

\mathcal{MD} -HBase: A Scalable Multi-dimensional Data Infrastructure for Location Aware Services

Shoji Nishimura^{*§}

Sudipto Das[†]

Divyakant Agrawal[†]

Amr El Abbadi[†]

**Service Platforms Research Laboratories
NEC Corporation
Kawasaki, Kanagawa 211-8666, Japan
s-nishimura@bk.jp.nec.com*

*†Department of Computer Science
University of California, Santa Barbara
Santa Barbara, CA 93106-5110, USA
{sudipto, agrawal, amr}@cs.ucsb.edu*

Abstract—The ubiquity of location enabled devices has resulted in a wide proliferation of location based applications and services. To handle the growing scale, database management systems driving such location based services (LBS) must cope with high insert rates for location updates of millions of devices, while supporting efficient real-time analysis on latest location. Traditional DBMSs, equipped with multi-dimensional index structures, can efficiently handle spatio-temporal data. However, popular open-source relational database systems are overwhelmed by the high insertion rates, real-time querying requirements, and terabytes of data that these systems must handle. On the other hand, Key-value stores can effectively support large scale operation, but do not natively support multi-attribute accesses needed to support the rich querying functionality essential for the LBSs. We present \mathcal{MD} -HBase, a scalable data management system for LBSs that bridges this gap between scale and functionality. Our approach leverages a multi-dimensional index structure layered over a Key-value store. The underlying Key-value store allows the system to sustain high insert throughput and large data volumes, while ensuring fault-tolerance, and high availability. On the other hand, the index layer allows efficient multi-dimensional query processing. We present the design of \mathcal{MD} -HBase that builds two standard index structures—the K-d tree and the Quad tree—over a range partitioned Key-value store. Our prototype implementation using HBase, a standard open-source Key-value store, can handle hundreds of thousands of inserts per second using a modest 16 node cluster, while efficiently processing multi-dimensional range queries and nearest neighbor queries in real-time with response times as low as hundreds of milliseconds.

Index Terms—location based services; key value stores; multi-dimensional data; real time analysis;

I. INTRODUCTION

The last few years have witnessed a significant increase in hand-held devices becoming location aware with the potential to continuously report up-to-date location information of their users. This has led to a large number of location based services (LBS) which customize a user's experience based on location. Some applications—such as customized recommendations and advertisements based on a user's current location and history—have immediate economic incentives, while some other applications—such as location based social networking or location aware gaming—enrich the user's experience in general. With major wireless providers serving hundreds of

millions of subscribers [1], millions of devices registering their location updates continuously is quite common. Database management systems (DBMS) driving these location based services must therefore handle millions of location updates per minute while answering near real time analysis and statistical queries that drive the different recommendation and personalization services.

Location data is inherently multi-dimensional, minimally including a user id, a latitude, a longitude, and a time stamp. A rich literature of multi-dimensional indexing techniques—for instance, K-d trees [2], Quad trees [3] and R-trees [4]—have empowered relational databases (RDBMS) to efficiently process multi-dimensional data. However, the major challenge posed by these location based services is in scaling the systems to sustain the high throughput of location updates and analyzing huge volumes of data to glean intelligence. For instance, if we consider only the insert throughput, a MySQL installation running on a commodity server becomes a bottleneck at loads of tens of thousands of inserts per second; performance is further impacted adversely when answering queries concurrently. Advanced designs, such as staging the database, defining views, database and application partitioning to scale out, a distributed cache, or moving to commercial systems, will increase the throughput; however, such optimizations are expensive to design, deploy, and maintain.

On the other hand, Key-value stores, both in-house systems such as Bigtable [5] and their open source counterparts like HBase [6], have proven to scale to millions of updates while being fault-tolerant and highly available. However, Key-value stores do not natively support efficient multi-attribute access, a key requirement for the rich functionality needed to support LBSs. In the absence of any filtering mechanism for secondary attribute accesses, such queries resort to full scan of the entire data. MapReduce [7] style processing is therefore a commonly used approach for analysis on Key-value stores. Even though the MapReduce framework provides abundant parallelism, a full scan is wasteful, especially when the selectivity of the queries is high. Moreover, many applications require near real-time query processing based on a user's current location. Therefore, query results based on a user's stale location is often useless. As a result, a design for batched query processing on data periodically imported into a data warehouse

[§]Work done as a visiting researcher at UCSB.

is inappropriate for the real-time analysis requirement.

RDBMSs

provide rich querying support for multi-dimensional data but are not scalable, while Key-value stores can scale but cannot handle

multi-dimensional data efficiently. Our solution, called *MD-HBase*, bridges this gap by layering a multi-dimensional index over a Key-value store to leverage the best of both worlds.¹ We use linearization techniques such as Z-ordering [8] to transform multi-dimensional location information into a one dimensional space and use a range partitioned Key-value store (HBase [6] in our implementation) as the storage back end. Figure 1 illustrates *MD-HBase*'s architecture showing the index layer and the data storage layer. We show how this design allows standard and proven multi-dimensional index structures, such as K-d trees and Quad trees, to be layered on top of the Key-value stores with minimal changes to the underlying store and negligible effect on the operation of the Key-value store. The underlying Key-value store provides the ability to sustain a high insert throughput and large data volumes, while ensuring fault-tolerance and high availability. The overlaid index layer allows efficient real-time processing of multi-dimensional range and nearest neighbor queries that comprise the basic data analysis primitives for location based applications. We evaluate different implementations of the data storage layer in the Key-value store and evaluate the trade-offs associated with these different implementations. In our experiments, *MD-HBase* achieved more than 200K inserts per second on a modest cluster spanning 16 nodes, while supporting real-time range and nearest neighbor queries with response times less than one second. Assuming devices reporting one location update per minute, this small cluster can handle updates from 10 – 15 million devices while providing between one to two orders of magnitude improvement over a MapReduce or Z-ordering based implementation for query processing. Moreover, our design does not introduce any scalability bottlenecks, thus allowing the implementation to scale with the underlying Key-value data store.

Contributions.

- We propose the design of *MD-HBase* that uses linearization to implement a scalable multi-dimensional index structure layered over a range-partitioned Key-value store.
- We demonstrate how this design can be used to implement a K-d tree and a Quad tree, two standard multi-dimensional index structures.
- We present three alternative implementations of the storage layer in the Key-value store and evaluate the tradeoffs

¹The name *MD-HBase* signifies adding multi-dimensional data processing capabilities to HBase, a range partitioned Key-value store.

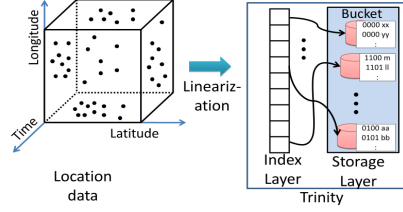


Fig. 1. Architecture of *MD-HBase*.

associated with each implementation.

- We provide a detailed evaluation of our prototype using synthetically generated location data and analyze *MD-HBase*'s scalability and efficiency.

Organization. Section II provides background on location based applications, linearization techniques, and multi-dimensional index structures. Section III describes the design and implementation of the index layer and Section IV describes the implementation of the data storage layer. Section V presents a detailed evaluation of *MD-HBase* by comparing the different design choices. Section VI surveys the related literature and Section VII concludes the paper.

II. BACKGROUND

A. Location based applications

Location data points are multi-dimensional with spatial attributes (e.g., longitude and latitude), a temporal attribute (e.g., timestamp), and an entity attribute (e.g., user's ID). Different applications use this information in a variety of ways. We provide two simple examples that illustrate the use of location data for providing location aware services.

APPLICATION 1: Location based advertisements and coupon distribution: Consider a restaurant chain, such as McDonalds, running a promotional discount for the next hour to advertise a new snack and wants to disseminate coupons to attract customers who are currently near any of their restaurants spread throughout the country. An LBS provider issues multi-dimensional range queries to determine all users within 5 miles from any restaurant in the chain and delivers a coupon to their respecting devices. Another approach to run a similar campaign with a limited budget is to limit the coupons to only the 100 users nearest to a restaurant location. In this case, the LBS provider issues nearest neighbors queries to determine the users. In either case, considering a countrywide (or worldwide) presence of this restaurant chain, the amount of data analyzed by such queries is huge.

APPLICATION 2: Location based social applications: Consider a social application that notifies a user of his/her friends who are currently nearby. The LBS provider in this case issues a range query around the user's current location and intersects the user's friend list with the results. Again, considering the scale of current social applications, an LBS provider must handle data for tens to hundreds of millions of its users to answer these queries for its users spread across a country.

Location information has two important characteristics. First, it is inherently skewed, both spatially and temporally. For instance, urban regions are more dense compared to rural regions, while business and school districts are dense during weekdays, while residential areas are dense during the night and weekends. Second, the time dimension is potentially unbounded and monotonically increasing. We later discuss how these characteristics influence many of the design choices.

B. Linearization Techniques

Linearization [9] is a method to transform multi-dimensional data points to a single dimension and is a key

aspect of the index layer in \mathcal{MD} -HBase. Linearization allows leveraging a single-dimensional database (a Key-value store in our case) for efficient multi-dimensional query processing. A space-filling curve [9] is one of the most popular approaches for linearization. A space filling curve visits all points in the multi-dimensional space in a systematic order. Z-ordering [8] is an example of a space filling curve. Z-ordering loosely preserves the locality of data-points in the multi-dimensional space and is also easy to implement.

Linearization alone is however not enough for efficient query processing; a linearization based multi-dimensional index layer, though simple in design, results in inefficient query processing. For instance, a range query on a linearized system is decomposed into several linearized sub-queries; however, a trade-off exists between the number of sub-queries and the number of false-positives. A reduction in the number of linearized sub-queries results in an increase in the overhead due to the large number of false-positives. On the other hand, eliminating false-positives results in a significant growth in the number of sub-queries. Furthermore, such a naïve technique is not robust to skew inherent in many real life applications.

C. Multi-dimensional Index Structures

The Quad tree [3] and the K-d tree [2] are two of the most popular multi-dimensional indexing structures. They split the multi-dimensional space recursively into **subspaces** in a systematic manner and organize these subspaces as a search tree. A Quad tree divides the n -dimensional space into 2^n subspaces along all dimensions whereas a K-d tree alternates the splitting of the dimensions. Each subspace has a maximum limit on the number of data points in it, beyond which the subspace is split. Two approaches are commonly used to split a subspace: a trie-based approach and a point-based approach [9]. The trie-based approach splits the space at the mid-point of a dimension, resulting in equal size splits; while the point-based technique splits the space by the median of data points, resulting in subspaces with equal number of data points. The trie-based approach is efficient to implement as it results in regular shaped subspaces. On the other hand, the point based approach is more robust to skew.

In addition to the performance issues, trie-based Quad trees and K-d trees have a property that allows them to be coupled with Z-ordering. A trie-based split of a Quad tree or a K-d tree results in subspaces where all z-values in any subspace are continuous. Figure 2 provides an illustration using a K-d tree for two dimensions; an example using a Quad tree is very similar. Different shades denote different subspaces and the dashed arrows denote the z-order traversal. As is evident, in any of the subspaces the Z-values are continuous. This observation forms the basis of the indexing layer of \mathcal{MD} -HBase. Compared to a naïve linearization based index structure or a B+-Tree index built on linearized values, K-d and Quad trees capture data

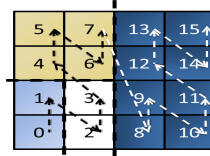


Fig. 2. Space splitting in a K-d tree.

distribution statistics. They partition the target space based on the data distribution and limit the number of data points in each subspace, resulting in more splits in hot regions. Furthermore, Quad trees and KD trees maintain the boundaries of subspaces in the original space. This allows efficient pruning of the space as well as reducing the number of false positive scans during query processing and therefore is more robust to skew. Moreover, when executing the best-first algorithm for kNN queries, the fastest practical algorithm, a B+ tree based index cannot be used due to the irregular sub-space shape. \mathcal{MD} -HBase therefore uses these multi-dimensional index structures instead of a simple single-dimension index based on linearization.

III. MULTI-DIMENSIONAL INDEX LAYER

We now present the design of the multi-dimensional index layer in \mathcal{MD} -HBase. Specifically, we show how standard index structures like K-d trees [2] and Quad trees [3] can be adapted to be layered on top of a Key-value store. The indexing layer assumes that the underlying data storage layer stores the items sorted by their key and range-partitions the key space. The keys correspond to the Z-value of the dimensions being indexed; for instance the location and timestamp. We use the trie-based approach for space splitting. The index partitions the space into conceptual **subspaces** that are in-turn mapped to a physical storage abstraction called **bucket**. The mapping between a conceptual subspace in the index layer and a bucket in the data storage layer can be one-to-one, many-to-one, or many-to-many depending on the implementation and requirements of the application. We develop a novel naming scheme for subspaces to simulate a trie-based K-d tree and a Quad tree. This naming scheme, called **longest common prefix naming**, has two important properties critical to efficient index maintenance and query processing; a description of the naming scheme and its properties are discussed below.

A. Longest Common Prefix Naming Scheme

If the multi-dimensional space is divided into equal sized subspaces and each dimension is enumerated using binary values, then the z-order of a given subspace is given by interleaving the bits from the different dimensions. Figure 3 illustrates this property. For example, the Z-value of the sub-space (00, 11) is represented as 0101. We name each subspace by the longest common prefix of the z-values of points contained in the subspace. Figure 4 provides an example for partitioning the space in a trie-based Quad tree and K-d tree. For example, consider a Quad tree built on the 2D space. The subspace at the top right of Figure 4(a), enclosed by a thick solid line, consists of z-values 1100, 1101, 1110, and 1111 with a longest common prefix of 11**;

11	0101	0111	1101	1111
10	0100	0110	1100	1110
01	0001	0011	1001	1011
00	0000	0010	1000	1010
	00	01	10	11

Fig. 3. Binary Z-ordering.

Similarly, the lower left subspace, enclosed by a broken line, only contains 0000, and is named 0000. Now consider the example of a K-d tree in Figure 4(b).

The subspace in the right half, enclosed by a solid line, is represented by the longest common prefix 1, while the subspace in the bottom left consisting of 0000 and 0001 is named as 000*. This naming scheme is similar to that of Prefix Hash Tree (PHT) [10], a variant of a distributed hash table.²

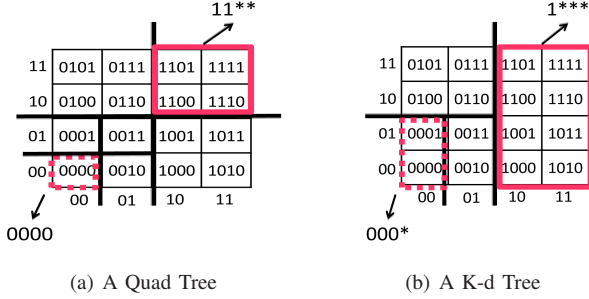


Fig. 4. The longest common prefix naming scheme

\mathcal{MD} -HBase leverages two important properties of this longest common prefix naming scheme. First, if subspace A encloses subspace B, the name of subspace A is a prefix of that of subspace B. For example, in Figure 4(a), a subspace which contains 0000–0011 encloses a subspace which contains only 0000. The former subspace’s name, 00**, is a prefix of the latter’s name 0000. Therefore, on a split, names of the new subspaces can be derived from the original subspace name by appending bits depending on which dimensions were split.

Second, the subspace name is enough to determine the region boundaries on all dimensions. This property derives from the fact that the z-values are computed by interleaving bits from the different dimensions and thus improves the pruning power of a range query. Determining the range boundary consists of the following two steps: (i) given the name, extract the bits corresponding to each dimension; and (ii) complete the extracted bits by appending 0s for the lower bound and 1s for the upper bound values, both bound values being inclusive. For instance, let us consider the Quad tree in Figure 4(a). The subspace at the top right enclosed by a solid lined box is named 11. Since this is a 2D space, the prefix for both the vertical and horizontal dimensions is 1. Therefore, using the rule for extension, the range for both dimensions is [10, 11]. Generalizing to the n -dimensional space, a Quad tree splits a space on all n dimensions resulting in 2^n subspaces. Therefore, the number of bits in the name of a subspace will be an exact multiple of the number of dimensions. This ensures that reconstructing the dimension values in step (i) will provide values for all dimensions. But since a K-d tree splits on alternate dimensions, the number of bits in the name is not guaranteed to be a multiple of the number of dimensions, in which case the prefix for different dimensions have different lengths. Considering the K-d tree example of Figure 4(b), for the lower left subspace named 000*, the prefix for the horizontal dimension is 00 while that of the vertical dimension is 0, resulting in the bounds for the horizontal dimension as [00, 00] and that of the vertical dimension as [00, 01].

²http://en.wikipedia.org/wiki/Distributed_hash_table

B. Index Layer

The index layer leverages the longest common prefix naming scheme to map multi-dimensional index structures into a single dimensional substrate. Figure 5 illustrates the mapping of a Quad tree partitioned space to the index layer; the mapping technique for K-d trees follows similarly. We now describe how the index layer can be used to efficiently implement some common operations and then discuss the space splitting algorithm for index maintenance.

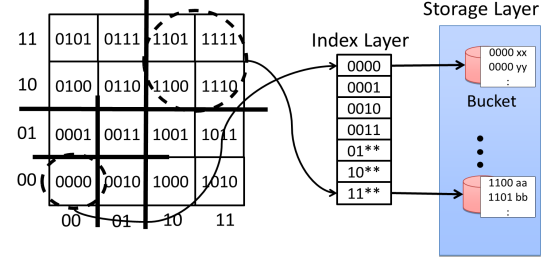


Fig. 5. Index layer mapping a trie-based Quad tree.

1) *Subspace Lookup and Point Queries*: Determining the subspace to which a point belongs forms the basis for point queries as well as for data insertion. Since the index layer comprises a sorted list of subspace names, determining the subspace to which a point belongs is efficient. Recall that the subspace name determines the bounds of the region that the subspace encloses. The search for the subspace finds the entry that has the maximum prefix matched with the z-value of the query point; this entry corresponds to the highest value smaller than the z-value of the query point. A prefix matching binary search, which substitutes exact matching comparison to prefix matching comparison for the termination condition of the binary search algorithm, is therefore sufficient. Algorithm 1 provides the algorithm for subspace lookup. To answer a point query, we first lookup the subspace corresponding to the z-value of the point. The point is then queried in the bucket to which the subspace maps.

Algorithm 1 Subspace Lookup

- 1: /* $\langle q \rangle$ be the query point. */
- 2: $Z_q \leftarrow \text{ComputeZ-value}(q)$
- 3: $Bkt_q \leftarrow \text{PrefixMatchingBinarySearch}(Z_p)$

2) *Insertion*: The algorithm to insert a new data point (shown in Algorithm 2) is similar to the point query algorithm. It first looks up the bucket corresponding to the subspace to which the point belongs, and then inserts the data point in the bucket. Since there is a maximum limit to the number of points a bucket can hold, the insertion algorithm checks the current size of the bucket to determine if a split is needed. The technique to split a subspace is explained in Section III-B5.

3) *Range Query*: A multi-dimensional range query is one of the most frequent queries for location based applications. Algorithm 3 provides the pseudo code for range query processing. Let $\langle q_l, q_h \rangle$ be the range for the query, q_l is the lower bound and q_h is the upper bound. The z-value of the

Algorithm 2 Insert a new location data point

```

1: /*  $\langle p \rangle$  be the new data point. */
2:  $Bkt_p \leftarrow \text{LookupSubspace}(p)$ 
3:  $\text{InsertToBucket}(Bkt_p, p)$ 
4: if ( $\text{Size}(Bkt_p) > \text{MaxBucketSize}$ ) then
5:    $\text{SplitSpace}(Bkt_p)$ 

```

lower bound determines the first subspace to scan. All subsequent subspaces until the one corresponding to the upper bound are potential candidate subspaces. Let \mathbb{S}_q denote the set of candidate subspaces. Since the z-order loosely preserves locality, some subspaces in \mathbb{S}_q might not be part of the range. For example, consider the Quad tree split shown in Figure 5. Consider the range query $\langle [01, 11], [10, 11] \rangle$. The z-value range for this query is $[0110, 1111]$ which results in \mathbb{S}_q equal to $\{01^{**}, 10^{**}, 11^{**}\}$. Since the query corresponds to only the top half of the space, the subspace named 10^{**} is a false positive. But such false positives are eliminated by a scan of the index. As the subspace name only is enough to determine the boundary of the region enclosed by the subspace, points in a subspace are scanned only if the range of the subspace intersects with the query range. This check is inexpensive and prunes out all the subspaces that are not relevant. For subspaces that are contained in the query, all points will be returned, while subspaces that only intersect with the query require further filtering. The steps for query processing are shown in Algorithm 3.

Algorithm 3 Range query

```

1: /*  $\langle q_l, q_h \rangle$  be the range for the query. */
2:  $Z_{low} \leftarrow \text{ComputeZ-value}(q_l)$ 
3:  $Z_{high} \leftarrow \text{ComputeZ-value}(q_h)$ 
4:  $\mathbb{S}_q \leftarrow \{Z_{low} \leq \text{SubspaceName} \leq Z_{high}\}$ 
5:  $\mathbb{R}_q \leftarrow \phi$  /* Initialize result set to empty set. */
6: for each Subspace  $S \in \mathbb{S}_q$  do
7:    $\langle S_l, S_h \rangle \leftarrow \text{ComputeBounds}(S)$ 
8:   if ( $\langle S_l, S_h \rangle \subseteq \langle q_l, q_h \rangle$ ) then
9:      $\mathbb{R}_q \cup \text{ScanBucketForSpace}(S)$ 
10:  else if ( $\langle S_l, S_h \rangle \cap \langle q_l, q_h \rangle$ ) then
11:    for each point  $p \in S$  do
12:      if ( $p \in \langle q_l, q_h \rangle$ ) then
13:         $\mathbb{R}_q \cup p$ 
14: return  $\mathbb{R}_q$ 

```

4) *Nearest Neighbor Query*: Nearest neighbor queries are also an important primitive operation for many location based applications. Algorithm 4 shows the steps for k nearest neighbors query processing in \mathcal{MD} -HBase. The algorithm is based on the best-first algorithm where the subspaces are scanned in order of the distance from the queried point [9].

The algorithm consists of two steps: *subspace search expansion* and *subspace scan*. During subspace search expansion we incrementally expand the search region and sort subspaces in the region in order of the minimum distance from the queried point. Algorithm 4 increases the search region width to the maximum distance from the queried point to the farthest corner of the scanned subspaces. The next step scans the nearest subspace that has not already been scanned and sorts points in order of the distance from the queried point. If the distance to

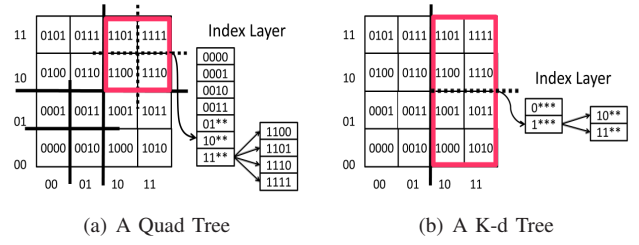


Fig. 6. Space split at the index layer.

the k -th point is less than the distance to the nearest unscanned subspace, the query process terminates.

Algorithm 4 k Nearest Neighbors query

```

1: /*  $q$  be the point for the query. */
2:  $PQ_{subspace} \leftarrow \text{CreatePriorityQueue}()$ 
3:  $PQ_{result} \leftarrow \text{CreatePriorityQueue}(k)$  /* queue capacity is  $k$ . */
4:  $width \leftarrow 0$  /* region size for subspace search */
5:  $\mathbb{S}_{scanned} \leftarrow \phi$  /* scanned subspaces */
6: loop
7:   /* expand search region. */
8:   if  $PQ_{subspace} = \phi$  then
9:      $\mathbb{S}_{next} \leftarrow \text{SubspacesInRegion}(q, width) - \mathbb{S}_{scanned}$ 
10:    for each Subspace  $S \in \mathbb{S}_{next}$  do
11:       $\text{Enqueue}(S, \text{MinDistance}(q, S), PQ_{subspace})$ 
12:    /* pick the nearest subspace. */
13:     $S \leftarrow \text{Dequeue}(PQ_{subspace})$ 
14:    /* search termination condition */
15:    if  $\text{KthDistance}(k, PQ_{result}) \leq \text{MinDistance}(q, S)$  then
16:      return  $PQ_{result}$ 
17:    /* scan and sort points by the distance from  $q$ . */
18:    for each Point  $p \in S$  do
19:       $\text{Enqueue}(p, \text{Distance}(q, p), PQ_{result})$ 
20:    /* maintain search status. */
21:     $\mathbb{S}_{scanned} \cup S$ 
22:     $width \leftarrow \max(width, \text{MaxDistance}(q, S))$ 

```

Algorithm 5 Subspace name generation

```

1: /* Quad tree */
2:  $\text{NewName} \leftarrow \phi$ 
3: for each dimension in  $[1, \dots, n]$  do
4:    $\text{NewName} \cup \{\text{OldName} \oplus 0, \text{OldName} \oplus 1\}$ 
5: /* K-d tree */
6:  $\text{NewName} \leftarrow \{\text{OldName} \oplus 0, \text{OldName} \oplus 1\}$ 

```

5) *Space Split*: Both K-d and Quad trees limit the number of points contained in each subspace; a subspace is split when the number of points in a subspace exceeds this limit. We determine the maximum number of points by the bucket size of the underlying storage layer. Since the index layer is decoupled from the data storage layer, a subspace split in the data storage layer is handled separately. A split in the index layer relies on the first property of the prefix naming scheme which guarantees that the subspace name is a prefix of the names of any enclosed subspace. A subspace split in the index layer therefore corresponds to replacing the row corresponding to the old subspace's name with the names of the new subspaces. The number of new subspaces created depends on the index structure used: a K-d tree splits a subspace only in one dimension, resulting in two new subspaces, while a Quad tree

splits a subspace in all dimensions, resulting in 2^n subspaces. For every dimension split, the name of the new subspaces is created by appending the old subspace name with a 0 and 1 at the position corresponding to the dimension being split. Algorithm 5 provides the pseudocode for subspace name generation following a split. Figure 6 provides an illustration of space splitting in the index layer for both Quad and K-d trees in a 2D space.

Even though conceptually simple, there are a number of implementation intricacies when splitting a space. These intricacies arise from the fact that the index layer is maintained as a table in the Key-value store and most current Key-value stores support transactions only at the single key level. Since a split touches at least three rows (one row for the old subspace name and at least two rows for the new subspaces), such a change is not transactional compared to other concurrent operations. For instance, an analysis query concurrent to a split might observe an inconsistent state. Furthermore, additional logic is needed to ensure the atomicity of the split. Techniques such as maintaining the index as a *key group* as suggested in [11] can be used to guarantee multi-key transactional consistency. We leave such extensions as future work. In addition, since deletions are not very common in LBSs, so we leave out the delete operation. However, deletion and resulting merger of subspaces can be handled as straightforward extension of the Space Split algorithm.

C. Implementing the Index Layer

The index layer is a sorted sequence of subspace names. Since the subspace names encode the boundaries, the index layer essentially imposes an order between the subspaces. In our prior discussion, we represented the

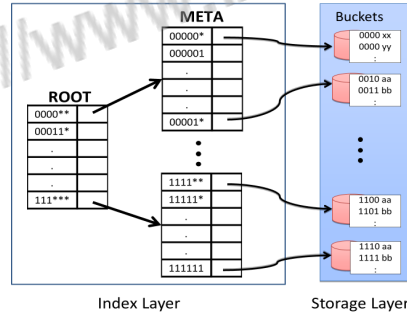


Fig. 7. Index layer Implementation on Bigtable.

index layer as a monolithic layer of sorted subspace names. A single partition index is enough for many application needs. Assume each bucket can hold about 10^6 data points. Each row in the index layer stores very little information: the subspace name, the mapping to the corresponding bucket, and some additional metadata for query processing. Therefore, 100 bytes per index table row is a good estimate. The underlying Key-value store partitions its data across multiple partitions. Considering the typical size of a partition in Key-value stores is about 100 MB [5], the index partition can maintain a mapping of about 10^6 subspaces. This translates to about 10^{12} data points using a single index partition. This estimate might vary depending on the implementation or the configurations used. But it provides a reasonable estimate of the size.

Our implementation also partitions the index layer for better

scalability and load balancing. We leverage the B+ tree style metadata access structure in Bigtable to partition the index layer without incurring any performance penalty. Bigtable, and its open source variant HBase used in our implementation, uses a two level structure to map keys to their corresponding tablets. The top level structure (called the ROOT) is never partitioned and points to another level (called the META) which points to the actual data. Figure 7 provides an illustration of this index implementation. This seamless integration of the index layer into the ROOT - META structure of Bigtable does not introduce any additional overhead as a result of index accesses. Furthermore, optimizations such as caching the index layer and connection sharing between clients on the same host reduces the number of accesses to the index. Adding this additional level in the index structure, therefore, strikes a good balance between scale and the number of indirection levels to access the data items. Using an analysis similar to that used above, a conservative estimate of the number of points indexed by this two level index structure is 10^{21} .

IV. DATA STORAGE LAYER

The data storage layer of \mathcal{MD} -HBase is a range partitioned Key-value store. We use HBase, the open source implementation of Bigtable, as the storage layer for \mathcal{MD} -HBase.

A. HBase Overview

A table in HBase consists of a collection of splits, called **regions**, where each region stores a range partition of the key space. Its architecture comprises a three layered B+ tree structure; the regions storing data constitute the lowest level. The two upper levels are special regions referred to as the ROOT and META regions, as described in Section III-C. An HBase installation consists of a collection of servers, called **region servers**, responsible for serving a subset of regions. Regions are dynamically assigned to the region servers; the META table maintains a mapping of the region to region servers. If a region's size exceeds a configurable limit, HBase splits it into two sub-regions. This allows the system to grow dynamically as data is inserted, while dealing with data skew by creating finer grained partitions for hot regions.

B. Implementation of the Storage Layer

A number of design choices exist to implement the data storage layer, and it is interesting to analyze and evaluate the trade-offs associated with each of these designs. Our implementation uses different approaches as described below.

1) *Table share model*: In this model, all buckets share a single HBase table where the data points are sorted by their corresponding Z-value which is used as the key. Our space splitting method guarantees that data points in a subspace are contiguous in the table since they share a common prefix. Therefore, buckets are managed by keeping their start and end keys. This model allows efficient space splitting that amounts to only updating the corresponding rows in the index table. On the other hand, this model restricts a subspace to be mapped to only a single bucket.

2) *Table per bucket model*: This model is another extreme where we allocate a table per bucket. The data storage layer therefore consists of multiple HBase tables. This model provides flexibility in mapping subspaces to buckets and allows greater parallelism by allowing operations on different tables to be dispatched in parallel. However, a subspace split in this technique is expensive since this involves moving data points from the subspaces to the newly created buckets.

3) *Hybrid model*: This hybrid model strikes a balance between the table share model and the table per bucket model. First, we partition the space and allocate a table to each resulting subspace. After this initial static partitioning, when a subspace is split as a result of an overflow, the newly created subspaces share the same table as that of the parent subspace.

4) *Region per bucket model*: A HBase table comprises of many regions. Therefore, another approach is to use a single table for the entire space and use each region as a bucket. A region split in HBase is quite efficient and is executed asynchronously. This design therefore has a low space split cost while being able to efficiently parallelize operations across the different buckets. However, contrary to the other three models discussed, this model is intrusive and requires changes to HBase; a hook is added to the split code to determine the appropriate split point based on the index structure being used.

C. Optimizations

Several optimizations are possible that further improve performance of the storage layer.

1) *Space Splitting pattern learning*: Space splitting introduces overhead that affects system performance. Therefore, advanced prediction of a split can be used to perform an asynchronous split that will reduce the impact on the inserts and queries executing during a split. One approach is to learn the split patterns to reduce occurrences of space split. Location data is inherently skewed, however, as noted earlier, such skew is often predictable by maintaining and analyzing historical information. Since *MD*-HBase stores historical data, the index structure inherently maintains statistics of the data distribution. To estimate the number of times a new bucket will be split in the future, we lookup how many buckets were allocated to the same spatial region in the past. For example, when we allocate a new bucket for region $([t_0, t_1], [x_0, x_1], [y_0, y_1])$, we lookup buckets for region $([t_0 - t_s, t_1 - t_s], [x_0, x_1], [y_0, y_1])$ in the index table. The intuition is that the bucket splitting pattern in the past is a good approximate predictor for the future.

2) *Random Projection*: Data skew makes certain subspaces hot, both for insertion and querying. An optimization is to map a subspace to multiple buckets with points in the subspace distributed amongst the multiple buckets. When the points are distributed randomly, this technique is called *random projection*. When inserting a data point in a subspace, we randomly select any one of the buckets corresponding to the subspace. A range query that includes the subspace must scan all the buckets mapped to the subspace. Thus, the random projection technique naturally extends our original algorithm;

it, however, presents a trade-off between load balancing and a potential increase in query cost.

V. EXPERIMENTAL EVALUATION

We now present a detailed evaluation of our prototype implementation of *MD*-HBase. We implemented our prototype using HBase 0.20.6 and Hadoop 0.20.2 as the underlying system. We evaluate the trade-offs associated with the different implementations for the storage layer and compare our technique with a MapReduce style analysis and query processing using only linearization over HBase. Our experiments were performed on an Amazon EC2 cluster whose size was varied from 4 to 16 nodes. Each node consists of 4 virtual cores, 15.7GB memory, 1,690 GB HDD, and 64bit Linux (v2.6.32). We evaluate the different implementations of the storage layer as described in Section IV: table per bucket design simulating the K-d and Quad trees (**TPB/Kd** and **TPB/Quad**), table sharing design simulating the K-d and Quad trees (**TS/Kd** and **TS/Quad**), and a region per bucket design for K-d trees (**RPB/Kd**)³. The baseline is an implementation using z-ordering for linearization (**ZOrder**) without any specialized index. We also implemented the queries in Map Reduce to evaluate the performance of a MapReduce system (**MR**) performing a full parallel scan of the data; our evaluation used the Hadoop runtime. Our evaluation uses synthetically generated data sets primarily due to the need for huge data sets (hundreds of gigabytes) and the need to control different aspects, such as skew and selectivity, to better understand the behavior of the system. Evaluation using real data is left for future work.

A. Insert throughput

Supporting high insert throughput for location updates is critical to sustain the large numbers of location devices. We evaluated the insert performance using five different implementations of the storage layer on a cluster with 4, 8, and 16 commodity nodes. Figure 8 plots the insert throughput as a function of the load on the system. We varied the number of load generators from 2 to 64; each generator created a load of 10,000 inserts per second. We used a synthetic spatially skewed data set using a Zipfian distribution with a Zipfian factor of 1.0 representing moderately skewed data. Using a Zipfian distribution allows us to control the skew while allowing quick generation of large data sets. Both the RPB/Kd system and the ZOrder systems showed good scalability; on the 16 node cluster, the RPB/Kd and ZOrder implementation sustained a peak throughput of about 220K location updates per second. In a system where location devices register a location update every minute, this deployment can handle 10–15 million devices. The main reason for the low scalability of the table per bucket and table sharing designs is the cost associated with the splitting a bucket. On the other hand, the region per bucket design splits buckets asynchronously using the HBase region split mechanism which is relatively

³Since in HBase, a region can only be split into two sub-regions, we could not implement RPB for Quad trees as our experiments are for a 3D space.

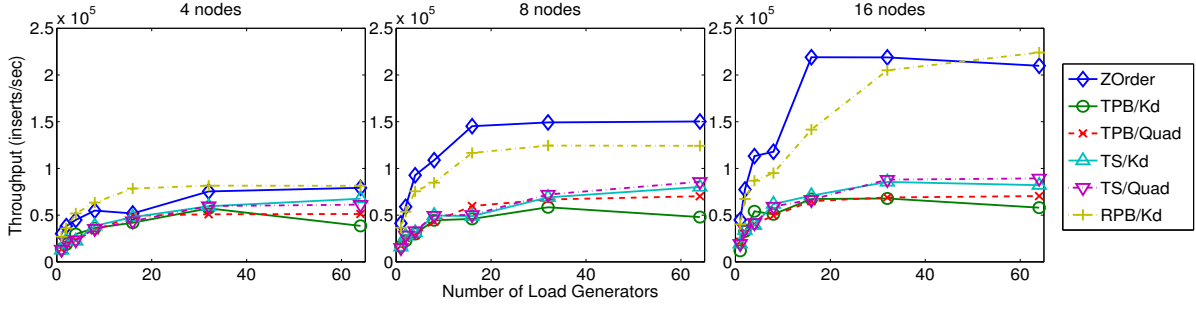


Fig. 8. Insert throughput (location updates per second) as a function of load on the system. Each load generator creates 10,000 inserts per second.

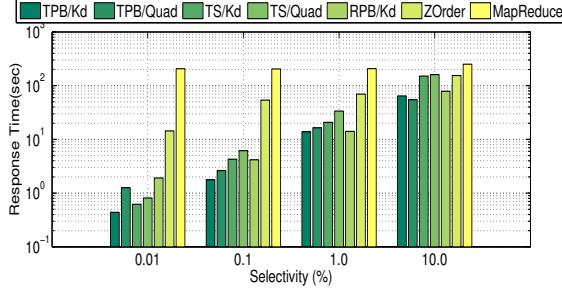


Fig. 9. Response times (in log scale) for range query as function of selectivity.

TABLE I
FALSE POSITIVE SCANS ON THE ZORDER SYSTEM

	Selectivity (%)			
	0.01	0.1	1.0	10.0
No. of buckets scanned	7	28	34	45
False positives	3	22	16	16
Percentage false positive	42.9%	78.5%	47.1%	35.5%

inexpensive. The TPB/TS systems block other operations until the bucket split completes. In our experiments, the TPB design required about 30 to 40 seconds to split a bucket and the TS design required about 10 seconds. Even though these designs result in more parallelism in distributing the insert load on the cluster, the bucket split overhead limits the peak throughput sustained by these designs.

B. Range Query

We now evaluate range query performance using the different implementations of the index structures and the storage layer and compare performance with ZOrder and MR. The MR system filters out points matching the queried range and reports aggregate statistics on the matched data points.

We generated about four hundred million points using a network-based model by Brinkhoff et al. [12]. Our data set simulates 40,000 objects moving 10,000 steps along the streets of the San Francisco bay area. Since the motion paths are based on a real map, this data set is representative of real world data distributions. Evaluation using other models to simulate motion is proposed future work. We executed the range queries on a four-node cluster in Amazon EC2.

Figure 9 plots the range query response times for the different designs as a function of the varying selectivity. As

is evident from Figure 9, all the \mathcal{MD} -HBase design choices outperform the baseline systems. In particular, for highly selective queries where our designs show a one to two orders of magnitude improvement over the baseline implementations using simple Z-ordering or using MapReduce. Moreover, the query response time of our proposed designs is proportional to the selectivity, which asserts the gains from the index layer when compared to brute force parallel scan of the entire data set as in the MR implementation whose response times are the worst amongst the alternatives considered and is independent of the query selectivity.

In the design using just Z-ordering (ZOrder), the system scans between the minimum and the maximum Z-value of the queried region and filters out points that do not intersect the queried ranges. If the scanned range spans several buckets, the system parallelizes the scans per bucket. Even though the ZOrder design shows better performance compared to MR, response time is almost ten times worse when compared to our proposed approaches, especially for queries with high selectivity. The main reason for the inferior performance of the ZOrder is the false positive scans resulting from the distortion in the mapping of the points introduced by linearization. Table I reports the number of false positives (buckets that matched the queried range but did not contain a matching data point) of the ZOrder design. For example, in case of 0.1 percent selectivity query, 22 out of 28 buckets scanned are false positive. The number of false positive scans depends on the alignment of the queried range with the bucket boundaries. In the ideal situation of a perfect match, query performance is expected to be close to that of our proposed designs; however, such ideal queries might not be frequently observed in practice. Our proposed designs can further optimize the scan range within each region. Since we partition the target space into the regular shaped subspaces, when the queried region partially intersects a subspace, \mathcal{MD} -HBase can compute the minimum scan range in the bucket from the boundaries of the queried region and the subspace. In contrast, the ZOrder design partitions the target space into irregular shaped subspaces and hence must scan all points in the selected region.

A deeper analysis of the different alternative designs for \mathcal{MD} -HBase shows that the TPB designs result in better performance. The TS designs often result in high disk contention since some buckets are co-located in the same data

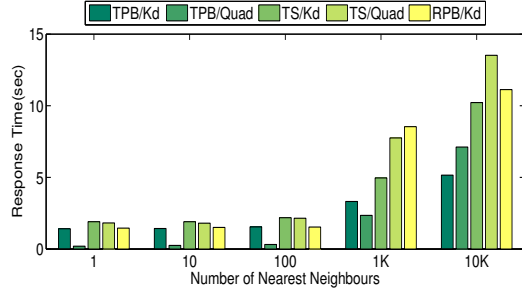


Fig. 10. Response time for kNN query processing as a function of the number of neighbors returned.

partition of the underlying Key-value store. Even though we randomize the bucket scan order in the TS design, significant reduction in disk access contention is not observed. The TPB designs therefore outperform the TS designs. For queries with high selectivity, the RPB design has a longer response time compared to other designs; this however is an artifact of our implementation choice for the RPB design. In the RPB design, we override the pivot calculation for a bucket split and delegate the bucket maintenance task to HBase which minimizes the insert overhead. However, since HBase stores bucket boundary information in the serialized form, the index lookup for query processing must de-serialize the region boundaries. As a result, for queries with very high selectivity, the index lookup cost dominates the total response time.

Comparing the K-d and the Quad trees, we observed better performance of K-d trees for queries with high selectivity. However, as the selectivity decreases, the response times for Quad trees are lower. The K-d tree creates a smaller number of buckets while the Quad tree provides better pruning compared to the K-d tree. In the case of high selectivity queries, sub-query creation cost dominates the total response time; each sub-query typically scans a short range and the time to scan a bucket comprises a smaller portion of the total response time. As a result, the K-d tree has better performance compared to the Quad tree. On the other hand, in case of low selectivity queries, the number of points to be scanned dominates the total response time since each sub-query tends to scan a long range. The query performance therefore depends on the pruning power of the index structure.

We also expect the bucket size to potentially effect overall performance. Selecting a bucket size is a trade-off between the insertion throughput and query performance. Larger buckets reduce the frequency of bucket splits thus improving insertion throughput but limits sub-space pruning power during query processing.

C. kNN Query

We now evaluate the performance of *MD*-HBase for k Nearest Neighbor (kNN) queries using the same dataset as the previous experiment. Figure 10 plots the response time (in seconds) for kNN queries as a function of the number of neighbors to be returned; we varied k as 1, 10, 100, 1,000, and 10,000. In this experiment, expansion of the search

region did not occur for $k \leq 100$, resulting in almost similar performance for all the designs. The TPB/Quad design has the best performance where the response time was around 250ms while that of other designs is around 1500ms to 2000ms. As we increase k to 1K and 10K, the response times also increase; however, the increase is not exponential, thanks to the best-first search algorithm used. For example, when we increased k from 1K to 10K, the response time increased only by three times in the worst case.

The TPB/Quad design outperformed the other designs for $k \leq 100$. Due to the larger fan-out, the bucket size of a Quad tree tends to be smaller than that of a K-d tree. As a result, TPB/Quad has smaller bucket scan times compared to that of the TPB/Kd design. However, both TS designs have almost similar response times for small k . We attribute this to other factors like the bucket seek time. Since the TPB designs use different tables for every bucket, a scan for all points in a bucket does not incur a seek to the first point in the bucket. On the other hand, the TS designs must seek to the first entry corresponding to a bucket when scanning all points of the bucket. As mentioned earlier, the RPB design has a high index lookup overhead which increases the response time, especially when k is large. This however is an implementation artifact; in the future, we plan to explore other implementation choices for RPB that can potentially improve its performance.

Comparing between the K-d tree and the Quad tree, Quad trees have smaller response times for smaller k , while K-d trees have smaller response times for larger k . For small values of k , the search region does not expand. Smaller sized buckets in a Quad tree result in better performance compared to larger buckets in a K-d tree. On the other hand, for larger values of k , the probability of a search region expansion is higher. Since a Quad tree results in smaller buckets, it results in more search region expansions compared to that of K-d tree, resulting in better performance of K-d tree.

Since the ZOrder and the MR systems do not have any index structure, kNN query processing is inefficient. One possible approach is to iteratively expand the queried region until the desired k neighbors are retrieved. This technique is however sensitive to a number of parameters such as the initial queried range and the range expansion ratio. Due to the absence of an index, there is no criterion to appropriately determine the initial query region. We therefore do not include these baseline systems in our comparison; however, the performance of ZOrder and MR is expected to be significantly worse than the proposed approaches.

VI. RELATED WORK

Scalable data processing—both for large update volumes and data analysis—has been an active area of interest in the last few years. When dealing with a large volume of updates, Key-value stores have been extremely successful in scaling to large data volumes; examples include Bigtable [5], HBase [6] etc. These data stores have a similar design philosophy where they have made scalability and high availability as the primary requirement, rich functionality being secondary. Even though

these systems are known to scale to terabytes of data, the lack of efficient multi-attribute based access limits their application for location based applications. On the other hand, when considering scalable data processing, MapReduce [7], [13] has been the most dominant technology, in addition to parallel databases [14]. Such systems have been proven to scale to petabytes of data while being fault-tolerant and highly available. However, such systems are suitable primarily in the context of batch processing. Furthermore, in the absence of appropriate multi-dimensional indices, MapReduce style processing has to scan through the entire dataset to answer queries. *MD*-HBase complements both Key-value stores and MapReduce style processing by providing an index structure for multi-dimensional data. As demonstrated in our prototype, our proposed indexing layer can be used directly with Key-value stores. Along similar lines, the index layer can also be integrated with MapReduce to limit the retrieval of false positives.

Our use of linearization for transforming multi-dimensional data points to a single dimensional space has been used recently in a number of other techniques. For instance, Jensen et al. [15] use Z-ordering and Hilbert curves as space filling curves and construct a B+ tree index using the linearized values. Tao et al. [16] proposed an efficient indexing technique for high dimensional nearest neighbor search using a collection of B-tree indices. The authors first use locality sensitive hashing to reduce the dimensionality of the data points and then apply Z-ordering to linearize the dataset, which is then indexed using a B-tree index. Our approach is similar to these approaches, the difference being that we build a K-d tree and a Quad tree based index using the linearized data points. The combination of the index layer and the data storage layer in *MD*-HBase however resembles a B+ tree, reminiscent of the Bigtable design. Subspace pruning in the index layer is key to speeding up the range query performance which becomes harder for data points with high dimensionality. In such cases, dimensionality reduction techniques, as used by Tao et al. [16], can be used to improve the pruning power.

Another class of approaches make the traditional multidimensional indices more scalable. Wang et al. [17] and Zhang et al. [18] proposed similar techniques where the systems have two index layers: a global index and a local index. The space is partitioned into several subspaces and each subspace is assigned a local storage. The global index organizes subspaces and the local index organizes data points in the subspace. Wang et al. [17] construct a content addressable network (CAN) over a cluster of R-tree indexed databases while Zhang et al. [18] use an R-tree as the global index and a K-d tree as the local index. Along these lines, *MD*-HBase only has a global index; it can however be extended to add local indices within the data storage layer.

VII. CONCLUSION

Scalable location data management is critical to enable the next generation of location based services. We proposed *MD*-

HBase, a scalable multi-dimensional data store supporting efficient multi-dimensional range and nearest neighbor queries. *MD*-HBase layers a multi-dimensional index structure over a range partitioned Key-value store. Using a design based on linearization, our implementation layers standard index structures like K-d trees and Quad trees. We implemented our design on HBase, a standard open-source key-value store, with minimal changes to the underlying system. The scalability and efficiency of the proposed design is demonstrated through a thorough experimental evaluation. Our evaluation using a cluster of nodes demonstrates the scalability of *MD*-HBase by sustaining insert throughput of over hundreds of thousands of location updates per second while serving multi-dimensional range queries and nearest neighbor queries in real time with response times less than a second. In the future, we plan to extend our design by adding more complex analysis operators such as skyline or cube, and exploring other alternative designs for the index and data storage layers.

ACKNOWLEDGMENTS

The authors would like to thank the anonymous reviewers for their helpful comments. This work is partly funded by NSF grants III 1018637 and CNS 1053594.

REFERENCES

- [1] http://en.wikipedia.org/wiki/List_of_mobile_network_operators, 2010.
- [2] J. L. Bentley, "Multidimensional binary search trees used for associative searching," *Commun. ACM*, vol. 18, no. 9, pp. 509–517, 1975.
- [3] R. A. Finkel and J. L. Bentley, "Quad trees: A data structure for retrieval on composite keys," *Acta Inf.*, vol. 4, pp. 1–9, 1974.
- [4] A. Guttman, "R-trees: a dynamic index structure for spatial searching," in *SIGMOD*, 1984, pp. 47–57.
- [5] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, and R. E. Gruber, "Bigtable: A Distributed Storage System for Structured Data," in *OSDI*, 2006, pp. 205–218.
- [6] "HBase: Bigtable-like structured storage for Hadoop HDFS," 2010, <http://hadoop.apache.org/hbase/>.
- [7] J. Dean and S. Ghemawat, "MapReduce: simplified data processing on large clusters," in *OSDI*, 2004, pp. 137–150.
- [8] G. M. Morton, "A computer oriented geodetic data base and a new technique in file sequencing," IBM Ottawa, Canada, Tech. Rep., 1966.
- [9] H. Samet, *Foundations of Multidimensional and Metric Data Structures*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2005.
- [10] S. Ramabhadran, S. Ratnasamy, J. M. Hellerstein, and S. Shenker, "Prefix hash tree: An indexing data structure over distributed hash tables," Intel Research, Berkeley, Tech. Rep., 2004.
- [11] S. Das, D. Agrawal, and A. El Abbadi, "G-Store: A Scalable Data Store for Transactional Multi key Access in the Cloud," in *SOCC*, 2010, pp. 163–174.
- [12] T. Brinkhoff and O. Str, "A framework for generating network-based moving objects," *Geoinformatica*, vol. 6, p. 2002, 2002.
- [13] "The Apache Hadoop Project," <http://hadoop.apache.org/core/>, 2010.
- [14] D. DeWitt and J. Gray, "Parallel database systems: the future of high performance database systems," *CACM*, vol. 35, no. 6, pp. 85–98, 1992.
- [15] C. S. Jensen, D. Lin, and B. C. Ooi, "Query and update efficient b+-tree based indexing of moving objects," in *VLDB*. VLDB Endowment, 2004, pp. 768–779.
- [16] Y. Tao, K. Yi, C. Sheng, and P. Kalnis, "Quality and efficiency in high dimensional nearest neighbor search," in *SIGMOD*, 2009, pp. 563–576.
- [17] J. Wang, S. Wu, H. Gao, J. Li, and B. C. Ooi, "Indexing multi-dimensional data in a cloud system," in *SIGMOD*, 2010, pp. 591–602.
- [18] X. Zhang, J. Ai, Z. Wang, J. Lu, and X. Meng, "An efficient multi-dimensional index for cloud data management," in *CloudDB*, 2009, pp. 17–24.



知网查重限时 7折 最高可优惠 120元

本科定稿，硕博定稿，查重结果与学校一致

立即检测

免费论文查重: <http://www.paperyy.com>

3亿免费文献下载: <http://www.ixueshu.com>

超值论文自动降重: http://www.paperyy.com/reduce_repetition

PPT免费模版下载: <http://ppt.ixueshu.com>
