1	2	3	4	4	5
v_2	2	/	1	2	3	4	4	5
v_3	1	1	/	1	2	3	3	4
v_4	2	2	1	/	1	2	2	3
v_5	3	3	2	1	/	1	1	2
v_6	4	4	3	2	1	/	2	3
v_7	4	4	3	2	1	2	/	1
v_8	5	5	4	3	2	3	1	/

(b) $E_{s,t}$ values

Figure 5: Example of $\chi_{s,t}$ and $E_{s,t}$ values for the graph

4.2 Extracting $\chi_{s,t}$ from Query Log

The frequency $\chi_{s,t}$ of query $Q_{s,t}$ plays an important role in our benefit model. According to a scientific study [7], the mobility patterns of human users follow a skewed distribution. For instance, hot regions (e.g., shopping malls, residential buildings) are expected to have high $\chi_{s,t}$ whereas rural regions are likely to have low $\chi_{s,t}$.

In this section, we study automatic techniques for deriving the values of $\chi_{s,t}$. In our caching system (see Figure 3), the server/proxy periodically collects the query log \mathcal{QL} and extracts the values of $\chi_{s,t}$. Literature on static web caching [2] suggest that the query frequency is stable within a month and it can be used as the periodic time interval. We first study a simple method to extract $\chi_{s,t}$, and then propose a more effective method for extracting $\chi_{s,t}$.

Timestamp	T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
Query	$Q_{3,6}$	$Q_{1,6}$	$Q_{2,7}$	$Q_{1,4}$	$Q_{4,8}$	$Q_{2,5}$	$Q_{3,6}$	$Q_{3,6}$

Table 4: Query log \mathcal{QL}

Node-pair frequency counting: In this method, we first create a *node-pair frequency table* χ with $|V| \times |V|$ entries, like the one in Figure 5a. The entry in the s -th row and the t -th column represents the value of $\chi_{s,t}$. Note that the storage space of the table is $O(|V|^2)$, regardless of how large the query log is.

At the beginning, all entries in the table are initialized to zero. Next, we examine each query $Q_{s,t}$ in the query log \mathcal{QL} and increment the entry $\chi_{s,t}$ (and $\chi_{t,s}$). Let's take the query log \mathcal{QL} in Table 4 as an example. For the first query $Q_{3,6}$ in \mathcal{QL} , we increment the entries $\chi_{3,6}$ and $\chi_{6,3}$. Continuing this process with the other queries in \mathcal{QL} , we obtain the table χ as shown in Figure 5a. The $\chi_{s,t}$ values in the table χ can then be readily used by our benefit model in Section 4.1.

Region-pair frequency counting: The node-pair frequency table χ requires $O(|V|^2)$ space, which cannot fit into main memory even

for a road network of moderate size (e.g., $|V| = 100,000$). To tackle this issue, we propose to: (i) partition the graph into L regions (where L is system parameter), and (ii) employ a compact table for storing the query frequency among pairs of regions only.

For the first step, we can apply any existing graph partitioning technique (e.g., kD-tree partitioning, spectral partitioning). kD-tree partitioning is applicable to the majority of road networks whose nodes are associated with coordinates. For the other graphs, we may apply the spectral partitioning which does not require node coordinates. In Figure 6, we apply a kD-tree on the coordinates of nodes in order to partition the graph into $L = 4$ regions: R_1, R_2, R_3, R_4 .

For the second step, we create a *region-pair frequency table* $\hat{\chi}$ with $L \times L$ entries, like the one in Figure 7b. The entry in the R_i -th row and the R_j -th column represents the value of $\hat{\chi}_{R_i, R_j}$. The storage space of this table is only $O(L^2)$ and it can be controlled by the parameter L . Initially, all entries in the table are set to zero. For each query $Q_{s,t}$ in the query log \mathcal{QL} , we first find the region (say, R_i) that contains node v_s and the region (say, R_j) that contains node v_t . Then, we increment the entry $\hat{\chi}_{R_i, R_j}$ and $\hat{\chi}_{R_j, R_i}$. As an example, we read the query log \mathcal{QL} in Table 4 and examine the first query $Q_{3,6}$. We can find that nodes v_3 and v_6 fall in the regions R_2 and R_3 respectively. Thus, we increment the entries $\hat{\chi}_{R_2, R_3}$ and $\hat{\chi}_{R_3, R_2}$. Continuing this process with the other queries in \mathcal{QL} , we obtain the table $\hat{\chi}$ as shown in Figure 7b.

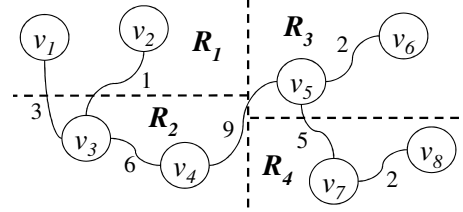


Figure 6: A graph partitioned into 4 regions

$R_1 :$	$\{v_1, v_2\}$
$R_2 :$	$\{v_3, v_4\}$
$R_3 :$	$\{v_5, v_6\}$
$R_4 :$	$\{v_7, v_8\}$

χ_{R_i, R_j}	R_1	R_2	R_3	R_4
R_1	0	1	2	1
R_2	1	0	3	1
R_3	2	3	0	0
R_4	1	1	0	0

(a) node sets of regions (b) region-pair frequency table $\hat{\chi}$

Figure 7: Counting region-pair frequency in table $\hat{\chi}$

To enable the computation of our benefit model in Section 4.1, we now discuss how to derive the value of $\chi_{s,t}$ from the region-pair frequency table $\hat{\chi}$. Note that the frequency of $\hat{\chi}_{R_i, R_j}$ is contributed by any pair of nodes (v_s, v_t) such that region R_i contains v_s and region R_j contains v_t . Thus, we obtain: $\hat{\chi}_{R_i, R_j} = \sum_{v_s \in R_i} \sum_{v_t \in R_j} \chi_{s,t}$. If we make the uniformity assumption within a region, then we have: $\hat{\chi}_{R_i, R_j} = |R_i| \cdot |R_j| \cdot \chi_{s,t}$ where $|R_i|$ and $|R_j|$ denote the number of nodes in the regions R_i and R_j respectively. In other words, we compute the value of $\chi_{s,t}$ from the $\hat{\chi}$ as follows:

$$\chi_{s,t} = \frac{\hat{\chi}_{R_i, R_j}}{|R_i| \cdot |R_j|} \quad (6)$$

Note that the value of $\chi_{s,t}$ is only computed when it is needed. No additional storage space is required to store $\chi_{s,t}$ in advance.

Another benefit of this region-pair method is that it can capture generic patterns for regions rather than specific patterns for nodes.

An example pattern could be that many users drive from a residential area to a shopping area. Since these drivers live in different apartment buildings, their starting points could be different, resulting in many dispersed entries in the node-pair frequency table χ . In contrast, they contribute to the same entry in the region-pair frequency table $\hat{\chi}$.

4.3 Benchmarking $E_{s,t}$ of Shortest Path API

In our caching system (see Figure 3), the shortest path API is invoked when there is a cache miss. We study how to capture the expense $E_{s,t}$ of query $Q_{s,t}$ in this section.

Recall that the cache can be placed at a proxy or a server. For the proxy scenario, the shortest path API triggers the issue of a query message to the server. The cost is dominated by the communication round-trip time, which remains the same for all queries. Thus, we define the expense $E_{s,t}$ of query $Q_{s,t}$ in this scenario as:

$$E_{s,t}(\text{proxy}) = 1 \quad (7)$$

Our subsequent discussion focuses on the server scenario. Let \mathcal{ALG} be a shortest path algorithm to be called by the shortest path API. We denote the running time of \mathcal{ALG} for query $Q_{s,t}$ as the expense $E_{s,t}(\mathcal{ALG})$.

In the following, we present a generic technique to estimate $E_{s,t}(\mathcal{ALG})$; it is applicable to any algorithm \mathcal{ALG} and to arbitrary graph topology. To our best knowledge, we are the first to explore this issue. There exists work on shortest path distance estimation [17] but no work on estimating the running time of arbitrary algorithm \mathcal{ALG} . Even for the existing shortest path indexes [8, 10, 11, 18, 19], only their worst-case query times have been analyzed. However, they cannot be used to estimate the running time for a specific query $Q_{s,t}$.

A brute-force approach is to precompute $E_{s,t}(\mathcal{ALG})$ by running \mathcal{ALG} for every pair of source node v_s and target node v_t . These values can be stored in a table, like the one in Figure 5b. However, this approach is prohibitively expensive as it requires running \mathcal{ALG} for $|V|^2$ times.

We plan to design an estimation technique that incurs small pre-computation overhead. Intuitively, the expense $E_{s,t}$ is strongly correlated with the distance of the shortest path $P_{s,t}$. Short-range queries are expected to incur small $E_{s,t}$ whereas long-range queries should produce high $E_{s,t}$. Our idea is to classify queries based on distances and then estimate the expense of a query according to its category.

Estimation structures: To enable estimation, we need to build two data structures: (i) a distance estimator, and (ii) an expense histogram.

The *distance estimator* aims at estimating the shortest path distance of a query $Q_{s,t}$. We simply adopt the landmark-based estimator [17] as the distance estimator. It requires selecting a set U of nodes as landmarks and precomputing the distances from each landmark node to every node in the graph. This incurs $O(|U||V|)$ storage space and $O(|U||E| \log |E|)$ construction time. Ref. [17] suggests that $|U| = 20$ is sufficient for accurate distance estimation. Figure 8a shows an example with two landmark nodes, i.e., $U = \{v_3, v_5\}$, together with their distances $d(u_j, v_i)$ to other nodes.

We propose to build an *expense histogram* for recording the average expense of queries with respect to their distances, as illustrated in Figure 8b. In general, the histogram consists of H categories of distances. Then, we execute the algorithm \mathcal{ALG} on a sample of S random queries to obtain their expenses, and update the corresponding buckets in the histogram. This histogram requires $O(H)$ storage space and $S \cdot O(\mathcal{ALG})$ construction time. We recommend

to use $H = 10$ and $S = 100$ in our experiments.

According to our experimental results, the construction of these structures takes less than 1 minute on typical road networks.

$d(u_j, v_i)$	v_1	v_2	v_3	v_4	v_5	v_6	v_7	v_8
$u_1 : v_3$	3	1	0	6	15	17	20	22
$u_2 : v_5$	18	16	15	9	0	2	5	7

(a) distance estimator

d	0-5	5-10	10-15	15-20	20-25
$E(d)$	1.2	2.5	3.8	5.3	6.8

(b) expense histogram

Figure 8: Example of estimating expense $\chi_{s,t}$

Estimation process: With the above structures, the value of $E_{s,t}$ can be estimated in two steps. First, we apply the distance estimator of Ref. [17] and estimate the shortest path distance of $P_{s,t}$ as: $\min_{i=1..|U|} d(u_j, v_s) + d(u_j, v_t)$. This step takes $O(|U|)$ time only. Second, we lookup the expense histogram and return the expense in the corresponding bucket as the estimated expense $E_{s,t}$.

Let's take the estimation of $E_{1,4}$ as an example. Using the distance estimator in Figure 8a, we estimate the shortest path distance of $P_{1,4}$ as: $\min\{3 + 6, 18 + 9\} = 9$. We then lookup the expense histogram in Figure 8b and thus estimate $E_{1,4}$ as 2.5.

4.4 Cache Selection Algorithm

Like static caching methods [1, 2, 15, 16], we exploit the query log \mathcal{QL} to identify promising results to be placed in the cache Ψ . Note that each query $Q_{a,b} \in \mathcal{QL}$ has a corresponding path result $P_{a,b}$. This section presents a cache selection algorithm for selecting these paths into cache Ψ such that the total cache benefit $\gamma(\Psi)$ is maximized, and the cache size $|\Psi|$ is bounded by a cache budget B .

In web search caching, Ozcan et al. [15] has proposed a greedy algorithm to fill the cache with results. Thus, we also adopt the greedy approach to solve our problem. Nevertheless, there are some challenges on applying the greedy approach to our problem.

Challenges of a greedy approach: It is tempting to fill the cache with paths by using a greedy approach. This approach would: (i) compute the benefit-per-size $\bar{\gamma}(P_{a,b})$ for each path $P_{a,b}$ and then (ii) iteratively fill the cache with the items having the highest $\bar{\gamma}(P_{a,b})$. Unfortunately, this approach does not necessarily produce a cache with high benefit.

As an example, we consider the graph in Figure 6 and the query log \mathcal{QL} in Table 4. Note that the result paths of queries of \mathcal{QL} are: $P_{1,6}, P_{2,7}, P_{1,4}, P_{4,8}, P_{2,5}, P_{3,6}$. To make the benefit calculation readable, we assume that $E_{s,t} = 1$ for each pair, and use the values of $\chi_{s,t}$ in Figure 5a. In this greedy approach, we first compute the benefit-per-size of each path above. For example, $P_{1,6}$ can answer five queries $Q_{3,6}, Q_{1,6}, Q_{1,4}, Q_{3,6}, Q_{3,6}$ in \mathcal{QL} , and its size $|P_{1,6}|$ is 5, so its benefit-per-size is: $\bar{\gamma}(P_{1,6}) = 5/5$. Since $P_{3,6}$ has a size 4 and it can answer three queries $Q_{3,6}, Q_{3,6}, Q_{3,6}$ in \mathcal{QL} , its benefit-per-size is: $\bar{\gamma}(P_{3,6}) = 3/4$. Repeating this process for the other paths, we obtain: $\bar{\gamma}(P_{1,4}) = 1/3, \bar{\gamma}(P_{1,6}) = 5/5, \bar{\gamma}(P_{2,5}) = 1/4, \bar{\gamma}(P_{2,7}) = 2/5, \bar{\gamma}(P_{3,6}) = 3/4, \bar{\gamma}(P_{4,8}) = 1/4$. Given the cache budget $B = 10$, the greedy approach would first pick $P_{1,6}$ and then pick $P_{3,6}$. Thus, we obtain the cache $\Psi = \{P_{1,6}, P_{3,6}\}$ with the size 9 (i.e., total number of nodes in the cache). No more paths can be inserted into the cache as it is full.

The problem with the greedy approach is that it ignores the existing content of the cache Ψ when it chooses a path $P_{a,b}$. Observe that, if many queries that can be answered by path $P_{a,b}$ can already be answered by some existing path in Ψ ; then it is not worthwhile to include $P_{a,b}$ into Ψ . In the above example, the greedy approach plans to pick the path $P_{3,6}$ next time after the path $P_{1,6}$ has been selected into the cache. Although path $P_{3,6}$ can answer three queries

$Q_{3,6}, Q_{3,6}, Q_{3,6}$ in \mathcal{QL} , all those queries can already be answered by the path $P_{1,6}$ in the cache. There is no need for path $P_{3,6}$ but the greedy approach will still pick it.

A revised greedy approach: To tackle the above issue, we study a notion that expresses the benefit of a path $P_{a,b}$ in terms of the queries that can be answered by $P_{a,b}$ but not by the existing paths in the cache Ψ .

DEFINITION 6. Incremental benefit-per-size of path $P_{a,b}$.

Given a shortest path $P_{a,b}$, its incremental benefit-per-size $\Delta\bar{\gamma}(P_{a,b}, \Psi)$ with respect to the cache Ψ , is defined as the additional benefit of placing $P_{a,b}$ into Ψ , per the size of $P_{a,b}$:

$$\begin{aligned}\Delta\bar{\gamma}(P_{a,b}, \Psi) &= \frac{\gamma(\Psi \cup \{P_{a,b}\}) - \gamma(\Psi)}{|P_{a,b}|} \\ &= \sum_{P_{s,t} \in \mathcal{U}(P_{a,b}) - \mathcal{U}(\Psi)} \frac{\chi_{s,t} \cdot E_{s,t}}{|P_{a,b}|}\end{aligned}\quad (8)$$

We propose a revised greedy algorithm that proceeds in rounds. A cache Ψ is initially empty. In each round, the algorithm computes the incremental benefit $\Delta\bar{\gamma}(P_{a,b}, \Psi)$ of each path $P_{a,b}$ with respect to the cache Ψ (with the selected paths so far). Then, the algorithm picks the path with the highest $\Delta\bar{\gamma}$ value and inserts it into Ψ . These rounds are repeated until the cache Ψ becomes full (i.e., reaching its budget \mathcal{B}).

We continue with the above running example and show the steps of this revised greedy algorithm in Table 5. In the first round, the cache Ψ is empty, so the incremental benefit $\Delta\bar{\gamma}(P_{a,b}, \Psi)$ of each path $P_{a,b}$ equals to its benefit $\bar{\gamma}(P_{a,b})$. From the previous example, we obtain: $\Delta\bar{\gamma}(P_{1,4}) = 1/3$, $\Delta\bar{\gamma}(P_{1,6}) = 5/5$, $\Delta\bar{\gamma}(P_{2,5}) = 1/4$, $\Delta\bar{\gamma}(P_{2,7}) = 2/5$, $\Delta\bar{\gamma}(P_{3,6}) = 3/4$, $\Delta\bar{\gamma}(P_{4,8}) = 1/4$. After choosing the path $P_{1,6}$ with the highest $\Delta\bar{\gamma}$ value, the cache becomes: $\Psi = \{P_{1,6}\}$. In the second round, we consider the cache when computing the $\Delta\bar{\gamma}$ value of a path. For the path $P_{3,6}$, all queries that can be answered by it can also be answered by the path $P_{1,6}$ in the cache. Thus, the $\Delta\bar{\gamma}$ value of $P_{3,6}$ is: $\Delta\bar{\gamma}(P_{3,6}) = 0$. Continuing with other queries, we obtain: $\Delta\bar{\gamma}(P_{1,4}) = 0$, $\Delta\bar{\gamma}(P_{1,6}) = 0$, $\Delta\bar{\gamma}(P_{2,5}) = 1/4$, $\Delta\bar{\gamma}(P_{2,7}) = 2/5$, $\Delta\bar{\gamma}(P_{3,6}) = 0$, $\Delta\bar{\gamma}(P_{4,8}) = 1/4$. The path $P_{2,7}$ with the highest $\Delta\bar{\gamma}$ value is chosen and then the cache becomes: $\Psi = \{P_{1,6}, P_{2,7}\}$. The total benefit of the cache $\gamma(\Psi)$ is 7. Now, the cache is full and no more paths can be inserted into the cache.

Round	Path						Cache Ψ	
	$P_{1,4}$	$P_{1,6}$	$P_{2,5}$	$P_{2,7}$	$P_{3,6}$	$P_{4,8}$	before round	after round
1	1/3	5/5	1/4	2/5	3/4	1/4	empty	$P_{1,6}$
2	0	0	1/4	2/5	0	1/4	$P_{1,6}$	$P_{1,6}, P_{2,7}$

Table 5: Incremental benefits of paths in our greedy algorithm (boxed values indicate the selected paths)

Our algorithm and time complexity: Algorithm 1 shows the pseudo-code of our revised greedy algorithm. It takes as input the graph $G(V, E)$, the cache budget \mathcal{B} , and the query log \mathcal{QL} . The cache budget \mathcal{B} denotes the capacity of the cache in terms of the number of nodes. The statistics of query frequency χ and query expense E are required for computing the incremental benefit of a path. The initialization phase corresponds to Lines 1–5. The cache Ψ is initially empty. A max-heap H is employed to organize result paths in descending order of their $\Delta\bar{\gamma}$ values. For each query $Q_{a,b}$ in the query log \mathcal{QL} , we retrieve its result path $P_{a,b}$, compute its $\Delta\bar{\gamma}$ value as $\Delta\bar{\gamma}(P_{a,b}, \Psi)$, and then insert $P_{a,b}$ into H .

Algorithm 1 Revised-Greedy(Graph $G(V, E)$, Cache budget \mathcal{B} , Query log \mathcal{QL} , Frequency χ , Expense E)

```

1: Cache  $\Psi \leftarrow$  create a new cache;
2:  $H \leftarrow$  create a new max-heap;           ▷ storing result paths
3: for each  $Q_{a,b} \in \mathcal{QL}$  do
4:    $P_{a,b}, \Delta\bar{\gamma} \leftarrow \Delta\bar{\gamma}(P_{a,b}, \Psi)$ ;       ▷ compute using  $\chi$  and  $E$ 
5:   insert  $P_{a,b}$  into  $H$ ;
6: while  $|\Psi| \leq \mathcal{B}$  and  $|H| > 0$  do
7:    $P_{a',b'} \leftarrow H.pop()$ ;                 ▷ potential best path
8:    $P_{a',b'}, \Delta\bar{\gamma} \leftarrow \Delta\bar{\gamma}(P_{a',b'}, \Psi)$ ;   ▷ update  $\Delta\bar{\gamma}$  value
9:   if  $P_{a',b'}, \Delta\bar{\gamma} \geq \text{top } \Delta\bar{\gamma} \text{ of } H$  then   ▷ actual best path
10:    if  $\mathcal{B} - |\Psi| \geq |P_{a',b'}|$  then           ▷ enough space
11:      insert  $P_{a',b'}$  into  $\Psi$ ;
12:    else                                       ▷ not the best path
13:      insert  $P_{a',b'}$  into  $H$ ;
14: return  $\Psi$ ;
```

The algorithm incorporates an optimization to reduce the number of incremental benefit computations in each round (i.e., the loop of Lines 6–13). First, the path $P_{a',b'}$ with the highest $\Delta\bar{\gamma}$ value is selected from H (Line 7) and its current $\Delta\bar{\gamma}$ value is computed (Line 8). According to Lemma 2, the $\Delta\bar{\gamma}$ value of a path $P_{a,b}$ in H , which was computed in some previous round, serves as an upper bound of its $\Delta\bar{\gamma}$ value in the current round. If $\Delta\bar{\gamma}(P_{a',b'}, \Psi)$ is above the top key of H (Line 9), then we can safely conclude that $P_{a',b'}$ is superior to all paths in H , without having to compute their exact $\Delta\bar{\gamma}$ values. Then, we insert the path $P_{a',b'}$ into the cache Ψ when it has sufficient remaining space $\mathcal{B} - |\Psi|$. In case $\Delta\bar{\gamma}(P_{a',b'}, \Psi)$ is smaller than the top key of H , we insert $P_{a',b'}$ back to H . Eventually, H becomes empty, the loop terminates, and the cache Ψ is returned.

LEMMA 2. $\Delta\bar{\gamma}$ is a decreasing function of round i .

Let Ψ_i be the cache just before the i -th round of the algorithm. It holds that: $\Delta\bar{\gamma}(P_{a,b}, \Psi_i) \geq \Delta\bar{\gamma}(P_{a,b}, \Psi_{i+1})$.

PROOF. All paths in Ψ_i must also be in Ψ_{i+1} , so we have: $\Psi_i \subseteq \Psi_{i+1}$. By Equation 3, we derive: $\mathcal{U}(\Psi_i) \subseteq \mathcal{U}(\Psi_{i+1})$ and then obtain: $\mathcal{U}(P_{a,b}) - \mathcal{U}(\Psi_i) \supseteq \mathcal{U}(P_{a,b}) - \mathcal{U}(\Psi_{i+1})$. By Definition 6, we have $\Delta\bar{\gamma}(P_{a,b}, \Psi) = \sum_{P_{s,t} \in \mathcal{U}(P_{a,b}) - \mathcal{U}(\Psi)} \chi_{s,t} \cdot E_{s,t} / |P_{a,b}|$. Combining the above facts, we get: $\Delta\bar{\gamma}(P_{a,b}, \Psi_i) \geq \Delta\bar{\gamma}(P_{a,b}, \Psi_{i+1})$. \square

We proceed to illustrate the power of the above optimization. Let's consider the second round of Table 5. Without the optimization, we must recompute the $\Delta\bar{\gamma}$ values of 5 paths $P_{1,4}, P_{2,5}, P_{2,7}, P_{3,6}, P_{4,8}$, before determining the path with the highest $\Delta\bar{\gamma}$ value. Using the optimization, we just need to pop the paths $P_{3,6}, P_{2,7}$ from the heap H and recompute their $\Delta\bar{\gamma}$ values. For the other paths (e.g., $P_{1,4}, P_{2,5}, P_{4,8}$), their upper bound $\Delta\bar{\gamma}$ values (1/3, 1/4, 1/4, from the first round) are smaller than the current $\Delta\bar{\gamma}$ value of $P_{2,7}$ (2/5), thus we do not need to recompute their current $\Delta\bar{\gamma}$ values.

We then analyze the time complexity of Algorithm 1, without using the optimization. Let $|\mathcal{QL}|$ be the number of result paths for queries in \mathcal{QL} . Let $|P|$ be the average size of above result paths. The number of paths in the cache is $\mathcal{B}/|P|$ so the algorithm completes in $\mathcal{B}/|P|$ rounds. In each round, we need to process $|\mathcal{QL}|$ result paths and recompute their $\Delta\bar{\gamma}$ values. Computing the $\Delta\bar{\gamma}$ value of a path $P_{a,b}$ requires examining each subpath of $P_{a,b}$ (see Definition 6); this takes $O(|P|^2)$ time as they are $O(|P|^2)$ subpaths in a path. Multiplying the above terms together, the time

the interval $\Psi_{1,3}$.

Prefix path ID compression: This technique first identifies inverted lists that share the same prefix, and then expresses an inverted list by using the other inverted list as prefix.

Let's consider the original inverted lists in Figure 11a. The inverted list of v_1 is the prefix of the inverted list of v_3 . Figure 11c shows the compressed inverted lists produced by this compression. In the compressed inverted list of v_3 , it suffices to store path IDs (e.g., Ψ_3) that do not appear in its prefix. The remaining path IDs of v_3 can be retrieved by following the parent (v_1) of its inverted list.

Both compression techniques can be combined together in order to achieve a higher compression power.

			content	parent
v_1	Ψ_1, Ψ_2	v_1	Ψ_1-2	
v_2	Ψ_2, Ψ_3	v_2	Ψ_2-3	
v_3	Ψ_1, Ψ_2, Ψ_3	v_3	Ψ_1-3	
v_4	Ψ_1, Ψ_3	v_4	Ψ_1, Ψ_3	
v_5	Ψ_3	v_5	Ψ_3	
		v_1	Ψ_1, Ψ_2	NIL
		v_2
		v_3	Ψ_3	v_1
		v_4
		v_5

(a) original

(b) interval compressed

(c) prefix compressed

Figure 11: Compressed inverted lists

6. EXPERIMENTS

In this section, we evaluate the performance of our methods with our competitors on real datasets. We have implemented two variants of our methods (SPC and SPC*). They share the same techniques in Section 4, and only differ in their cache structures: (i) SPC uses a path array cache (Section 5.1), and (ii) SPC* uses the compressed graph cache (Sections 5.2,5.3). Our competitors are LRU (a dynamic caching method) and HQF (a static caching method). They have been introduced in Section 3.3. All the above methods are written in C++. We conduct our experiments on a machine with Intel i7 3.4GHz CPU.

We will evaluate the above methods for a cache located at the proxy, and for a cache located at the server. For the proxy scenario, the performance measure is the hit ratio. For the server scenario, the performance measures are: (i) the total running time of the server, and (ii) the total number of road network nodes visited.

6.1 Experimental Setting

Which historical query datasets do we have, and what is their size and origin.

- Aalborg: Query workload from and around the Danish city of Aalborg
- Beijing: Query workload from Beijing

Due to the double-blind policy of SIGMOD, we do not disclose the

according to XXX I

Our 3 caching methods, SPC, HQF, and LRU, share a number of common settings which, unless stated otherwise, will be set to their default values: The number of levels in the kD-tree are (i.e. 16.384 regions). We will use the list cache representation as the default cache representation, where each vertex use one byte. The default cache size is set to 160.000 bytes.

Except for tests on synthetic data, the size of training and test query datasets are given already. The size of each dataset, as well as the map they are captured on, is given in table 6.

We do not give a default size for the query-datasets Maps, as we will execute all of our tests, described in section 6.2 and 6.3, for

Dataset	# Queries	# Nodes in Map
City X	3,286	136,347
City Y	12,956	79,603

Table 6: Size of query datasets. Training and testing datasets are of equal size.

each dataset.

GPS trajectories

6.2 Caching in the Proxy Scenario

On the proxy [REF FIGURE/SECTION] we won't be doing SP calculations, so the most important performance measure is the cache hit ratio.

For both datasets we vary the cache size and number of levels to show the impact on the cache hit ratio. We have implemented a number of optimizations to the cache storage [REF PREV SEC. NEEDED] and we will show their impact on the cache hit ratio.

Effect of the cache size:

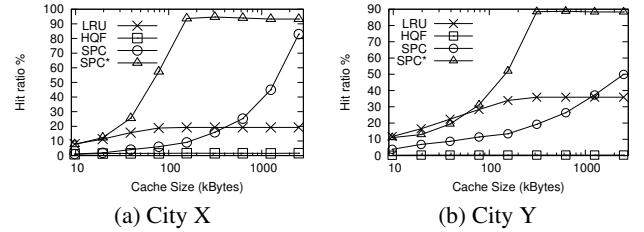


Figure 12: Hit ratio vs. cache size

- We can see SPC^* performs much better than competitors
- LRU perform well on small cache sizes, but the cache hit ratio quickly level out and stop growing, thus LRU is not scalable.
- HQF performance is very consistent, with the hit ratio varying only slightly. The hit ratio does however stay at just over 1% for City X experiments and around 0.2% for City Y experiments. It is not usable
- We observe that SPC initially does not grow as fast as LRU, but the growth continues as we increase cache size, unlike for LRU.
- SPC^* grows faster than LRU and gets better cache hit ratio at all cache sizes than LRU and SPC, except for the smallest cache size.
- The test on the City X vs. City Y sets show that SPC performs much better on the City X set, where as LRU performs about 50% worse. This can be explained by the fact that both methods rely on users behaving consistently over time, but where SPC relies on global consistency in user behavior, LRU need local consistency too. This makes SPC/ SPC^* more robust and we can expect them to always outperform LRU.

Effect of the kD-tree level:

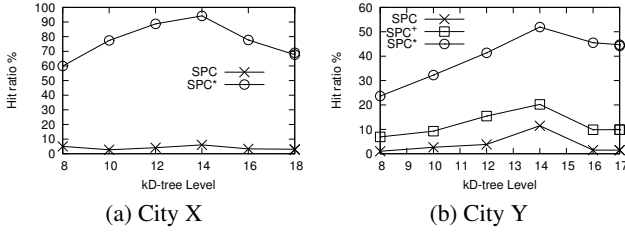


Figure 13: Hit ratio vs. Levels

$ U $	City X	City Y
5	35.5	33.3
10	25.4	28.2
20	22.2	23.9
40	19.4	23.0
80	19.9	21.9

(a) varying number of landmark $|U|$

S	City X	City Y
25	23.4	28.8
50	20.7	28.4
100	22.2	23.9
200	20.4	22.7
400	21.1	21.3

(b) varying sample size S

Table 7: Average error percentage of cost estimation

6.3 Caching in the Server Scenario

On the server [REF TO FIGURE] our aim is slightly different than on the Proxy, as we also have to consider that there may be some paths that are so small that it may be cheaper to simply recalculate the result, instead of caching it.

We use the Dijkstra's algorithm as the shortest path algorithm.

Estimation of shortest path running cost:

Effect of the kD-tree level:

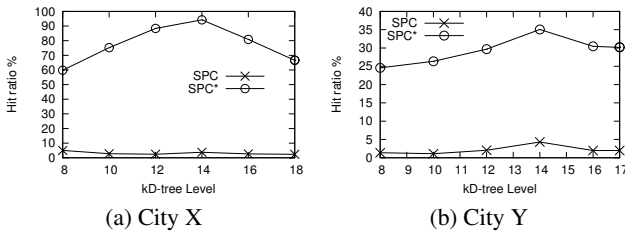


Figure 14: Hit ratio vs. Levels

Effect of the cache size:

7. CONCLUSION

We have studied ...

8. REFERENCES

- [1] I. S. Altingövde, R. Ozcan, and Ö. Ulusoy. A cost-aware strategy for query result caching in web search engines. In *ECIR*, pages 628–636, 2009.
- [2] R. A. Baeza-Yates, A. Gionis, F. Junqueira, V. Murdock, V. Plachouras, and F. Silvestri. The impact of caching on search engines. In *SIGIR*, pages 183–190, 2007.
- [3] R. A. Baeza-Yates and F. Saint-Jean. A three level search engine index based in query log distribution. In *SPIRE*, pages 56–65, 2003.
- [4] T. H. Cormen, C. E. Leiserson, and C. Stein. *Introduction to Algorithms*. MIT Press, 3rd edition, 2009.

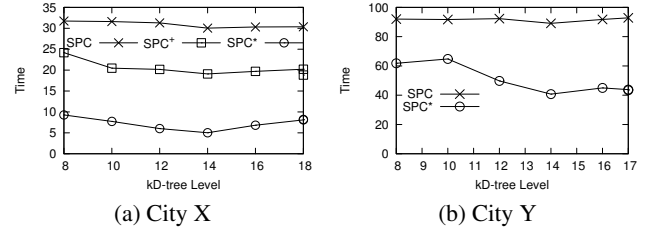


Figure 15: Runtime vs. Levels

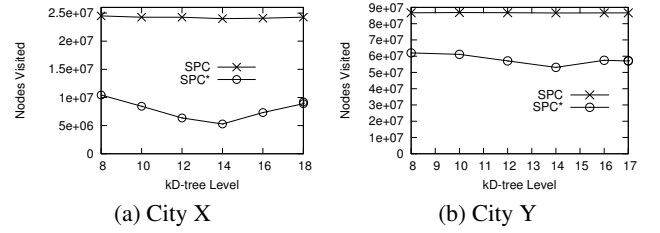


Figure 16: Nodes Visited vs. Levels

- [5] S. Dar, M. J. Franklin, B. T. Jónsson, D. Srivastava, and M. Tan. Semantic data caching and replacement. In *VLDB*, pages 330–341, 1996.
- [6] Q. Gan and T. Suel. Improved techniques for result caching in web search engines. In *WWW*, pages 431–440, 2009.
- [7] M. C. González, C. A. Hidalgo, and A. L. Barabási. Understanding Individual Human Mobility Patterns. *Nature*, 453(7196):779–782, 2008.
- [8] H. Hu, D. L. Lee, and V. C. S. Lee. Distance Indexing on Road Networks. In *VLDB*, pages 894–905, 2006.
- [9] H. Hu, J. Xu, W. S. Wong, B. Zheng, D. L. Lee, and W.-C. Lee. Proactive caching for spatial queries in mobile environments. In *ICDE*, pages 403–414, 2005.
- [10] S. Jung and S. Pramanik. An Efficient Path Computation Model for Hierarchically Structured Topographical Road Maps. *IEEE TKDE*, 14(5):1029–1046, 2002.
- [11] H.-P. Kriegel, P. Kröger, M. Renz, and T. Schmidt. Hierarchical graph embedding for efficient query processing in very large traffic networks. In *SSDBM*, pages 150–167, 2008.
- [12] K. Lee, W.-C. Lee, B. Zheng, and J. Xu. Caching complementary space for location-based services. In *EDBT*, pages 1020–1038, 2006.
- [13] X. Long and T. Suel. Three-level caching for efficient query

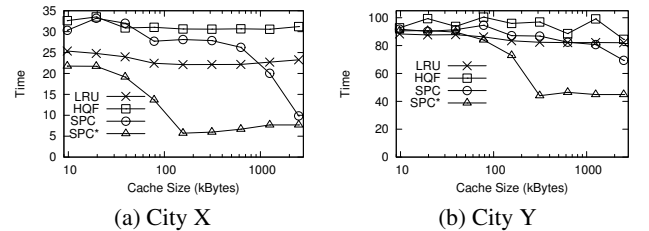


Figure 17: Runtime vs. Cache Size

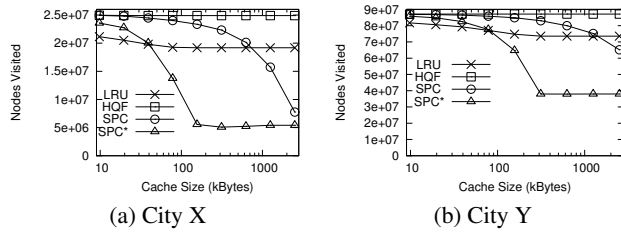


Figure 18: Nodes Visited vs. Cache Size

processing in large web search engines. In *WWW*, pages 257–266, 2005.

- [14] E. P. Markatos. On caching search engine query results. *Computer Communications*, 24(2):137–143, 2001.
- [15] R. Ozcan, I. S. Altingövde, B. B. Cambazoglu, F. P. Junqueira, and Ö. Ulusoy. A five-level static cache architecture for web search engines. *Information Processing & Management*, (0):–, 2011.
- [16] R. Ozcan, I. S. Altingövde, and Ö. Ulusoy. Static query result caching revisited. In *WWW*, pages 1169–1170, 2008.
- [17] M. Potamias, F. Bonchi, C. Castillo, and A. Gionis. Fast shortest path distance estimation in large networks. In *CIKM*, pages 867–876, 2009.
- [18] H. Samet, J. Sankaranarayanan, and H. Alborzi. Scalable network distance browsing in spatial databases. In *SIGMOD*, pages 43–54, 2008.
- [19] F. Wei. TEDI: Efficient Shortest Path Query Answering on Graphs. In *SIGMOD*, pages 99–110, 2010.
- [20] B. Zheng and D. L. Lee. Semantic caching in location-dependent query processing. In *SSTD*, pages 97–116, 2001.